

Does hedging with implied volatility factors improve the hedging efficiency of barrier options? *

Szymon Borak[†]
Matthias R. Fengler
Wolfgang K. Härdle

CASE – Center for Applied Statistics and Economics
Humboldt-Universität zu Berlin,
Spandauer Straße 1, 10178 Berlin, Germany

February 11, 2009

Forthcoming: The Journal of Risk Model Validation

*We gratefully acknowledge financial support by the Deutsche Forschungsgemeinschaft and the Sonderforschungsbereich 649 “Ökonomisches Risiko”.

[†]Corresponding author: borak@wiwi.hu-berlin.de, TEL ++49 30 2093 5630 FAX ++49 30 2093 5649.

Does hedging with implied volatility factors improve the hedging efficiency of barrier options?

Abstract

The price of a barrier option depends on the shape of the entire implied volatility surface which is a high-dimensional dynamic object. Barrier options are hence exposed to nontrivial volatility risk. We extract the key risk factors of implied volatility surface fluctuations by means of a semiparametric factor model. Based on the factors we define a practical hedging procedure within a local volatility framework. The hedging performance is evaluated using DAX index options.

JEL classification codes: G11

Keywords: implied volatility surface, smile, local volatility, exotic options, semiparametric factor model, hedging

1 Introduction

In equity derivative markets barrier options are appealing instruments for investors looking for a partial protection of their equity allocation. From the perspective of an institution issuing barrier options this demand raises the need of efficient hedging strategies. This is a challenging task for at least two reasons. First, reverse barrier options, such as down-and-out puts and up-and-out calls, have discontinuous payoff profiles and knock out deep in-the-money thereby losing the maximum possible intrinsic value. Second, barrier options, as many other exotic options, are exposed to nontrivial volatility risk, since the knock-out probability strongly depends on the skew of the implied volatility smile. The latter effect also prevents simple Black-Scholes type formulae, such as those by Rubinstein and Reiner (1991), from being usable in practice.

Nowadays there is a plethora of models available that take the shape of the implied volatility surface (IVS) into account for option valuation. Potential candidates are: the local volatility (LV) model proposed by Dupire (1994), Derman and Kani (1994), and Rubinstein (1994), which introduces a nonparametric local volatility function that deterministically depends on the asset price and time; stochastic volatility models like Hull and White (1987), Stein and Stein (1991), Heston (1993), Carr et al. (2003); jump-diffusion models, such as Merton (1976), Bates (1996), and Kou (2002). When calibrated to the IVS, all these models are able to replicate the plain vanilla market to a similar extent, whereas their prices for barrier options may differ due to the different properties of the underlying asset price dynamics, see Hull and Suo (2002) and Hirt et al. (2003) on model risk for barrier options. The more challenging part is hedging. For it is straight forward to compute derivatives for the parameters of these models, but it is intricate to give the parameter greeks a meaning by mapping them on tradable instruments provided by the plain vanilla market. More seriously, since the prices of the hedging instruments, either over-the-counter or as listed options, are given in terms of implied volatility, they necessarily follow the dynamics of the IVS. Indeed it is in question whether the IVS dynamics inherent in the model that is calibrated to a static surface and used for pricing truly match the stylized facts of IVS dynamics, see Hagan et al. (2002) and Bergomi (2005) for such a discussion in context of the LV model and the Heston model, respectively. In contrast, the dynamics of the IVS are empirically well understood, see Skiadopoulos et al. (1999), Alexander (2001), Cont and da Fonseca (2002), Fong et al.

(2003), Hafner (2004), Fengler et al. (2007) among others. The typical approach extracts the main driving factors like level, slope, or term structure movements and models these factors. It therefore appears natural to exploit this knowledge for hedging and portfolio risk management.

The aim of this paper is to study dynamic hedges of reverse barrier options built on factor functions of empirically observed IVS dynamics. We project the complex, high dimensional dynamics of the IVS on a low and finite dimensional space spanned by the semiparametric factor model (SFM)

$$\hat{\sigma}_t(\kappa, \tau) = \exp \left\{ \sum_{l=0}^L Z_{t,l} m_l(\kappa, \tau) \right\}, \quad (1)$$

where $\hat{\sigma}_t(\kappa, \tau)$ denotes the implied volatility of a certain moneyness κ and maturity τ observed in time t . The functions m are nonparametric components and invariant in time, while the time evolution is modelled by the latent factor series $Z_{t,l}$. In order to estimate (1) we apply an estimation technique suggested in Fengler et al. (2007). The SFM estimates the prevalent movements of the IVS in an $(L + 1)$ -dimensional function space.

Given the estimated factor functions \hat{m} , we construct hedges for barrier options priced in a LV model. We use a LV model, since by the nonparametric nature of the local volatility function it can match any arbitrage-free set of option prices to an arbitrarily precise degree. It will hence replicate the deformations of the IVS defined by the estimated factor functions and allow for a precise computation of factor greeks not prone to calibration error. Moreover, the LV model is numerically very efficient and allows for fast and accurate price valuations using the finite difference method. The factor hedges we obtain are more general than the usual vega hedges which are defined by a parallel shift of the IVS since they will take into account nontrivial surface movements, such as nonparallel up-and-down shifts, slope and term structure risks. Depending on the payoff profile of an exotic option, these risks can be substantial. Our approach is hence similar in spirit to Diebold et al. (2006) who define factor based duration measures and study the efficacy of these measures for the insurance of bond portfolios.

We note that strictly speaking it may not be necessary to vega hedge in an LV framework,

since it defines a complete market. This however is a theoretical perspective which does not correspond to market practice. When minimizing portfolio risk, traders are likely to set up vega hedges as soon as a liquid over-the-counter or listed option markets allow them to do so. In this sense our approach is e.g. similar to the practice of hedging a long dated plain vanilla option which are priced by means of a smile-adjusted Black-Scholes model by adding a short dated option to the portfolio.

The dynamic hedging performance of plain vanilla options in a LV model is studied in Dumas et al. (1998), Coleman et al. (2001), McIntyre (2001) and Vähämaa (2004), while the case of reverse barrier options is treated in Engelmann et al. (2006). Engelmann et al. (2006) implement hedging strategies that are delta ($\partial/\partial S$), vega ($\partial/\partial\sigma$) and vanna ($\partial^2/\partial\sigma\partial S$) neutral where vega and vanna are obtained by parallel shifts of the IVS and computing the difference quotient. We complement this analysis by defining sensitivities with respect to the most prevalent IVS movements motivated by model (1), namely ($\partial/\partial Z_1$), ($\partial/\partial Z_2$) and by constructing portfolios neutral to these greeks. For this purpose we establish a portfolio containing a reverse barrier option and hedge it on a daily basis with plain vanillas and the underlying asset using DAX data from January 3rd, 2000 to June 30th, 2004. We then study the distribution of the hedging errors across the different hedging strategies.

For completeness we remark that static hedging of barrier options is a competing way of portfolio insurance, see Derman et al. (1995), Carr and Chou (1997), Carr et al. (1998), Andersen et al. (2002), Tompkins (2002), Nalholm and Poulsen (2006a), Nalholm and Poulsen (2006b). For a static hedge one sets up a portfolio of plain vanillas which replicates the payoff of the barrier option as close as possible. The hedge is unwound in case of a knock-out or at expiry and no other adjustment of the hedge is necessary. In fact, Engelmann et al. (2007) and Maruhn et al. (2008) show that there are static hedges outperforming dynamic hedges. However, the practical use of static hedges is limited, since they may not always be implementable due to insufficient market depth of listed plain vanilla options.

The paper is structured as follows. In Section 2 we present the framework on which the empirical procedure is based. Section 3 concentrates on the description of the hedging method. In Section 4 we present the data, describe the empirical hedging design and discuss the empirical results. Section 5 concludes.

2 Models

2.1 Local Volatility Model

In the LV model the risk neutral price of the underlying asset is governed by the stochastic differential equation:

$$dS_t = r_t S_t dt + \sigma(S_t, t) S_t dW_t, \quad (2)$$

where W_t is a Wiener process and r_t denotes the instantaneous interest rate. Dividends are assumed to be zero, since the DAX, on which our empirical study is based, is a performance index. $\sigma(S_t, t)$ is the local volatility function which depends on the underlying price and time. This function has a unique representation if an arbitrage-free set of call options is given for all strikes and maturities, Dupire (1994). It can be shown that

$$\sigma^2(S_t, t) = \frac{2 \frac{\partial \hat{\sigma}(K, T)}{\partial T} + \frac{\hat{\sigma}(K, T)}{T} + 2K \int_0^T r_s ds \frac{\partial \hat{\sigma}(K, T)}{\partial K}}{K^2 \left\{ \frac{\partial^2 \hat{\sigma}(K, T)}{\partial K^2} - d_1 \sqrt{T} \left(\frac{\partial \hat{\sigma}(K, T)}{\partial K} \right)^2 + \frac{1}{\hat{\sigma}(K, T)} \left(\frac{1}{K\sqrt{T}} + d_1 \frac{\partial \hat{\sigma}(K, T)}{\partial K} \right)^2 \right\}} \Bigg|_{K=S_t, T=t} \quad (3)$$

where $d_1 = \frac{\log(S_0/K) + \int_0^T r_s ds + 0.5 \hat{\sigma}^2(K, T) T}{\hat{\sigma}(K, T) \sqrt{T}}$ and where $\hat{\sigma}(K, T)$ is the implied volatility at strike K and expiry T . Formula (3) gives a correspondence between local and implied volatility surfaces.

The LV model received much attention in the finance community since it achieves an almost exact fit of the observed vanilla market and is numerically and computationally very tractable. The price of the barrier option denoted by V with barrier B and expiry date T is obtained by numerically solving the partial differential equation

$$r_t V(S, t) = \frac{\partial V(S, t)}{\partial t} + \frac{1}{2} \sigma^2(S, t) S^2 \frac{\partial^2 V(S, t)}{\partial S^2} + r_t S \frac{\partial V(S, t)}{\partial S} \quad (4)$$

with additional boundary conditions, i.e. $V(B, t) = 0$ for $t < T$ and $V(S, T)$ equal to the payoff at expiry. For calibration of the model a number of methods are available, see Bouchouev and Isakov (1999) for comprehensive review. For example one may directly apply

the formula (3). Here we adopt the approach of Andersen and Brotherton-Ratcliffe (1997) which determines r and σ so that forwards, zero coupon bonds and plain vanilla options are priced correctly on each grid point. The finite difference method then gives barrier option prices and sensitivities very efficiently.

Yet the LV is also subject to criticism, see Fengler (2005, Chapter 3.11) for the details of this discussion. The severest objection was brought forward by Hagan et al. (2002) by showing that the LV model implies unrealistic smile dynamics and consequently wrong spot greeks. In practice this problem can be addressed by enforcing the desired smile dynamics when computing the greeks. Instead of calculating model-consistent LV greeks, one fixes the IVS in strikes (sticky-strike) or in moneyness (sticky-moneyness) and recalibrates the LV surface under the spot movements. Engelmann et al. (2006) find that the empirical performance of the dynamic hedges is negligible under different stickiness assumptions, if a vega hedge is implemented. Overall they find that the sticky-strike approach, which we will adopt here, performs best. We therefore believe that the LV model serves well for the purpose of this study.

2.2 The Semiparametric Factor Model

To model the IVS dynamics we employ the SFM which yields estimates of the IVS for each day of the sample and explains its dynamic behavior by extracting a small number of key driving factors of the surface movements. For this aim one could use any other factor model like the functional principal components model of Cont and da Fonseca (2002) or the parametric model of Hafner (2004). An alternative definition of the skew shifts can be also found in Taleb (1997). Our choice for the SFM is motivated by the flexible nonparametric structure, which allows to extract the most important factors along with a dimension reduction, and its adaptedness to the expiry behavior of implied volatility data, see Fengler et al. (2007) for details.

To describe the SFM denote by $Y_{t,j}$ the log-implied volatility observed on day $t = 1, \dots, T$. The index $j = 1, \dots, J_t$ counts the implied volatilities observed on day t . Let $X_{t,j}$ be a two-dimensional variable containing (forward) moneyness $\kappa_{t,j}$ and time to maturity $\tau_{t,j}$. We define the moneyness $\kappa_{t,j} \stackrel{\text{def}}{=} K_{t,j}/F_{\tau_{t,j}}$, where $K_{t,j}$ is a strike and $F_{\tau_{t,j}}$ the forward price of

the underlying asset at time t . The SFM regresses $Y_{t,j}$ on $X_{t,j}$ by:

$$Y_{t,j} = \sum_{l=0}^L Z_{t,l} m_l(X_{t,j}) + \varepsilon_{t,j}, \quad (5)$$

where m_l ($l = 1, \dots, L$) are nonparametric components and the $Z_{t,l}$ form a latent factor series depending on time t . The estimation error is denoted by $\varepsilon_{t,j}$. The basis functions m_0, \dots, m_L are constant in time, while the dynamic propagation of the IVS is modelled by the time varying weights $Z_{t,l}$.

The estimation procedure is based on minimizing the following least squares criterion ($\widehat{Z}_{t,0} \equiv 1$ for identification):

$$\sum_{t=1}^T \sum_{j=1}^{J_t} \int \left\{ Y_{t,j} - \sum_{l=0}^L \widehat{Z}_{t,l} \widehat{m}_l(u) \right\}^2 K_h(u - X_{t,j}) du, \quad (6)$$

where K_h denotes a two-dimensional kernel function. A possible choice for a two-dimensional kernel is a product of one-dimensional kernels $K_h(u) = k_{h_1}(u_1) \times k_{h_2}(u_2)$, where $h = (h_1, h_2)^\top$ are bandwidths and $k_h(v) = h^{-1}k(h^{-1}v)$ is a one dimensional kernel function. The minimization procedure searches across all functions $\widehat{m}_l : \mathbb{R}^2 \rightarrow \mathbb{R}$ ($l = 0, \dots, L$) and time series $\widehat{Z}_{t,l} \in \mathbb{R}$ ($t = 1, \dots, T; l = 1, \dots, L$). Details concerning the estimation algorithm can be found in Fengler et al. (2007) and Park et al. (2009). In the final step of the procedure one orthogonalizes the functions $\widehat{m}_1, \dots, \widehat{m}_L$ and orders them with respect to the variance explained. As a consequence the largest portion of variance is explained by the quantity $\widehat{Z}_{t,1}\widehat{m}_1$ and the second largest by $\widehat{Z}_{t,1}\widehat{m}_1 + \widehat{Z}_{t,2}\widehat{m}_2$ and so forth.

In order to illustrate the decomposition of the IVS dynamics achieved by the SFM we present in Figure 1 the results on DAX option data from January 3rd, 2000 till June 30th, 2004. The figure presents the estimated $\widehat{Z}_{t,l}$ time series in the upper panel and the estimates of the basis functions in the lower panel. The function \widehat{m}_0 is not presented to save space. It has no effect on the dynamics of the IVS but has to be included to set the correct level of the surface. The function \widehat{m}_1 is relatively flat and corresponds to the most important shocks. Changes in $\widehat{Z}_{t,1}$ result in up-and-down type of movements of the whole surface, but the deviations from a

flat basis function give different weight for each maturity-moneyness location. This effect is illustrated in Figure 2, where we plot several surfaces and one particular smile with different values of $\widehat{Z}_{t,1}$. The second factor function can be interpreted as a tilting of the smile. This can be inferred from the shape of \widehat{m}_2 and its influence on the IVS in the plots. The variation in $\widehat{Z}_{t,2}$ results in changing the slope of the smile by making it steeper or flatter while keeping roughly the same implied volatility levels.

We finally remark that the SFM has spurred further research on IVS dynamics and beyond. Brüggemann et al. (2008) study the statistical properties of the estimated factor series using a vector autoregressive framework and analyze the associated movements of macroeconomic variables. Giacomini and Härdle (2008) apply the modelling idea for an explanation of the dynamics of risk neutral densities. The CO₂ allowance term structure is studied in Trück et al. (2006) and electricity forward curves in Borak and Weron (2009).

3 Hedging Framework

Dynamic hedging of the asset V , in our case the reverse barrier option, is based on frequent adjustments of the hedge portfolio. This hedging strategy requires to construct a portfolio which is to first (or higher) order neutral to the relevant risk factors. Apart from standard delta hedging, a successful strategy requires hedging the vega, and possibly higher order greeks as pointed out by Ederington and Guan (2007).

For the LV framework Engelmann et al. (2006) study delta, delta-vega and delta-vega-vanna hedges. One knock-out option is hedged with the underlying asset and a set of plain vanilla options. Let the value of the barrier option be denoted by V and let HP_1 and HP_2 be portfolios of plain vanilla options. The corresponding hedge ratios are then given by solving

$$\begin{pmatrix} 1 & \frac{\partial HP_1}{\partial S} & \frac{\partial HP_2}{\partial S} \\ 0 & \frac{\partial HP_1}{\partial \sigma} & \frac{\partial HP_2}{\partial \sigma} \\ 0 & \frac{\partial^2 HP_1}{\partial \sigma \partial S} & \frac{\partial^2 HP_2}{\partial \sigma \partial S} \end{pmatrix} \cdot \begin{pmatrix} a_0 \\ a_1 \\ a_2 \end{pmatrix} = \begin{pmatrix} \frac{\partial V}{\partial S} \\ \frac{\partial V}{\partial \sigma} \\ \frac{\partial^2 V}{\partial \sigma \partial S} \end{pmatrix}. \quad (7)$$

Equation (7) reflects the full delta-vega-vanna hedge. Putting $a_2 = 0$ reduces (7) to the

delta-vega hedge and $a_1 = a_2 = 0$ to the pure delta hedge. Since good hedges have a large exposure to the risk factors to be hedged, one could use an at-the-money plain vanilla option for the HP_1 and for HP_2 a risk reversal. A risk reversal is a combination of a long out-of-the-money call and a short out-of-the-money put (or vice versa).

In order to compute the sensitivities one reprices the option under different scenarios and computes the greeks by a finite difference quotient. Following Engelmann et al. (2006), we make a sticky strike assumption for our greeks, i.e. the IVS remains constant in strikes. Vega and vanna are computed shifting the IVS in a parallel fashion. To be more specific, we compute

$$\frac{\partial V}{\partial S} \stackrel{\text{def}}{\approx} \frac{V(S + \Delta S, \hat{\sigma}) - V(S - \Delta S, \hat{\sigma})}{2\Delta S}, \quad (8)$$

$$\frac{\partial V}{\partial \hat{\sigma}} \stackrel{\text{def}}{\approx} \frac{V(S, \hat{\sigma} + \Delta \hat{\sigma}) - V(S, \hat{\sigma} - \Delta \hat{\sigma})}{2\Delta \hat{\sigma}}, \quad (9)$$

$$\frac{\partial^2 V}{\partial S \partial \hat{\sigma}} \stackrel{\text{def}}{\approx} \frac{\{V(S + \Delta S, \hat{\sigma} + \Delta \hat{\sigma}) - V(S + \Delta S, \hat{\sigma}) - V(S - \Delta S, \hat{\sigma} + \Delta \hat{\sigma}) + V(S - \Delta S, \hat{\sigma})\}}{2\Delta S \Delta \hat{\sigma}}. \quad (10)$$

With small abuse of notation $V(S, \hat{\sigma})$ denotes here the price obtained with spot S and IVS $\hat{\sigma}$, where we omit its arguments for simplicity. $\hat{\sigma} + \Delta \hat{\sigma}$ means the parallel shift of the whole surface.

It is empirically widely confirmed that parallel shifts are the most prevalent movements of the IVS. It would be misleading, however, to conclude from this observation that other types of surface variations do only negligibly influence the prices of exotic derivatives, such as barrier options. Contrariwise a higher slope leads to a smaller price of an in-the-money down-and-out put. Consider an artificial example of two one year down-and-out put with strike 110, barrier 80 at the current spot level of 100. The first option is priced with the IVS observed on January 3rd, 2000 and the second one on January 2nd, 2001. Figure 3 shows the surfaces of these days. The LV prices of these options are 1.91% and 2.37% respectively (in percentage of the spot price), which is quite a difference. From the upper panel of Figure 1 one observes that the level related factor assumes similar values on these days, while the slope factor differs significantly. This price discrepancy stems mainly from the slope effect,

which is an exposure not directly hedged in traditional approaches. Our procedure will hedge such volatility shocks.

In our hedging framework we define new sensitivities with respect to the variation of the (log)-IVS, which we call ζ -greeks. Based on the results discussed in Section 2.2, the ζ_1 -greek ($\partial/\partial Z_{t,1}$) reflects an adjusted up-and-down shift, while the ζ_2 -greek ($\partial/\partial Z_{t,2}$) corresponds to the slope effect. Similarly to (7) we obtain the hedge ratios by

$$\begin{pmatrix} 1 & \frac{\partial HP_1}{\partial S} & \frac{\partial HP_2}{\partial S} \\ 0 & \frac{\partial HP_1}{\partial Z_{t,1}} & \frac{\partial HP_2}{\partial Z_{t,1}} \\ 0 & \frac{\partial HP_1}{\partial Z_{t,2}} & \frac{\partial HP_2}{\partial Z_{t,2}} \end{pmatrix} \cdot \begin{pmatrix} a_0 \\ a_1 \\ a_2 \end{pmatrix} = \begin{pmatrix} \frac{\partial V}{\partial S} \\ \frac{\partial V}{\partial Z_{t,1}} \\ \frac{\partial V}{\partial Z_{t,2}} \end{pmatrix}. \quad (11)$$

We call the full setting a $\zeta_1\zeta_2$ -hedge, the reduced one with $a_2 = 0$ a ζ_1 -hedge. As in the traditional hedge we use an at-the-money plain vanilla for HP_1 , again due to the high vega. For HP_2 , we employ risk reversals because they primarily respond to changes in the wings of the IVS. Moreover, by selecting appropriate strikes it can even be set up in a vega-neutral, i.e. ζ_1 -neutral, way.

We calculate the ζ -greeks by means of a difference quotient. As pricing input for the barrier options we do not use the estimate of the IVS obtained by the SFM, as it is necessarily subject to an estimation error. Instead, in order to avoid mispricings, we use the truly observed ones. Thus, by the definition of the ζ -greeks, the approximations are given by

$$\frac{\partial V}{\partial Z_{t,l}} \stackrel{\text{def}}{\approx} \frac{V(S, \hat{\sigma} \exp(\Delta Z_{t,l} \hat{m}_l)) - V(S, \hat{\sigma} \exp(-\Delta Z_{t,l} \hat{m}_l))}{2\Delta Z_{t,l}}. \quad (12)$$

In the practical implementation of (12) one faces a couple of numerical issues, which need to be addressed. First, the size of the $\Delta Z_{t,l}$ has to be chosen. An increment too small or too large can distort the meaning of the greeks. Moreover it cannot be unique for all $Z_{t,l}$, since the shift size depends on the basis functions \hat{m}_l and on the IVS on a particular day. Therefore we choose for each t a $\Delta Z_{t,l}$ such that the (absolute) mean upward (downward) shift amounts approximately to one volatility-point. Note that we do *not* use $\hat{Z}_{t,l}$ for these perturbations. Another challenge is an accurate calculation of the barrier greeks. To reduce

numerical errors we employ a constant grid in the pricing algorithm for calculating the ζ -greeks. Furthermore, the IVS $\hat{\sigma}$ needs to be arbitrage-free. However, the shifted surfaces do not necessarily possess this property. We thus additionally check no-arbitrage conditions before calculating the ζ -greeks and apply an algorithm due to Fengler (2008) in case of violations. This method estimates the option price function by means of a natural smoothing spline under no-arbitrage constraints, i.e. under convexity, monotonicity and bounds on the price function and on the first order strike derivatives. The resulting estimate is then converted back to implied volatility. The algorithm is not applied when computing vega and vanna since parallel shifts do typically not result into arbitrage violations.

The aforementioned greeks are demonstrated in Figure 4 for the down-and-out put with half a year to expiry. The plot displays the greeks as a function of spot and keeps other characteristics of the barrier option unchanged. It has to be noted that the SFM, i.e. $\hat{Z}_{t,l}$ and \hat{m}_l , can only be identified up to sign. The sign of the ζ -greeks therefore has no particular meaning. Hence vega and ζ_1 display similar patterns. For the spot values close to the barrier level vega is negative and approaches zero as it becomes a delta product. For out-of-the money options vega is positive since the option then resembles a plain vanilla contract. A similar behavior is observed for ζ_2 and vanna, but the vanna is discontinuous at the barrier as it is derived from the delta.

4 Empirical Results

4.1 Data

The data set covers DAX index options traded at the EUREX from January 3rd, 2000 till June 30th, 2004 which give 1135 trading days. We use settlement prices, which are prices published by the EUREX based on the last intra-day trades. The DAX index is a capital weighted performance index comprising 30 German blue chips. Since dividends less corporate tax are reinvested into the index, they do not need to be taken into account for option valuation.

We preprocess the data by eliminating implied volatilities bigger than 80% and maturities

smaller than 10 days. Arbitrage violations in the option data are removed by the arbitrage free smoothing procedure described in Fengler (2008). After smoothing, the data are converted into a regular grid of moneyness and time to maturities. For option pricing, the zero rates from EURIBOR quotes are linearly interpolated, see Dumas et al. (1998) for this practice.

4.2 Experimental Design

In our empirical study we assume no transaction costs, no restrictions on short selling and the possibility of trading each asset at arbitrary size. Each security is priced using the LV model calibrated to daily market data. We implement the hedging strategies described in Section 3, i.e. we focus exclusively on volatility and spot risks, leaving other risks like interest rate exposure unhedged.

In the first step of our experiment we estimate the SFM. As kernel function we use a product quartic kernel, where $k(u) = 15/16(1 - u^2)^2$ for $|u| < 1$ and 0 otherwise. For a data driven bandwidth choice and the model size selection, we refer to Fengler et al. (2007). The basic idea is to estimate the model for different combinations of L and h and compare various information criteria. For the moneyness direction we finally use a bandwidth of 0.04, but we slightly oversmooth the surfaces in the time to maturity direction in order to reduce numerical errors for the subsequent price computations. More precisely, we use a local bandwidth modelled by an arctangent function which increases monotonously from 0.02 to 0.15 (expressed in years). Since in the hedging procedure only two main factors are included, we set $L = 2$. With this choice the model describes sufficiently well the IVS dynamics, since the measure of explained variation is close to 98%.

For each day up to one year before the last observation date in the sample, a long position in the reverse barrier option is created. This is to evaluate all initiated hedges at market prices within the sample. We use up-and-out calls with strikes at 80% of the spot and barriers at 140% and down-and-out put with strikes at 80% and barriers at 110%. These specifications correspond to typically traded contracts. Based on the calibrated LV model, ζ -greeks, delta, vega and vanna are calculated and the hedging strategies as described in Section 3 are set up. We concentrate on vega, vanna, ζ_1 and $\zeta_1\zeta_2$ strategies since the pure delta hedge is of

inferior quality. As HP_1 we use at-the-money puts for the up-and-out calls and at-the-money calls for the down-and-out puts. The risk reversal are structured by taking 80% and 120% strikes of the current spot.

Positions that have not knocked are updated on a daily basis. This choice is motivated by the results of Engelmann et al. (2006) who do not obtain different rankings of the strategies for other re-balancing frequencies. For each day we calculate the greeks to solve (7) and (11) and adjust the hedge ratios a_0, a_1, a_2 . The hedges are financed from the cash account and if the barrier is breached or the barrier option expires we unwind the hedge and record the hedging error. All positions are traded at market prices. In case of a knock-out event, the hedging error pays or earns interest until expiry in order to render the results comparable. Also the cash account bears interest or is financed at the riskless short rate of the concurrent trading day. Summing up, we have a collection of hedging errors for the two types of barrier options with four different hedging strategies for each of them.

One could object that the experimental design suffers from an in-sample problem, since the SFM is estimated on the same data set as the hedging experiment. It is however a common finding in the empirical literature, either on interest rates or on the IVS, that eigenvectors or eigenfunctions are remarkably stable across time. Formal tests on IVS data between the years 1995 to 2001 confirming this hypothesis are provided by Fengler (2005, Chapter 5.2.3). Even if we made use of a training-sample, we would therefore recover very similar factor functions. Thus the issue will not seriously affect the results.

4.3 Results

For evaluating the performance we use a pool of 885 hedging errors (1135 trading days less 250 days, since products issued thereafter would not expire within the sample). In order to make them comparable we normalize by the spot price at the time when the hedge is initiated. This normalization is common in practice and is meant to remove the dependence from the underlying's level. Another normalizing factor could be the option price itself, but since the risk reversal has a market price close to zero, measuring errors with respect to the spot appears to be more natural.

The aim of hedging is to replicate the payoff of the option. In the ideal case the hedge portfolio should have zero variance and zero mean, but for obvious reasons this cannot be realized in practice. Our aim is to give a comparative analysis of the hedging error distributions in order to check how the volatility factors affect the hedging performance. We use traditional descriptive statistics to assess the location and dispersion of the errors. Clearly, a superior method would keep these quantities close to zero in absolute terms.

The empirical results are summarized in Tables 2 and 3 for up-and-out calls and down-and-out puts respectively. We present the minimum, maximum, mean, median, standard deviation, and the absolute deviation around the median. The terminal hedging error distributions are given in the rows marked with a ‘0’. As can be inferred from the tables, the center of all distributions is located around zero, with means slightly below zero for the up-and-out calls and slightly above zero for the down-and-out puts. Thus the different hedges are hardly distinguishable in terms of the center of the distribution. This finding corresponds to our expectations: the volatility risk is removed, both for the vega and the ζ_1 -hedges, and vanna and $\zeta_1\zeta_2$ -hedges do not add any additional drift, since they are almost costless.

For evaluating the dispersion of the hedging errors we focus on the standard deviation and the absolute deviation around the median (madev.). The first observation is that hedges relying on higher order greeks tend to exhibit lower variance. In case of the down-and-out puts the vanna hedge has a slightly smaller dispersion than the $\zeta_1\zeta_2$ -hedge, and the traditional vega hedge performs very similar to the ζ_1 -hedge. For the up-and-out calls the ranking is reversed: the standard hedges are clearly outperformed by the factor hedges. How can this asymmetry be explained and how is the quality of the factor hedges to be judged?

There are two major sources of bias in the hedging strategies due to the behavior of the underlying. Observe that during the analyzed time period the DAX had a downward trend: 81% out of the down-and-out put options knocked out, but only 10% of the up-and-out call options, while 5% of the puts and 39% calls expired in-the-money, see Table 1. As a first issue consider the huge amount of up-and-out calls ending in-the-money. This gives rise to what is known among practitioners as ‘theta risk’. For explanation reconsider the case in Section 3, where we demonstrated that the prices for one-year down-and-out puts with a strike of 110% and barrier at 80% were less than 3% in the two scenarios. In contrast, when the put ends in-the-money it will pay out up to 30%. Consequently, the value of an in-the-money reverse

barrier option increases sharply the nearer the expiry date draws (i.e. has a strong theta), rendering it more and more difficult for traders to earn the payoff by trading the gamma. Theta risk can thus lead to a more dispersed error distribution. A second issue is gap risk. We do not unwind the hedges at the barriers, but at the observed spots, since this is the more realistic scenario in practice. When a barrier is breached, one still owns the hedge and incurs unbalanced gains or losses. Again this leads to a more dispersed hedging error distribution. As is clear from Table 1, theta risk is dominating the risk in case of the calls and gap risk in case of the puts.

To receive a deeper insight, we refer once more to Tables 2 and 3. We report the statistics of the hedging experiment stopped at 1 day, 5 days and 25 days before the expiry. As is seen the dispersion measures increase the nearer expiry draws, and the distributions become less skewed and less heavy-tailed, while the location measures prove to remain stable. In terms of dispersion the relative order of the hedging strategies across the two products remains the same: for the down-and-out puts the strategies are comparable, while factor hedging remains superior for the up-and-out calls. This finding is confirmed in Figure 5, which displays the standard deviations of the hedging errors as a function through the options' life time. It is intuitive to expect this function to increase. Moreover there is a sharp jump just before the expiry date contributing a large portion of the overall cumulative hedging error in particular for the up-and-out calls. All these observations highlight the importance of the expiry effect relative to gap risk when interpreting the data.

We overall conclude two main findings. First, factor hedging is at least of similar quality as traditional hedging approaches. In particular the hedging efficiency does not deteriorate. This is a reassuring result given the huge computational effort that must be spent and that could easily come at the costs of accuracy. This result is obtained when the barrier options expire worthless or knock out early in life time. Second, when the option needs to be hedged till expiry and ends in-the-money, the factor hedging approach dominates clearly. From a trader's perspective the first situation is the 'easy one' unless the knock-out occurs close to expiry. The second one is much more intricate, because the intrinsic value needs to be earned. This is a strong case for volatility factor hedging.

5 Conclusion

We provide an empirical study on hedging reverse barrier options in the local volatility model. The main focus of this study is on risk factors arising from a decomposition of the dynamic behavior of the implied volatility surface, which are identified with a flexible semiparametric technique. The hedging framework is constructed as a natural extension to traditional vega hedging, where the sensitivity is measured with respect to the more complex surface movements.

Our empirical investigation shows that hedging higher order risk with risk reversals brings improvements to hedging with at-the-money plain vanillas only. This is consistent across the vanna hedge and the more complex factor based hedges, thus confirming evidence of Ederington and Guan (2007). Intuitively the vega hedge resembles a single factor based hedge since the first dynamic factor corresponds to a parallel type of shift. Adding a vanna hedge or another factor to the portfolio removes similar risks as can be inferred from the comparable hedging performance.

Measured in terms of the hedging error variance, factor hedging performs at least as good as the corresponding vega and vanna hedges, in certain cases it is superior. As is confirmed by hedging up-and-out call options and down-and-out put options, the first case occurs when options knock out early in life time or expire worthless, while the second occurs when the options need to be hedged up to expiry and end in-the-money. This evidence is present not only in the terminal hedging errors but also through the option's life time. From a trader's perspective the second case is the more interesting, making factor hedging a powerful alternative to traditional hedging.

These findings, however, are not necessarily similar for other complex derivatives sensitive to IVS movements, such as cliquets or long-dated forward starting options. Also a portfolio context may yield different findings. In particular, when a book of options contains assets with several maturities it could be beneficial to consider additional factors, such as those related to the term structure of the IVS. This exposure can be hedged by constructing the corresponding calendar spreads. Another application in a portfolio context could be stress test scenarios based on the volatility factors. This would provide a good understanding of

the volatility exposure of the portfolio. We leave these issues to future research.

References

- Alexander, C. (2001). Principles of the skew. *RISK*, 14(1):S29–S32.
- Andersen, L. B. G., Andreasen, J., and Eliezer, D. (2002). Static replication of barrier options: Some general results. *Journal of Computational Finance*, 5(4):1–25.
- Andersen, L. B. G. and Brotherton-Ratcliffe, R. (1997). The equity option volatility smile: An implicit finite-difference approach. *Journal of Computational Finance*, 1(2):5–37.
- Bates, D. S. (1996). Jumps and stochastic volatility: Exchange rate processes implicit in Deutsche Mark options. *Review of Financial Studies*, 9:69–107.
- Bergomi, L. (2005). Smile dynamics II. *RISK*, 18(10):67–73.
- Borak, S. and Weron, R. (2009). A semiparametric factor model for electricity forward curve dynamics. *The Journal of Energy Markets*, 1(3).
- Bouchouev, I. and Isakov, V. (1999). Uniqueness, stability and numerical methods for the inverse problem that arises in financial markets. *Inverse Problems*, 15:R95–R116.
- Brüggemann, R., Härdle, W., Mungo, J., and Trenkler, C. (2008). VAR modeling for dynamic loadings driving volatility strings. *Journal of Financial Econometrics*, 6:361–381.
- Carr, P. and Chou, A. (1997). Breaking Barriers. *Risk Magazine*, 10:139–145.
- Carr, P., Ellis, K., and Gupta, V. (1998). Static hedging of exotic options. *Journal of Finance*, 53(3):1165–1190.
- Carr, P., Geman, H., Madan, D., and Yor, M. (2003). Stochastic volatility for Lévy processes. *Mathematical Finance*, 13:345–382.
- Coleman, T. F., Kim, Y., Li, Y., and Verma, A. (2001). Dynamic hedging with a deterministic local volatility function model. *Journal of Risk*, 4(1):63–89.
- Cont, R. and da Fonseca, J. (2002). The dynamics of implied volatility surfaces. *Quantitative Finance*, 2(1):45–60.

- Derman, E., Ergener, D., and Kani, I. (1995). Static options replication. *Journal of Derivatives*, 2(4):78–95.
- Derman, E. and Kani, I. (1994). Riding on a smile. *RISK*, 7(2):32–39.
- Diebold, F., Ji, L., and Li, C. (2006). A three-factor yield curve model: Non-affine structure, systematic risk sources, and generalized duration. In Klein, L., editor, *Long-Run Growth and Short-Run Stabilization: Essays in Memory of Albert Ando*. Edward Elgar, Cheltenham, U.K.
- Dumas, B., Fleming, J., and Whaley, R. E. (1998). Implied volatility functions: Empirical tests. *Journal of Finance*, 80(6):2059–2106.
- Dupire, B. (1994). Pricing with a smile. *RISK*, 7(1):18–20.
- Ederington, L. and Guan, W. (2007). Higher order greeks. *Journal of Derivatives*, 14:7–34.
- Engelmann, B., Fengler, M., Nalholm, M., and Schwendner, P. (2007). Static versus Dynamic Hedges: An Empirical Comparison for Barrier Options. *Review of Derivatives Research*, 9(3):239–264.
- Engelmann, B., Fengler, M., and Schwendner, P. (2006). Better than its reputation: An empirical hedging analysis of the local volatility model for barrier options. Working paper, Available at SRRN.
- Fengler, M. R. (2005). *Semiparametric Modeling of Implied Volatility*. Lecture Notes in Finance. Springer-Verlag, Berlin, Heidelberg.
- Fengler, M. R. (2008). Arbitrage-free smoothing of the implied volatility surface. *Quantitative Finance*. Forthcoming.
- Fengler, M. R., Härdle, W., and Mammen, E. (2007). A semiparametric factor model for implied volatility surface dynamics. *Journal of Financial Econometrics*, 5(2):189–218.
- Fengler, M. R., Härdle, W., and Villa, C. (2003). The dynamics of implied volatilities: A common principle components approach. *Review of Derivatives Research*, 6:179–202.

- Giacomini, E. and Härdle, W. (2008). Dynamic Semiparametric Factor Models in Pricing Kernels Estimation. In Dabo-Niang, S. and Ferraty, F., editors, *Functional and Operatorial Statistics*, pages 181–187. Physica-Verlag HD.
- Hafner, R. (2004). *Stochastic Implied Volatility*. Springer, Berlin.
- Hagan, P., Kumar, D., Lesniewski, A., and Woodward, D. (2002). Managing smile risk. *Wilmott magazine*, 1:84–108.
- Heston, S. (1993). A closed-form solution for options with stochastic volatility with applications to bond and currency options. *Review of Financial Studies*, 6:327–343.
- Hirsa, A., Courtadon, G., and Madan, D. (2003). The effect of model risk on the valuation of barrier options. *Journal of Risk Finance*, 4:47–55.
- Hull, J. and White, A. (1987). The pricing of options on assets with stochastic volatilities. *Journal of Finance*, 42:281–300.
- Hull, J. C. and Suo, W. (2002). A methodology for assessing model risk and its application to the implied volatility function model. *Journal of Financial and Quantitative Analysis*, 37(2):297–318.
- Kou, S. G. (2002). A jump-diffusion model for option pricing. *Management Science*, 48:1086–1101.
- Maruhn, J., Nalholm, M., and Fengler, M. R. (2008). Empirically robust static uncertain skew hedges for reverse barrier options. Working paper.
- McIntyre, M. L. (2001). Performance of Dupire’s implied diffusion approach under sparse and incomplete data. *Journal of Computational Finance*, 4(4):33–84.
- Merton, R. C. (1976). Option pricing when underlying stock returns are discontinuous. *Journal of Financial Economics*, 3:125–144.
- Nalholm, M. and Poulsen, R. (2006a). Static hedging and model risk for barrier options. *Journal of Future Markets*, 26:449–463.
- Nalholm, M. and Poulsen, R. (2006b). Static hedging of barrier options under general asset dynamics: Unification and application. *Journal of Derivatives*, 13:46–60.

- Park, B., Mammen, E., Härdle, W., and Borak, S. (2009). Time Series Modelling with Semiparametric Factor Dynamics. *Journal of the American Statistical Association*. Forthcoming.
- Rubinstein, M. (1994). Implied binomial trees. *Journal of Finance*, 49:771–818.
- Rubinstein, M. and Reiner, E. (1991). Breaking down the barrier. *RISK*, 4(9):28–35.
- Skiadopoulos, G., Hodges, S., and Clewlow, L. (1999). The dynamics of the S&P 500 implied volatility surface. *Review of Derivatives Research*, 3:263–282.
- Stein, E. M. and Stein, J. C. (1991). Stock price distributions with stochastic volatility: An analytic approach. *Review of Financial Studies*, 4:727–752.
- Taleb, N. (1997). *Dynamic Hedging: Managing Vanilla and Exotic Options*. John Wiley & Sons.
- Tompkins, R. (2002). Static versus dynamic hedging of exotic option: An evaluation of hedge performance via simulation. *The Journal of Risk Finance*, 3:6–34.
- Trück, S., Borak, S., Härdle, W., and Weron, R. (2006). Convenience yields for CO₂ emission allowance futures contracts. Discussion Paper 2006-076, SFB 649, Humboldt-Universität zu Berlin.
- Vähämaa, S. (2004). Delta hedging with the smile. *Financial Markets and Portfolio Management*, 18(3):241–255.

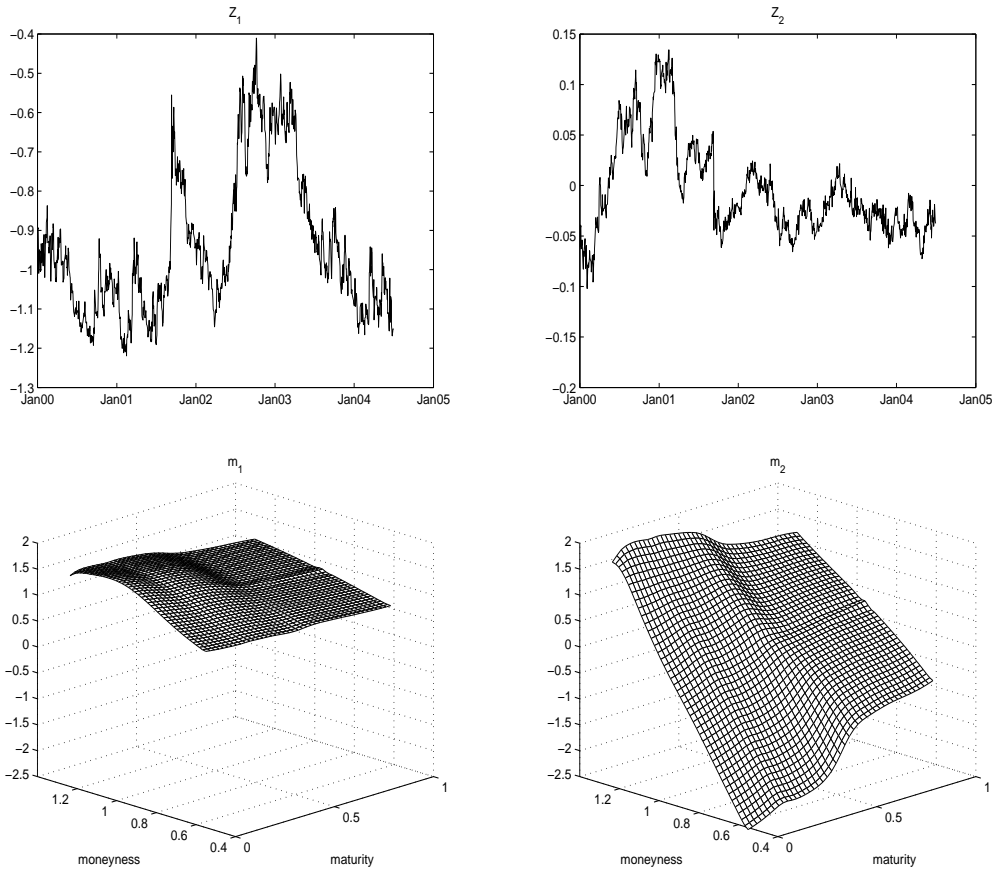


Figure 1: *The estimates of the SFM obtained from IVS data from January 3rd, 2000 till June 30th, 2004 for $L = 2$. Upper panel: estimated latent factor series \widehat{Z}_1 and \widehat{Z}_2 . Lower panel: estimates of \widehat{m}_1 , the non-uniform up-and-down shift, and \widehat{m}_2 , the slope risk.*

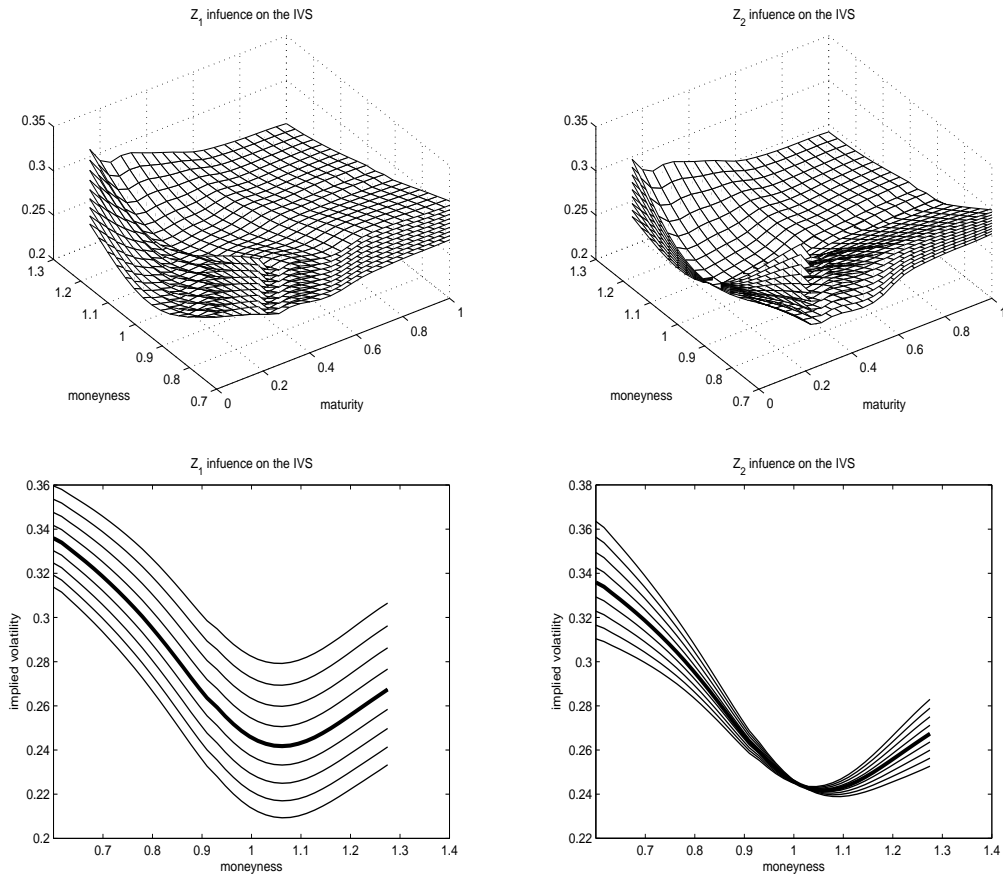


Figure 2: Impact of \widehat{Z}_1 and \widehat{Z}_2 on the IVS. Shocks in \widehat{Z}_1 trigger up-and-down movements while shocks in \widehat{Z}_2 tilt the smile around at-the-money point. Upper panel: a visualization of the shocks for the entire surface. Lower panel: the impact presented on one particular smile.

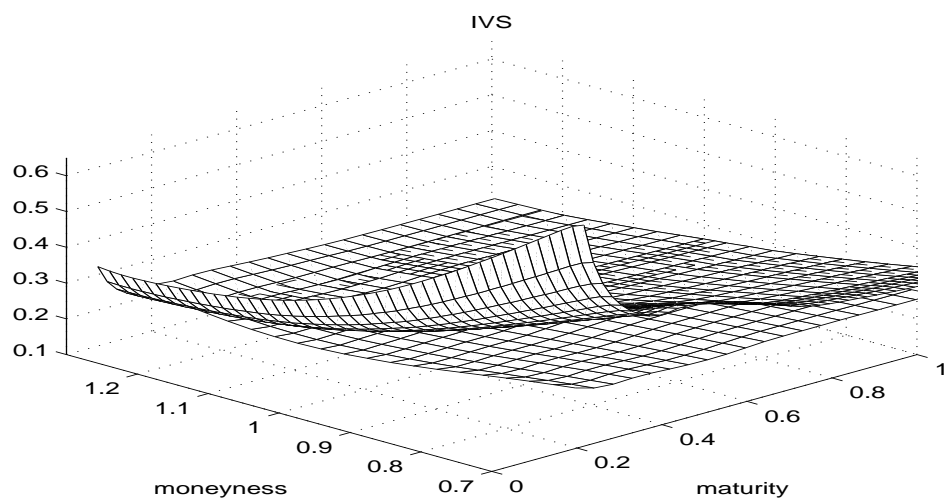


Figure 3: *IVS* observed on January 3rd, 2000 (the steeper surface) and January 2nd, 2001 (the flatter one). DAX levels on these days were 6751 and 6290 respectively.

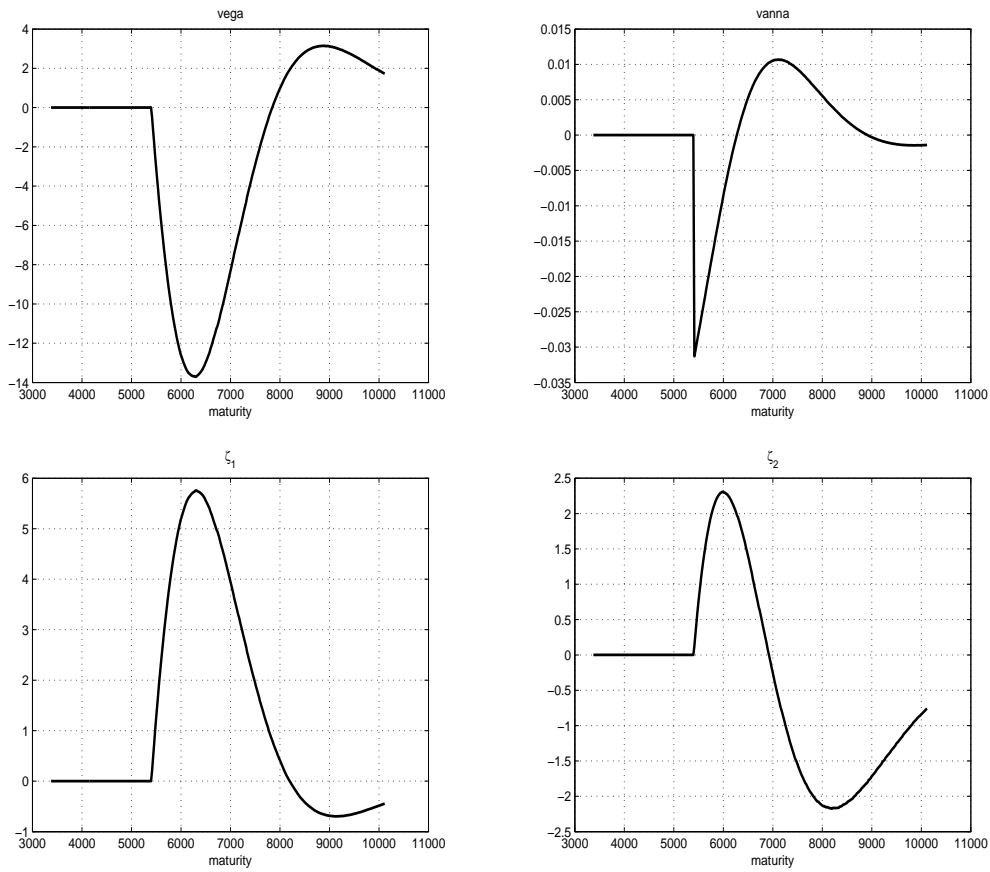


Figure 4: Greeks for a down-and-out put option with maturity 0.5 years with barrier 5400 strike 7425 as a function of the spot. Upper left panel: vega. Upper right panel: vanna. Lower right panel: ζ_1 . Lower right panel: ζ_2

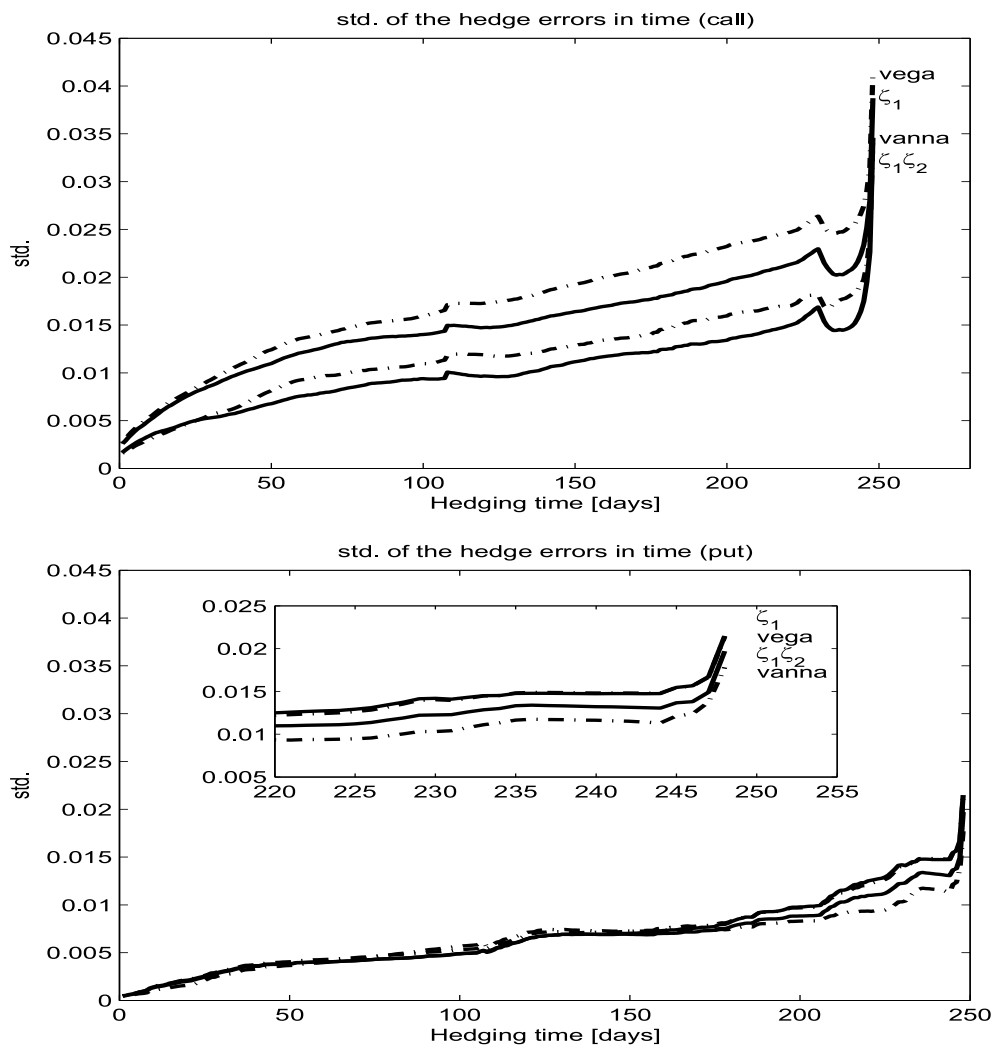


Figure 5: *Standard deviations of the hedging errors as a function of time from option issuance. Solid lines represent the factor hedging methods motivated by the SFM. Dashed lines represent the vega and vanna hedges. Upper panel: up-and-out call. Lower panel: down-and-out put.*

option type	barrier	strike	knock-outs	in-the-money
up-and-out call	140%	80%	10%	39%
down-and-out put	80%	110%	81%	5%

Table 1: *Characteristics of the analyzed barrier options. Strikes and barriers are in percentage of spot at issuance. The column ‘knock-outs’ refers to the contracts that breached the barrier and ‘in-the-money’ to those yielding a positive payoff at expiry.*

	days	min	max	mean	median	std.	madev.	skew.	kurt.
vega	0	-0.1038	0.5813	-0.0165	-0.0175	0.0413	0.0209	7.2801	97.60
	1	-0.1038	0.2581	-0.0172	-0.0174	0.0314	0.0199	2.3526	19.32
	5	-0.1037	0.0970	-0.0181	-0.0169	0.0260	0.0183	0.3636	4.91
	25	-0.0827	0.0649	-0.0174	-0.0164	0.0249	0.0178	0.0587	3.74
ζ_1	0	-0.0752	0.5768	-0.0118	-0.0136	0.0387	0.0183	8.4877	119.04
	1	-0.0751	0.2332	-0.0125	-0.0134	0.0279	0.0172	2.8026	22.27
	5	-0.0749	0.0755	-0.0134	-0.0121	0.0216	0.0155	0.2846	4.60
	25	-0.0761	0.0573	-0.0134	-0.0127	0.0215	0.0161	0.0343	3.55
vanna	0	-0.1340	0.5310	-0.0081	-0.0138	0.0345	0.0151	8.6289	124.62
	1	-0.1340	0.1842	-0.0089	-0.0136	0.0239	0.0140	2.1325	17.40
	5	-0.1339	0.0807	-0.0099	-0.0131	0.0187	0.0121	0.2157	9.54
	25	-0.0582	0.0772	-0.0096	-0.0141	0.0173	0.0118	1.3367	6.22
$\zeta_1\zeta_2$	0	-0.0830	0.5684	-0.0066	-0.0119	0.0345	0.0137	10.6470	161.00
	1	-0.0829	0.2091	-0.0073	-0.0117	0.0226	0.0126	4.0718	31.51
	5	-0.0829	0.0710	-0.0083	-0.0113	0.0157	0.0106	1.3447	7.12
	25	-0.0370	0.0629	-0.0086	-0.0118	0.0152	0.0108	1.3559	5.72

Table 2: Hedging error distributions of the up-and-out calls. Given are descriptive statistics for the various hedging strategies. The rows present the statistics at 0, 1, 5 and 25 days before expiration.

	days	min	max	mean	median	std.	madev.	skew.	kurt.
vega	0	-0.0264	0.2799	0.0058	-0.0004	0.0213	0.0105	5.4903	51.91
	1	-0.0756	0.1172	0.0050	-0.0004	0.0166	0.0098	2.4531	12.38
	5	-0.0187	0.0882	0.0041	-0.0008	0.0147	0.0090	2.4682	10.92
	25	-0.0186	0.0749	0.0038	-0.0004	0.0124	0.0083	2.0267	8.73
ζ_1	0	-0.0210	0.2808	0.0080	0.0016	0.0214	0.0107	5.4775	51.59
	1	-0.0702	0.1215	0.0072	0.0015	0.0167	0.0100	2.4501	12.14
	5	-0.0137	0.0882	0.0063	0.0013	0.0147	0.0091	2.4112	10.29
	25	-0.0113	0.0798	0.0059	0.0014	0.0127	0.0085	2.0632	8.54
vanna	0	-0.0608	0.2072	0.0022	-0.0016	0.0178	0.0081	5.7326	53.32
	1	-0.0955	0.1309	0.0014	-0.0016	0.0137	0.0074	2.7735	23.76
	5	-0.0323	0.0649	0.0006	-0.0018	0.0114	0.0069	1.9306	9.98
	25	-0.0205	0.0582	0.0004	-0.0016	0.0093	0.0059	2.0273	9.23
$\zeta_1\zeta_2$	0	-0.0332	0.2676	0.0065	0.0008	0.0196	0.0092	6.0258	60.47
	1	-0.0824	0.1146	0.0057	0.0008	0.0149	0.0085	2.4463	14.06
	5	-0.0234	0.0774	0.0048	0.0007	0.0130	0.0079	2.4134	10.42
	25	-0.0121	0.0727	0.0045	0.0010	0.0110	0.0069	2.3842	10.27

Table 3: Hedging error distributions of the down-and-out puts. Given are descriptive statistics for the various hedging strategies. The rows present the statistics at 0, 1, 5 and 25 days before the expiration.

Adaptive pointwise estimation in time-inhomogeneous conditional heteroscedasticity models

P. ČÍŽEK[†], W. HÄRDLE[‡] AND V. SPOKOINY[§]

[†]*Department of Econometrics & OR, Tilburg University, P.O. Box 90153, 5000LE Tilburg, The Netherlands*

E-mail: P.Cizek@uvt.nl

[‡]*Humboldt-Universität zu Berlin and CASE, Spandauerstrasse 1, 10178 Berlin, Germany*

E-mail: haerdle@wiwi.hu-berlin.de

[§]*Weierstrass-Institute, Humboldt-Universität zu Berlin and CASE, Mohrenstrasse 39, 10117 Berlin, Germany*

E-mail: spokoiny@wias-berlin.de

First version received: April 2008; final version accepted: April 2009

Summary This paper offers a new method for estimation and forecasting of the volatility of financial time series when the stationarity assumption is violated. Our general, local parametric approach particularly applies to general varying-coefficient parametric models, such as GARCH, whose coefficients may arbitrarily vary with time. Global parametric, smooth transition and change-point models are special cases. The method is based on an adaptive pointwise selection of the largest interval of homogeneity with a given right-end point by a local change-point analysis. We construct locally adaptive estimates that can perform this task and investigate them both from the theoretical point of view and by Monte Carlo simulations. In the particular case of GARCH estimation, the proposed method is applied to stock-index series and is shown to outperform the standard parametric GARCH model.

Keywords: *Adaptive pointwise estimation, Autoregressive models, Conditional heteroscedasticity models, Local time-homogeneity.*

1. INTRODUCTION

A growing amount of econometrical and statistical research is devoted to modelling financial time series and their volatility, which measures dispersion at a point in time (i.e. conditional variance). Although many economies and financial markets have been recently experiencing many shorter and longer periods of instability or uncertainty such as the Asian crisis (1997), the Russian crisis (1998), the start of the European currency (1999), the ‘dot-Com’ technology-bubble crash (2000–02) or the terrorist attacks (September, 2001), the war in Iraq (2003) and the current global recession (2008), mostly used econometric models are based on the assumption of time homogeneity. This includes linear and non-linear autoregressive (AR) and moving-average models and conditional heteroscedasticity (CH) models such as ARCH (Engel, 1982)

and GARCH (Bollerslev, 1986), stochastic volatility models (Taylor, 1986), as well as their combinations such as AR-GARCH.

On the other hand, the market and institutional changes have long been assumed to cause structural breaks in financial time series, which was confirmed, e.g. in data on stock prices (Andreou and Ghysels, 2002, and Beltratti and Morana, 2004) and exchange rates (Herwatz and Reimers, 2001). Moreover, ignoring these breaks can adversely affect the modelling, estimation and forecasting of volatility as suggested e.g. by Diebold and Inoue (2001), Mikosch and Starica (2004), Pesaran and Timmermann (2004) and Hillebrand (2005). Such findings led to the development of the change-point analysis in the context of CH models; see e.g. Chen and Gupta (1997), Kokoszka and Leipus (2000) and Andreou and Ghysels (2006).

An alternative approach lies in relaxing the assumption of time homogeneity and allowing some or all model parameters to vary over time (Chen and Tsay, 1993, Cai et al., 2000, and Fan and Zhang, 2008). Without structural assumptions about the transition of model parameters over time, time-varying coefficient models have to be estimated non-parametrically, e.g. under the identification condition that their parameters are smooth functions of time (Cai et al., 2000). In this paper, we follow a different strategy based on the assumption that a time series can be locally, i.e. over short periods of time, approximated by a parametric model. As suggested by Spokoiny (1998), such a local approximation can form a starting point in the search for the longest period of stability (homogeneity), i.e. for the longest time interval in which the series is described well by the parametric model. In the context of the local constant approximation, this strategy was employed for volatility modelling by Härdle et al. (2003), Mercurio and Spokoiny (2004) and Spokoiny (2009a). Our aim is to generalize this approach so that it can identify intervals of homogeneity for any parametric CH model regardless of its complexity.

In contrast to the local constant approximation of the volatility of a process (Mercurio and Spokoiny, 2004), the main benefit of the proposed generalization consists in the possibility to apply the methodology to a much wider class of models and to forecast over a longer time horizon. The reason is that approximating the mean or volatility process by a constant is in many cases too restrictive or even inappropriate and it is fulfilled only for short time intervals, which precludes its use for longer-term forecasting. On the contrary, parametric models like GARCH mimic the majority of stylized facts about financial time series and can reasonably fit the data over rather long periods of time in many practical situations. Allowing for time dependence of model parameters offers then much more flexibility in modelling real-life time series, which can be both with or without structural breaks since global parametric models are included as a special case.

Moreover, the proposed adaptive local parametric modelling unifies the change-point and varying-coefficient models. First, since finding the longest time-homogeneous interval for a parametric model at any point in time corresponds to detecting the most recent change-point in a time series, this approach resembles the change-point modelling as in Bai and Perron (1998) or Mikosch and Starica (1999, 2004), for instance, but it does not require prior information such as the number of changes. Additionally, the traditional structural-change tests require that the number of observations before each break point is large (and can grow to infinity) as these tests rely on asymptotic results. On the contrary, the proposed pointwise adaptive estimation does not rely on asymptotic results and does not thus place any requirements on the number of observations before, between or after any break point. Second, since the adaptively selected time-homogeneous interval used for estimation necessarily differs at each time point, the model coefficients can arbitrarily vary over time. In comparison to varying-coefficient models assuming

smooth development of parameters over time (Cai et al., 2000), our approach however allows for structural breaks in the form of sudden jumps in parameter values.

Although seemingly straightforward, extending Mercurio and Spokoiny's (2004) procedure to the local parametric modelling is a non-trivial problem, which requires new tools and techniques. We concentrate here on the change-point estimation of financial time series, which are often modelled by data-demanding models such as GARCH. While the benefits of a flexible change-point analysis for time series spanning several years are well known, its feasibility (which stands in the focus of this work) is much more difficult to achieve. The reason is thus that, at each time point, the procedure starts from a small interval, where a local parametric approximation holds, and then iteratively extends this interval and tests it for time-homogeneity until a structural break is found or data exhausted. Hence, a model has to be initially estimated on very short time intervals (e.g. 10 observations). Using standard testing methods, such a procedure might be feasible for simple parametric models, but it is hardly possible for more complex parametric models such as GARCH that generally require rather large samples for reasonably good estimates.

Therefore, we use an alternative and more robust approach to local change-point analysis that relies on a finite-sample theory of testing a growing sequence of historical time intervals on homogeneity against a change-point alternative. The proposed adaptive pointwise estimation procedure applies to a wide class of time-series models, including AR and CH models. Concentrating on the latter, we describe in details the adaptive procedure, derive its basic properties, and focusing on the feasibility of adaptive estimation for CH models, study the performance in comparison to the parametric (G)ARCH by means of simulations and real-data applications. The main conclusion is two-fold: on one hand, the adaptive pointwise estimation is feasible and beneficial also in the case of data-demanding models such as GARCH; on the other hand, the adaptive estimates based on various parametric models such as constant, ARCH or GARCH models are much closer to each other (while being better than the usual parametric estimates), which eliminates to some extent the need for using too complex models in adaptive estimation.

The rest of the paper is organized as follows. In Section 2, the parametric estimation of CH models and its finite-sample properties are introduced. In Section 3, we define the adaptive pointwise estimation procedure and discuss the choice of its parameters. Theoretical properties of the method are discussed in Section 4. In the specific case of the ARCH(1) and GARCH(1,1) models, a simulation study illustrates the performance of the new methodology with respect to the standard parametric and change-point models in Section 5. Applications to real stock-index series data are presented in Section 6. The proofs are provided in the Appendix.

2. PARAMETRIC CONDITIONAL HETEROSCEDASTICITY MODELS

Consider a time series Y_t in discrete time, $t \in N$. The CH assumption means that $Y_t = \sigma_t \varepsilon_t$, where $\{\varepsilon_t\}_{t \in N}$ is a white noise process and $\{\sigma_t\}_{t \in N}$ is a predictable volatility (conditional variance) process. Modelling of the volatility process σ_t typically relies on some parametric CH specification such as the ARCH (Engle, 1982) and GARCH (Bollerslev, 1986) models:

$$\sigma_t^2 = \omega + \sum_{i=1}^p \alpha_i Y_{t-i}^2 + \sum_{j=1}^q \beta_j \sigma_{t-j}^2, \quad (2.1)$$

where $p \in N, q \in N$ and $\theta = (\omega, \alpha_1, \dots, \alpha_p, \beta_1, \dots, \beta_q)^\top$ is the parameter vector. An attractive feature of this model is that, even with very few coefficients, one can model most stylized facts of financial time series like volatility clustering or excessive kurtosis, for instance. A number of (G)ARCH extensions were proposed to make the model even more flexible; e.g. EGARCH (Nelson, 1991), QGARCH (Sentana, 1995) and TGARCH (Glosten et al., 1993) that account for asymmetries in a volatility process.

All such CH models can be put into a common class of generalized linear volatility models:

$$Y_t = \sigma_t \varepsilon_t = \sqrt{g(X_t)} \varepsilon_t, \quad (2.2)$$

$$X_t = \omega + \sum_{i=1}^p \alpha_i h(Y_{t-i}) + \sum_{j=1}^q \beta_j X_{t-j}, \quad (2.3)$$

where g and h are known functions and X_t is a (partially) unobserved process (structural variable) that models the volatility coefficient σ_t^2 via transformation $g : \sigma_t^2 = g(X_t)$. For example, the GARCH model (2.1) is described by $g(u) = u$ and $h(r) = r^2$.

Models (2.2)–(2.3) are time homogeneous in the sense that the process Y_t follows the same structural equation at each time point. In other words, the parameter θ and hence the structural dependence in Y_t is constant over time. Even though models like (2.2)–(2.3) can often fit data well over a longer period of time, the assumption of homogeneity is too restrictive in practical applications: to guarantee a sufficient amount of data for sufficiently precise estimation, these models are often applied over time spans of many years. On the contrary, the strategy pursued here requires only local time homogeneity, which means that at each time point t there is a (possibly rather short) interval $[t - m, t]$, where the process Y_t is well described by models (2.2)–(2.3). This strategy aims then both at finding an interval of homogeneity (preferably as long as possible) and at the estimation of the corresponding parameter values θ , which then enable predicting Y_t and X_t .

Next, we discuss the parameter estimation for models (2.2)–(2.3) using observations Y_t from some time interval $I = [t_0, t_1]$. The conditional distribution of each observation Y_t given the past \mathcal{F}_{t-1} is determined by the structural variable X_t , whose dynamics are described by the parameter vector $\theta : X_t = X_t(\theta)$ for $t \in I$ due to (2.3). We denote the underlying value of θ by θ_0 .

For estimating θ_0 , we apply the quasi-maximum likelihood (quasi-MLE) approach using the estimating equations generated under the assumption of Gaussian errors ε_t . This guarantees efficiency under the normality of innovations and consistency under rather general moment conditions (Hansen and Lee, 1994, and Francq and Zakoian, 2007). The log-likelihood for models (2.2)–(2.3) on an interval I can be represented in the form

$$L_I(\theta) = \sum_{t \in I} \ell\{Y_t, g[X_t(\theta)]\}$$

with log-likelihood function $\ell(y, v) = -0.5\{\log(v) + y^2/v\}$. We define the quasi-MLE estimate $\tilde{\theta}_I$ of the parameter θ by maximizing the log-likelihood $L_I(\theta)$,

$$\tilde{\theta}_I = \underset{\theta \in \Theta}{\operatorname{argmax}} L_I(\theta) = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{t \in I} \ell\{Y_t, g[X_t(\theta)]\}, \quad (2.4)$$

and denote by $L_I(\tilde{\theta}_I)$ the corresponding maximum.

To characterize the quality of estimating the parameter vector $\theta_0 = (\omega, \alpha_1, \dots, \alpha_p, \beta_1, \dots, \beta_q)^\top$ by $\tilde{\theta}_I$, we now present an exact (non-asymptotic) exponential risk bound. This bound concerns the value of maximum $L_I(\tilde{\theta}_I) = \max_{\theta \in \Theta} L_I(\theta)$ rather than the point of maximum $\tilde{\theta}_I$. More precisely, we consider the difference $L_I(\tilde{\theta}_I, \theta_0) = L_I(\tilde{\theta}_I) - L_I(\theta_0)$. By definition, this value is non-negative and represents the deviation of the maximum of the log-likelihood process from its value at the ‘true’ point θ_0 . Later, we comment on how the accuracy of estimation of the parameter θ_0 by $\tilde{\theta}_I$ relates to the value $L_I(\tilde{\theta}_I, \theta_0)$. We will also see that the bound for $L_I(\tilde{\theta}_I, \theta_0)$ yields the confidence set for the parameter θ_0 , which will be used for the proposed change-point test. Now, the non-asymptotic risk bound is specified in the following theorem, which formulates corollaries 4.2 and 4.3 of Spokoiny (2009b) for the case of the quasi-MLE estimation of a CH model (2.2)–(2.3) at $\theta = \theta_0$. The result can be viewed as an extension of the Wilks phenomenon that the distribution of $L_I(\tilde{\theta}_I, \theta_0)$ for a linear Gaussian model is $\chi_p^2/2$, where p is the number of estimated parameters in the model.

THEOREM 2.1. *Assume that the process Y_t follows models (2.2)–(2.3) with the parameter $\theta_0 \in \Theta$, where the set Θ is compact. The function $g(\cdot)$ is assumed to be continuously differentiable with the uniformly bounded first derivative and $g(x) \geq \delta > 0$ for all x . Further, let the process $X_t(\theta)$ be sub-ergodic in the sense that for any smooth function $f(\cdot)$ there exists f^* such that for any time interval I*

$$E_{\theta_0} \left| \sum_I \{f(X_t(\theta)) - E_{\theta_0} f(X_t(\theta))\} \right|^2 \leq f^* |I|, \quad \theta \in \Theta.$$

Finally, let $E \exp\{\varkappa(\varepsilon_t^2 - 1) | \mathcal{F}_{t-1}\} \leq c(\varkappa)$ for some $\varkappa > 0, c(\varkappa) > 0$, and all $t \in N$. Then there are $\lambda > 0$ and $\epsilon(\lambda, \theta_0) > 0$ such that for any interval I and $\mathfrak{z} > 0$

$$P_{\theta_0}(L_I(\tilde{\theta}_I, \theta_0) > \mathfrak{z}) \leq \exp\{\epsilon(\lambda, \theta_0) - \lambda \mathfrak{z}\}. \tag{2.5}$$

Moreover, for any $r > 0$, there is a constant $\mathfrak{R}_r(\theta_0)$ such that

$$E_{\theta_0} |L_I(\tilde{\theta}_I, \theta_0)|^r \leq \mathfrak{R}_r(\theta_0). \tag{2.6}$$

REMARK 2.1. The condition $g(x) \geq \delta > 0$ guarantees that the variance process cannot reach zero. In the case of GARCH, it is sufficient to assume $\omega > 0$, for instance.

One attractive feature of Theorem 2.1, formulated in the following corollary, is that it enables constructing the non-asymptotic confidence sets and testing the parametric hypothesis on the basis of the fitted log-likelihood $L_I(\tilde{\theta}_I, \theta)$. This feature is especially important for our procedure presented in Section 3.

COROLLARY 2.1. *Under the assumptions of Theorem 2.1, let the value \mathfrak{z}_α fulfil $\epsilon(\lambda, \theta_0) - \lambda \mathfrak{z}_\alpha < \log \alpha$ for some $\alpha < 1$. Then the random set $\mathcal{E}_I(\mathfrak{z}_\alpha) = \{\theta : L_I(\tilde{\theta}_I, \theta) \leq \mathfrak{z}_\alpha\}$ is an α -confidence set for θ_0 in the sense that $P_{\theta_0}(\theta_0 \notin \mathcal{E}_I(\mathfrak{z}_\alpha)) \leq \alpha$.*

Theorem 2.1 also gives a non-asymptotic and fixed upper bound for the risk of estimation $L_I(\tilde{\theta}_I, \theta_0)$ that applies to an arbitrary sample size $|I|$. To understand the relation of this result to the classical rate result, we can apply the standard arguments based on the quadratic expansion

of the log-likelihood $L(\tilde{\theta}, \theta)$. Let $\nabla^2 L(\theta)$ denote the Hessian matrix of the second derivatives of $L(\theta)$ with respect to the parameter θ . Then

$$L_I(\tilde{\theta}_I, \theta_0) = 0.5(\tilde{\theta}_I - \theta_0)^\top \nabla^2 L_I(\theta'_I)(\tilde{\theta}_I - \theta_0), \tag{2.7}$$

where θ'_I is a convex combination of θ_0 and $\tilde{\theta}_I$. Under usual regularity assumptions and for sufficiently large $|I|$, the normalized matrix $|I|^{-1} \nabla^2 L_I(\theta)$ is close to some matrix $V(\theta)$, which depends only on the stationary distribution of Y_t and is continuous in θ . Then (2.5) approximately means that $\|\sqrt{V(\theta_0)}(\tilde{\theta}_I - \theta_0)\|^2 \leq \mathfrak{z}/|I|$ with probability close to 1 for large \mathfrak{z} . Hence, the large deviation result of Theorem 2.1 yields the root- $|I|$ consistency of the MLE estimate $\tilde{\theta}_I$. See Spokoiny (2009b) for further details.

3. POINTWISE ADAPTIVE NON-PARAMETRIC ESTIMATION

An obvious feature of models (2.2)–(2.3) is that the parametric structure of the process is assumed constant over the whole sample and cannot thus incorporate changes and structural breaks at unknown times in the models. A natural generalization leads to models whose coefficients may change over time (Fan and Zhang, 2008). One can then assume that the structural process X_t satisfies the relation (2.3) at any time, but the vector of coefficients θ may vary with the time t , $\theta = \theta(t)$. The estimation of the coefficients as general functions of time is possible only under some additional assumptions on these functions. Typical assumptions are (i) varying coefficients are smooth functions of time (Cai et al., 2000) and (ii) varying coefficients are piecewise constant functions (Bai and Perron, 1998, and Mikosch and Starica, 1999, 2004).

Our local parametric approach differs from the commonly used identification assumptions (i) and (ii). We assume that the observed data Y_t are described by a (partially) unobserved process X_t due to (2.2), and at each point T , there exists a historical interval $I(T) = [t_0, T]$ in which the process X_t ‘nearly’ follows the parametric specification (2.3) (see Section 4 for details on what ‘nearly’ means). This local structural assumption enables us to apply well-developed parametric estimation for data $\{Y_t\}_{t \in I(T)}$ to estimate the underlying parameter $\theta = \theta(T)$ by $\hat{\theta} = \hat{\theta}(T)$. (The estimate $\hat{\theta} = \hat{\theta}(T)$ can then be used for estimating the value \hat{X}_T of the process X_t at T from equation (2.3) and for further modelling such as forecasting Y_{T+1} .) Moreover, this assumption includes the above-mentioned ‘smooth transition’ and ‘switching regime’ assumptions (i) and (ii) as special cases: parameters $\hat{\theta}(T)$ vary over time as the interval $I(T)$ changes with T and, at the same time, discontinuities and jumps in $\hat{\theta}(T)$ as a function of time are possible.

To estimate $\hat{\theta}(T)$, we have to find the historical interval of homogeneity $I(T)$, i.e. the longest interval I with the right-end point T , where data do not contradict a specified parametric model with fixed parameter values. Starting at each time T with a very short interval $I = [t_0, T]$, we search by successive extending and testing of interval I on homogeneity against a change-point alternative: if the hypothesis of homogeneity is not rejected for a given I , a larger interval is taken and tested again. Contrary to Bai and Perron (1998) and Mikosch and Starica (1999), who detect all change points in a given time series, our approach is local: it focuses on the local change-point analysis near point T of estimation and tries to find only one change closest to the reference point.

In the rest of this section, we first discuss the test statistics employed to test the time-homogeneity of an interval I against a change-point alternative in Section 3.1. Later, we rigorously describe the pointwise adaptive estimation procedure in Section 3.2. Its

implementation and the choice of parameters entering the adaptive procedure are described in Sections 3.2–3.4. Theoretical properties of the method are studied in Section 4.

3.1. Test of homogeneity against a change-point alternative

The pointwise adaptive estimation procedure crucially relies on the test of local time-homogeneity of an interval $I = [t_0, T]$. The null hypothesis for I means that the observations $\{Y_t\}_{t \in I}$ follow the parametric models (2.2)–(2.3) with a fixed parameter θ_0 , leading to the quasi-MLE estimate $\tilde{\theta}_I$ from (2.4) and the corresponding fitted log-likelihood $L_I(\tilde{\theta}_I)$.

The change-point alternative for a given change-point location $\tau \in I$ can be described as follows: process Y_t follows the parametric models (2.2)–(2.3) with a parameter θ_J for $t \in J = [t_0, \tau]$ and with a different parameter θ_{J^c} for $t \in J^c = [\tau + 1, T]$; $\theta_J \neq \theta_{J^c}$. The fitted log-likelihood under this alternative reads as $L_J(\tilde{\theta}_J) + L_{J^c}(\tilde{\theta}_{J^c})$. The test of homogeneity can be performed using the likelihood ratio (LR) test statistic $T_{I,\tau}$:

$$T_{I,\tau} = \max_{\theta_J, \theta_{J^c} \in \Theta} \{L_J(\theta_J) + L_{J^c}(\theta_{J^c})\} - \max_{\theta \in \Theta} L_I(\theta) = \{L_J(\tilde{\theta}_J) + L_{J^c}(\tilde{\theta}_{J^c}) - L_I(\tilde{\theta}_I)\}.$$

Since the change-point location τ is generally not known, we consider the supremum of the LR statistics $T_{I,\tau}$ over some subset $\tau \in \mathcal{T}(I)$; cf. Andrews (1993):

$$T_{I,\mathcal{T}(I)} = \sup_{\tau \in \mathcal{T}(I)} T_{I,\tau}. \quad (3.1)$$

A typical example of a set $\mathcal{T}(I)$ is $\mathcal{T}(I) = \{\tau : t_0 + m' \leq \tau \leq T - m''\}$ for some fixed $m', m'' > 0$.

3.2. Adaptive search for the longest interval of homogeneity

This section presents the proposed adaptive pointwise estimation procedure. At each point T , we aim at estimating the unknown parameters $\theta(T)$ from historical data $Y_t, t \leq T$; this procedure repeats for every current time point T as new data arrive. At the first step, the procedure selects on the base of historical data an interval $\hat{I}(T)$ of homogeneity in which the data do not contradict the parametric models (2.2)–(2.3). Afterwards, the quasi-MLE estimation is applied using the selected historical interval $\hat{I}(T)$ to obtain estimate $\hat{\theta}(T) = \tilde{\theta}_{\hat{I}(T)}$. From now on, we consider an arbitrary, but fixed time point T .

Suppose that a growing set $I_0 \subset I_1 \subset \dots \subset I_K$ of historical interval-candidates $I_k = [T - m_k + 1, T]$ with the right-end point T is fixed. The smallest interval I_0 is accepted automatically as homogeneous. Then the procedure successively checks every larger interval I_k on homogeneity using the test statistic $T_{I_k, \mathcal{T}(I_k)}$ from (3.1). The selected interval \hat{I} corresponds to the largest accepted interval $I_{\hat{k}}$ with index \hat{k} such that

$$T_{I_k, \mathcal{T}(I_k)} \leq \mathfrak{z}_k, \quad k \leq \hat{k}, \quad (3.2)$$

and $T_{I_{\hat{k}+1}, \mathcal{T}(I_{\hat{k}+1})} > \mathfrak{z}_{\hat{k}+1}$, where the critical values \mathfrak{z}_k are discussed later in this section and specified in Section 3.3. This procedure then leads to the adaptive estimate $\hat{\theta} = \tilde{\theta}_{\hat{I}}$ corresponding to the selected interval $\hat{I} = I_{\hat{k}}$.

The complete description of the procedure includes two steps. (A) Fixing the set-up and the parameters of the procedure. (B) Data-driven search for the longest interval of homogeneity.

(A) Set-up and parameters:

- 1 Select specific parametric models (2.2)–(2.3) [e.g. constant volatility, ARCH(1), GARCH(1,1)].
- 2 Select the set $\mathcal{I} = (I_0, \dots, I_K)$ of interval-candidates, and for each $I_k \in \mathcal{I}$, the set $\mathcal{T}(I_k)$ of possible change points $\tau \in I_k$ used in the LR test (3.1).
- 3 Select the critical values $\mathfrak{z}_1, \dots, \mathfrak{z}_K$ in (3.2) as described in Section 3.3.

(B) Adaptive search and estimation: Set $k = 1$, $\hat{I} = I_0$ and $\hat{\theta} = \tilde{\theta}_{I_0}$.

- 1 Test the hypothesis $H_{0,k}$ of no change point within the interval I_k using test statistics (3.1) and the critical values \mathfrak{z}_k obtained in (A3). If a change point is detected ($H_{0,k}$ is rejected), go to (B3). Otherwise proceed with (B2).
- 2 Set $\hat{\theta} = \tilde{\theta}_{I_k}$ and $\hat{I}_k = \tilde{I}_k$. Further, set $k := k + 1$. If $k \leq K$, repeat (B1); otherwise go to (B3).
- 3 Define $\hat{I} = I_{k-1}$ = ‘the last accepted interval’ and $\hat{\theta} = \tilde{\theta}_{\hat{I}}$. Additionally, set $\hat{\theta}_{I_k} = \dots = \hat{\theta}_{I_K} = \hat{\theta}$ if $k \leq K$.

In step (A), one has to select three main ingredients of the procedure. First, the parametric model used locally to approximate the process Y_t has to be specified in (A1), e.g. the constant volatility or GARCH(1,1) in our context. Next, in step (A2), the set of intervals $\mathcal{I} = \{I_k\}_{k=0}^K$ is fixed, each interval with the right-end point T , length $m_k = |I_k|$, and the set $\mathcal{T}(I_k)$ of tested change points. Our default proposal is to use a geometric grid $m_k = \lceil m_0 a^k \rceil$, $a > 1$, and to set $I_k = [T - m_k + 1, T]$ and $\mathcal{T}(I_k) = [T - m_{k-1} + 1, T - m_{k-2}]$. Although our experiments show that the procedure is rather insensitive to the choice of m_0 and a (e.g. we use $m_0 = 10$ and $a = 1.25$ in simulations), the length m_0 of interval I_0 should take into account the parametric model selected in (A1). The reason is that I_0 is always assumed to be time-homogeneous and m_0 thus has to reflect flexibility of the parametric model; e.g. while $m_0 = 20$ might be reasonable for the GARCH(1,1) model, $m_0 = 5$ could be a reasonable choice for the locally constant approximation of a volatility process. Finally, in step (A3), one has to select the K critical values \mathfrak{z}_k in (3.2) for the LR test statistics $T_{I_k, \mathcal{T}(I_k)}$ from (3.1). The critical values \mathfrak{z}_k will generally depend on the parametric model describing the null hypothesis of time-homogeneity, the set \mathcal{I} of intervals I_k and corresponding sets of considered change points $\mathcal{T}(I_k)$, $k \leq K$, and additionally, on two constants r and ρ that are counterparts of the usual significance level. All these determinants of the critical values can be selected in step (A) and the critical values are thus obtained before the actual estimation takes place in step (B). Due to its importance, the method of constructing critical values $\{\mathfrak{z}_k\}_{k=1}^K$ is discussed separately in Section 3.3.

The main step (B) performs the search for the longest time-homogeneous interval. Initially, I_0 is assumed to be homogeneous. If I_{k-1} is negatively tested on the presence of a change point, one continues with I_k by employing test (3.1) in step (B1), which checks for a potential change point in I_k . If no change point is found, then I_k is accepted as time-homogeneous in step (B2); otherwise the procedure terminates in step (B3). We sequentially repeat these tests until we find a change point or exhaust all intervals. The latest (longest) interval accepted as time-homogeneous is used for estimation in step (B3). Note that the estimate $\hat{\theta}_{I_k}$ defined in (B2) and (B3) corresponds to the latest accepted interval \hat{I}_k after the first k steps, or equivalently, the interval selected out of I_1, \dots, I_k .

Moreover, the whole search and estimation step (B) can be repeated at different time points T without reiterating the initial step (A) as the critical values \mathfrak{z}_k depend only on the approximating parametric model and interval lengths $m_k = |I_k|$, not on the time point T (see Section 3.3).

3.3. Choice of critical values \mathfrak{z}_k

The presented method of choosing the interval of homogeneity \hat{I} can be viewed as multiple testing procedure. The critical values for this procedure are selected using the general approach of testing theory: to provide a prescribed performance of the procedure under the null hypothesis, i.e. in the pure parametric situation. This means that the procedure is trained on the data generated from the pure parametric time-homogeneous model from step (A1). The correct choice in this situation is the largest considered interval I_K and a choice $I_{\hat{k}}$ with $\hat{k} < K$ can be interpreted as a ‘false alarm’. We select the minimal critical values ensuring a small probability of such a false alarm. Our condition slightly differs though from the classical level condition because we focus on parameter estimation rather than on hypothesis testing.

In the pure parametric case, the ‘ideal’ estimate corresponds to the largest considered interval I_K . Due to Theorem 2.1, the quality of estimation of the parameter θ_0 by $\tilde{\theta}_{I_K}$ can be measured by the log-likelihood ‘loss’ $L_{I_K}(\tilde{\theta}_{I_K}, \theta_0)$, which is stochastically bounded with exponential and polynomial moments: $E_{\theta_0} |L_{I_K}(\tilde{\theta}_{I_K}, \theta_0)|^r \leq \mathfrak{R}_r(\theta_0)$. If the adaptive procedure stops earlier at some intermediate step $k < K$, we select instead of $\tilde{\theta}_{I_K}$ another estimate $\hat{\theta} = \tilde{\theta}_{I_k}$ with a larger variability. The loss associated with such a false alarm can be measured by the value $L_{I_k}(\tilde{\theta}_{I_k}, \hat{\theta}) = L_{I_k}(\tilde{\theta}_{I_k}) - L_{I_k}(\hat{\theta})$. The corresponding condition bounding the loss due to the adaptive estimation reads as

$$E_{\theta_0} |L_{I_k}(\tilde{\theta}_{I_k}, \hat{\theta})|^r \leq \rho \mathfrak{R}_r(\theta_0). \quad (3.3)$$

This is in fact an implicit condition on the critical values $\{\mathfrak{z}_k\}_{k=1}^K$, which ensures that the loss associated with the false alarm is at most the ρ -fraction of the log-likelihood loss of the ‘ideal’ or ‘oracle’ estimate $\tilde{\theta}_{I_K}$ for the parametric situation. The constant r corresponds to the power of the loss in (3.3), while ρ is similar in meaning to the test level. In the limit case when r tends to zero, this condition (3.3) becomes the usual level condition: $P_{\theta_0}(I_K \text{ is rejected}) = P_{\theta_0}(\tilde{\theta}_{I_k} \neq \hat{\theta}) \leq \rho$. The choice of the metaparameters r and ρ is discussed in Section 3.4.

A condition similar to (3.3) is imposed at each step of the adaptive procedure. The estimate $\hat{\theta}_{I_k}$ coming after the k steps of the procedure should satisfy

$$E_{\theta_0} |L_{I_k}(\tilde{\theta}_{I_k}, \hat{\theta}_{I_k})|^r \leq \rho_k \mathfrak{R}_r(\theta_0), \quad k = 1, \dots, K, \quad (3.4)$$

where $\rho_k = \rho k/K \leq \rho$. The following theorem presents some sufficient conditions on the critical values $\{\mathfrak{z}_k\}_{k=1}^K$ ensuring (3.4); recall that $m_k = |I_k|$ denotes the length of I_k .

THEOREM 3.1. *Suppose that $r > 0, \rho > 0$. Under the assumptions of Theorem 2.1, there are constants a_0, a_1, a_2 such that the condition (3.4) is fulfilled with the choice*

$$\mathfrak{z}_k = a_0 r \log(\rho^{-1}) + a_1 r \log(m_K/m_{k-1}) + a_2 \log(m_k), \quad k = 1, \dots, K.$$

Since K and $\{m_k\}_{k=1}^K$ are fixed, the \mathfrak{z}_k 's in Theorem 3.1 have a form $\mathfrak{z}_k = C + D \log(m_k)$ for $k = 1, \dots, K$ with some constant C and D . However, a practically relevant choice of these constants has to be done by Monte Carlo simulations. Note first that every particular choice of the coefficients C and D determines the whole set of the critical values $\{\mathfrak{z}_k\}_{k=1}^K$ and thus the local change-point procedure. For the critical values given by fixed (C, D) , one can run the procedure and observe its performance on the simulated data using the data-generating process (2.2)–(2.3); in particular, one can check whether the condition (3.4) is fulfilled. For any (sufficiently large) fixed value of C , one can thus find the minimal value $D(C) < 0$ of D that ensures (3.4).

Every corresponding set of critical values in the form $\beta_k = C + D(C) \log(m_k)$ is admissible. The condition $D(C) < 0$ ensures that the critical values decreases with k . This reflects the fact that a false alarm at an early stage of the algorithm is more crucial because it leads to the choice of a highly variable estimate. The critical values β_k for small k should thus be rather conservative to provide the stability of the algorithm in the parametric situation. To determine C , the value β_1 can be fixed by considering the false alarm at the first step of the procedure, which leads to estimation using the smallest interval I_0 instead of the ‘ideal’ largest interval I_K . The related condition (used in Section 5.1) reads as

$$E_{\theta_0} |L_{I_K}(\tilde{\theta}_{I_K}, \tilde{\theta}_{I_0})|^r \mathbf{1}(T_{I_1, \mathcal{I}(I_1)} > \beta_1) \leq \rho \mathfrak{R}_r(\theta_0)/K. \tag{3.5}$$

Alternatively, one could select a pair (C, D) that minimizes the resulting prediction error; see Section 3.4.

3.4. Selecting parameters r and ρ

The choice of critical values using inequality (3.4) additionally depends on two ‘metaparameters’ r and ρ . A simple strategy is to use conservative values for these parameters and the corresponding set of critical values (e.g. our default is $r = 1$ and $\rho = 1$). On the other hand, the two parameters are global in the sense that they are independent of T . Hence, one can also determine them in a data-driven way by minimizing some global forecasting error (Cheng et al., 2003). Different values of r and ρ may lead to different sets of critical values and hence to different estimates $\hat{\theta}^{(r, \rho)}(T)$ and to different forecasts $\hat{Y}_{T+h|T}^{(r, \rho)}$ of the future values Y_{T+h} , where h is the forecasting horizon. Now, a data-driven choice of r and ρ can be done by minimizing the following objective function:

$$(\hat{r}, \hat{\rho}) = \arg \min_{r > 0, \rho > 0} PE_{\Lambda, \mathcal{H}}(r, \rho) = \arg \min_{r, \rho} \sum_T \sum_{h \in \mathcal{H}} \Lambda(Y_{T+h}, \hat{Y}_{T+h|T}^{(r, \rho)}), \tag{3.6}$$

where Λ is a loss function and \mathcal{H} is the forecasting horizon set. For example, one can take $\Lambda_r(v, v') = |v - v'|^r$ for $r \in [1/2, 2]$. For daily data, the forecasting horizon could be one day, $\mathcal{H} = \{1\}$, or two weeks, $\mathcal{H} = \{1, \dots, 10\}$.

4. THEORETIC PROPERTIES

In this section, we collect basic results describing the quality of the proposed adaptive procedure. First, the definition of the procedure ensures the performance prescribed by (3.4) in the parametric situation. We however claimed that the adaptive pointwise estimation applies even if the process Y_t is only locally approximated by a parametric model. Therefore, we now define a locally ‘nearly parametric’ process, for which we derive an analogy of Theorem 2.1 (Section 4.1). Later, we prove certain ‘oracle’ properties of the proposed method (Section 4.2).

4.1. Small modelling bias condition

This section discusses the concept of a ‘nearly parametric’ case. To define it rigorously, we have to quantify the quality of approximating the true latent process X_t , which drives the observed data Y_t due to (2.2), by the parametric process $X_t(\theta)$ described by (2.3) for some $\theta \in \Theta$. Below

we assume that the innovations ε_t in the model (2.2) are independent and identically distributed and denote the distribution of $\sqrt{v}\varepsilon_t$ by P_v so that the conditional distribution of Y_t given \mathcal{F}_{t-1} is $P_{g(X_t)}$. To measure the distance of a data-generating process from a parametric model, we introduce for every interval $I_k \in \mathcal{I}$ and every parameter $\theta \in \Theta$ the random quantity

$$\Delta_{I_k}(\theta) = \sum_{t \in I_k} \mathcal{K}\{g(X_t), g[X_t(\theta)]\},$$

where $\mathcal{K}(v, v')$ denotes the Kullback–Leibler distance between P_v and $P_{v'}$. For CH models with Gaussian innovations ε_t , $\mathcal{K}(v, v') = -0.5\{\log(v/v') + 1 - v/v'\}$. In the parametric case with $X_t = X_t(\theta_0)$, we clearly have $\Delta_{I_k}(\theta_0) = 0$. To characterize the ‘nearly parametric case’, we introduce a {small modelling bias} (SMB) condition, which simply means that, for some $\theta \in \Theta$, $\Delta_{I_k}(\theta)$ is bounded by a small constant with a high probability. Informally, this means that the ‘true’ model can be well approximated on the interval I_k by the parametric one with the parameter θ . The best parametric fit (2.3) to the underlying model (2.2) on I_k can be defined by minimizing the value $E\Delta_{I_k}(\theta)$ over $\theta \in \Theta$ and $\tilde{\theta}_{I_k}$ can be viewed as its estimate.

The following theorem claims that the results on the accuracy of estimation given in Theorem 2.1 can be extended from the parametric case to the general non-parametric situation under the SMB condition. Let $\varrho(\hat{\theta}, \theta)$ be any loss function for an estimate $\hat{\theta}$.

THEOREM 4.1. *Let for some $\theta \in \Theta$ and some $\Delta \geq 0$*

$$E\Delta_{I_k}(\theta) \leq \Delta. \tag{4.1}$$

Then it holds for an estimate $\hat{\theta}$ constructed from the observations $\{Y_t\}_{t \in I_k}$ that

$$E \log(1 + \varrho(\hat{\theta}, \theta)/E\theta\varrho(\hat{\theta}, \theta)) \leq 1 + \Delta.$$

This general result applied to the quasi-MLE estimation with the loss function $L_I(\tilde{\theta}_I, \theta)$ yields the following corollary.

COROLLARY 4.1. *Let the SMB condition (4.1) hold for some interval I_k and $\theta \in \Theta$. Then*

$$E \log\left(1 + |L_{I_k}(\tilde{\theta}_{I_k}, \theta)|^r / \mathfrak{R}_r(\theta)\right) \leq 1 + \Delta,$$

where $\mathfrak{R}_r(\theta)$ is the parametric risk bound from (2.6).

This result shows that the estimation loss $|L_I(\tilde{\theta}_I, \theta)|^r$ normalized by the parametric risk $\mathfrak{R}_r(\theta)$ is stochastically bounded by a constant proportional to e^Δ . If Δ is not large, this result extends the parametric risk bound (Theorem 2.1) to the non-parametric situation under the SMB condition. Another implication of Corollary 4.1 is that the confidence set built for the parametric model (Corollary 2.1) continues to hold, with a slightly smaller coverage probability, under SMB.

4.2. The ‘oracle’ choice and the ‘oracle’ result

Corollary 4.1 suggests that the ‘optimal’ or ‘oracle’ choice of the interval I_k from the set I_1, \dots, I_K can be defined as the largest interval for which the SMB condition (4.1) still holds (for a given small $\Delta > 0$). For such an interval, one can neglect deviations of the underlying

process from a parametric model with a fixed parameter θ . Therefore, we say that the choice k^* is the ‘oracle’ choice if there exists $\theta \in \Theta$ such that

$$E \Delta_{I_{k^*}}(\theta) \leq \Delta \tag{4.2}$$

for a fixed $\Delta > 0$ and that (4.2) does not hold for $k > k^*$. Unfortunately, the underlying process X_t and, hence, the value Δ_{I_k} is unknown and the oracle choice cannot be implemented. The proposed adaptive procedure tries to mimic this oracle on the basis of available data using the sequential test of homogeneity. The final oracle result claims that the adaptive estimate provides the same (in order) accuracy as the oracle one.

By construction, the pointwise adaptive procedure described in Section 3 provides the prescribed performance if the underlying process follows the parametric model (2.2). Now, condition (3.4) combined with Theorem 4.1 implies similar performance in the first k^* steps of the adaptive estimation procedure.

THEOREM 4.2. *Let $\theta \in \Theta$ and $\Delta > 0$ be such that $E \Delta_{I_{k^*}}(\theta) \leq \Delta$ for some $k^* \leq K$. Also let $\max_{k \leq k^*} E_{\theta} |L_{I_k}(\tilde{\theta}_{I_k}, \theta)|^r \leq \mathfrak{R}_r(\theta)$. Then*

$$E \log \left(1 + \frac{|L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \theta)|^r}{\mathfrak{R}_r(\theta)} \right) \leq 1 + \Delta \quad \text{and} \quad E \log \left(1 + \frac{|L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \hat{\theta}_{I_{k^*}})|^r}{\mathfrak{R}_r(\theta)} \right) \leq \rho + \Delta.$$

Similarly to the parametric case, under the SMB condition $E \Delta_{I_{k^*}}(\theta) \leq \Delta$, any choice $\hat{k} < k^*$ can be viewed as a false alarm. Theorem 4.2 documents that the loss induced by such a false alarm at the first k^* steps and measured by $L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \hat{\theta}_{I_{k^*}})$ is of the same magnitude as the loss $L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \theta)$ of estimating the parameter θ from the SMB (4.2) by $\tilde{\theta}_{I_{k^*}}$. Thus, under (4.2) the adaptive estimation during steps $k \leq k^*$ does not induce larger errors into estimation than the quasi-MLE estimation itself.

For further steps of the algorithm with $k > k^*$, where (4.2) does not hold, the value $\Delta' = E \Delta_{I_k}(\theta)$ can be large and the bound for the risk becomes meaningless due to the factor $e^{\Delta'}$. To establish the result about the quality of the final estimate, we thus have to show that the quality of estimation cannot be destroyed at the steps $k > k^*$. The next ‘oracle’ result states the final quality of our adaptive estimate $\hat{\theta}$.

THEOREM 4.3. *Let $E \Delta_{I_{k^*}}(\theta) \leq \Delta$ for some $k^* \leq K$. Then $L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \hat{\theta}) \mathbf{1}(\hat{k} \geq k^*) \leq \mathfrak{z}_{k^*}$ yielding*

$$E \log \left(1 + \frac{|L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \hat{\theta})|^r}{\mathfrak{R}_r(\theta)} \right) \leq \rho + \Delta + \log \left(1 + \frac{\mathfrak{z}_{k^*}^r}{\mathfrak{R}_r(\theta)} \right).$$

Due to this result, the value $L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \hat{\theta})$ is stochastically bounded. This can be interpreted as the oracle property of $\hat{\theta}$ because it means that the adaptive estimate $\hat{\theta}$ belongs with a high probability to the confidence set of the oracle estimate $\tilde{\theta}_{I_{k^*}}$.

5. SIMULATION STUDY

In the last two sections, we present simulation study (Section 5) and real data applications (Section 6) documenting the performance of the proposed adaptive estimation procedure. To verify the practical applicability of the method in a complex setting, we concentrate on the volatility estimation using parametric and adaptive pointwise estimation of constant volatility, ARCH(1) and GARCH(1,1) models (for the sake of brevity, referred to as the local constant,

local ARCH and local GARCH). The reason is that the estimation of GARCH models requires generally hundreds of observations for a reasonable quality of estimation, which puts the adaptive procedure working with samples as small as 10 or 20 observations to a hard test. Additionally, the critical values obtained as described in Section 3.3 depend on the underlying parameter values in the case of (G)ARCH.

Here we first study the finite-sample critical values for the test of homogeneity by means of Monte Carlo simulations and discuss practical implementation details (Section 5.1). Later, we demonstrate the performance of the proposed adaptive pointwise estimation procedure in simulated samples (Section 5.2). Note that, throughout this section, we identify the GARCH(1,1) models by triplets (ω, α, β) : e.g. (1, 0.1, 0.3)-model. Constant volatility and ARCH(1) are then indicated by $\alpha = \beta = 0$ and $\beta = 0$, respectively. The GARCH estimation is done using the GARCH 3.0 package (Laurent and Peters, 2006) and Ox 3.30 (Doornik, 2002). Finally, since the focus is on modelling the volatility σ_t^2 in (2.2), the performance measurement and comparison of all models at time t is done by the absolute prediction error (PE) of the volatility process over a prediction horizon \mathcal{H} : $\text{APE}(t) = \sum_{h \in \mathcal{H}} |\sigma_{t+h}^2 - \hat{\sigma}_{t+h|t}^2| / |\mathcal{H}|$, where $\hat{\sigma}_{t+h|t}^2$ represents the volatility prediction by a particular model.

5.1. Finite-sample critical values for the test of homogeneity

A practical application of the pointwise adaptive procedure requires critical values for the test of local homogeneity of a time series. Since they are obtained under the null hypothesis that a chosen parametric model (locally) describes the data, see Section 3, we need to obtain the critical values for the constant volatility, ARCH(1) and GARCH(1,1) models. Furthermore, for given r and ρ , the average risk (3.4) between the adaptive and oracle estimates can be bounded for critical values that linearly depend on the logarithm of interval length $|I_k|$: $\mathfrak{z}(|I_k|) = \mathfrak{z}_k = C + D \log(|I_k|)$ (see Theorem 3.1). As described in Section 3.3, we choose here the smallest C satisfying (3.5) and the corresponding minimum admissible value $D = D(C) < 0$ that guarantees the conditions (3.4).

We simulated the critical values for ARCH(1) and GARCH(1,1) models with different values of underlying parameters; see Table 1 for the critical values corresponding to $r = 1$ and $\rho = 1$. Their simulation was performed sequentially on intervals with lengths ranging from $|I_0| = m_0 = 10$ to $|I_K| = 570$ observations using a geometric grid with multiplier $a = 1.25$; see Section 3.2. (The results are, however, not sensitive to the choice of a .)

Unfortunately, the critical values depend on the parameters of the underlying (G)ARCH model (in contrast to the constant-volatility model). They generally seem to increase with the values of the ARCH and GARCH parameters keeping the other one fixed; see Table 1. To deal with this dependence on the underlying model parameters, we propose to choose the largest (most conservative) critical values corresponding to any estimated parameter in the analysed data. For example, if the largest estimated parameters of GARCH(1,1) are $\hat{\alpha} = 0.3$ and $\hat{\beta} = 0.8$, one should use $\mathfrak{z}(10) = 26.4$ and $\mathfrak{z}(570) = 14.5$, which are the largest critical values for models with $\alpha = 0.3, \beta \leq 0.8$ and with $\alpha \leq 0.3, \beta = 0.8$. (The proposed procedure is, however, not overly sensitive to this choice, as we shall see later.)

Finally, let us have a look at the influence of the tuning constants r and ρ in (3.4) on the critical values for several selected models (Table 2). The influence is significant, but can be classified in the following way. Whereas increasing ρ generally leads to an overall decrease of critical values (cf. Theorem 3.1), but primarily for the longer intervals, increasing r leads to an increase of

Table 1. Critical values $\mathfrak{z}_k = \mathfrak{z}(|I_k|)$ of the supremum LR test.

$\mathfrak{z}(I_k)$	α	$ I_k $	β									
			0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
0.0		10	15.5	15.5	16.4	16.8	17.9	17.3	17.0	17.0	16.9	16.0
		570	5.5	7.2	7.0	7.0	7.5	7.5	7.4	7.3	7.0	6.7
0.1		10	16.3	14.5	15.1	15.9	16.4	15.9	16.1	16.0	16.0	
		570	8.6	9.0	9.1	9.6	9.8	10.7	11.5	12.5	14.0	
0.2		10	16.7	15.2	15.7	16.2	16.9	18.9	20.1	25.1		
		570	9.4	10.6	11.2	11.4	11.4	12.5	13.3	14.2		
0.3		10	18.5	16.4	16.7	16.9	18.1	21.8	26.4			
		570	9.7	10.8	12.0	12.4	12.9	13.5	14.5			
0.4		10	22.1	16.5	18.3	19.3	22.8	30.9				
		570	9.9	12.0	13.0	13.4	13.9	14.7				
0.5		10	26.2	19.1	19.5	25.4	38.1					
		570	10.7	12.6	13.8	14.0	14.6					
0.6		10	33.0	22.8	25.9	32.4						
		570	12.7	12.7	13.9	15.3						
0.7		10	41.1	24.8	29.1							
		570	16.8	14.7	16.1							
0.8		10	66.2	26.4								
		570	31.5	15.8								
0.9		10	88.6									
		570	60.9									

Note: $\omega = 1, r = 1$ and $\rho = 1$.

critical values mainly for the shorter intervals; cf. (3.4). In simulations and real applications, we verified that a fixed choice such as $r = 1$ and $\rho = 1$ performs well. To optimize the performance of the adaptive methods, one can however determine constants r and ρ in a data-dependent way as described in Section 3.3. We use here this strategy for a small grid of $r \in \{0.5, 1.0\}$ and $\rho \in \{0.5, 1.0, 1.5\}$ and find globally optimal r and ρ . We will document, though, that the differences in the average absolute PE (3.6) for various values of r and ρ are relatively small.

5.2. Simulation study

We aim (i) to examine how well the proposed estimation method is able to adapt to long stable (time-homogeneous) periods and to less stable periods with more frequent volatility changes and (ii) to see which adaptively estimated model—local volatility, local ARCH or local GARCH—performs best in different regimes. To this end, we simulated 100 series from two change-point GARCH models with a low GARCH effect ($\omega, 0.2, 0.1$) and a high GARCH effect ($\omega, 0.2, 0.7$). Changes in constant ω are spread over a time span of 1000 days; see Figure 1. There is a long stable period at the beginning (500 days ≈ 2 years) and end (250 days ≈ 1 year) of time series with several volatility changes between them.

Table 2. Critical values $\mathfrak{z}(|I_k|)$ of the supremum LR test for various values r and ρ .

Model (ω, α, β)		(0.1, 0.0, 0.0)		(0.1, 0.2, 0.0)		(0.1, 0.1, 0.8)	
r	ρ	$\mathfrak{z}(10)$	$\mathfrak{z}(570)$	$\mathfrak{z}(10)$	$\mathfrak{z}(570)$	$\mathfrak{z}(10)$	$\mathfrak{z}(570)$
1.0	0.5	16.3	7.3	17.4	11.2	18.7	17.1
1.0	1.0	15.4	5.5	16.7	9.4	16.0	14.0
1.0	1.5	14.9	4.5	15.9	8.3	15.2	13.4
0.5	0.5	10.7	7.1	11.7	10.1	11.7	10.1
0.5	1.0	8.9	5.5	10.3	8.5	10.3	8.5
0.5	1.5	7.7	4.6	9.3	7.5	9.3	7.5

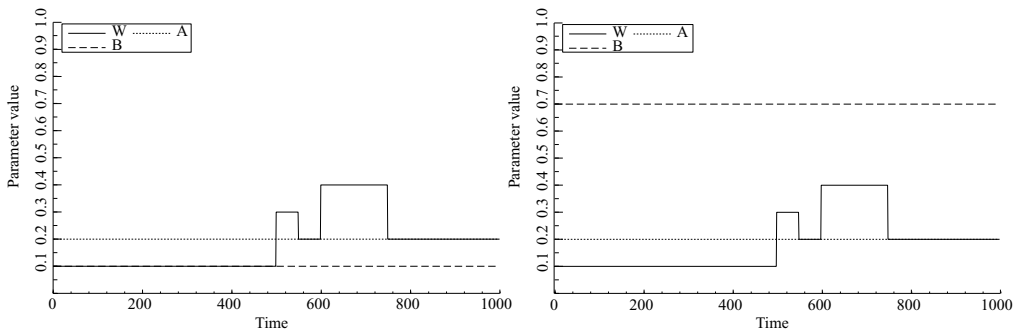


Figure 1. GARCH(1,1) parameters of low (left panel) and high (right panel) GARCH-effect simulations.

5.2.1. *Low GARCH effect.* Let us now discuss simulation results from the low GARCH-effect model. First, we mention the effect of structural changes in time series on the parameter estimation. Later, we compare the performance of all methods in terms of absolute PE.

Estimating a parametric model from data containing a change point will necessarily lead to various biases in estimation. For example, Hillebrand (2005) demonstrates that a change in volatility level ω within a sample drives the GARCH parameter β very close to 1. This is confirmed when we analyse the parameter estimates for parametric and adaptive GARCH at each time point $t \in [250, 1000]$ as depicted on Figure 2, where the mean (solid line), the 10% and 90% quantiles (dotted lines), and the true values (thick dotted line) of the model parameters are provided. The parametric estimates are consistent before breaks starting at $t = 500$, but the GARCH parameter β becomes inconsistent and converges to 1 once data contain breaks, $t > 500$. The locally adaptive estimates are similar to parametric ones before the breaks and become rather imprecise after the first change point, but they are not too far from the true value on average and stay consistent (in the sense that the confidence interval covers the true values). The low precision of estimation can be attributed to rather short intervals used for estimation (cf. Figure 2 for $t < 500$).

Next, we would like to compare the performance of parametric and adaptive estimation methods by means of absolute PE: first for the prediction horizon of one day, $\mathcal{H} = \{1\}$, and later for prediction two weeks ahead, $\mathcal{H} = \{1, \dots, 10\}$. To make the results easier to decipher,

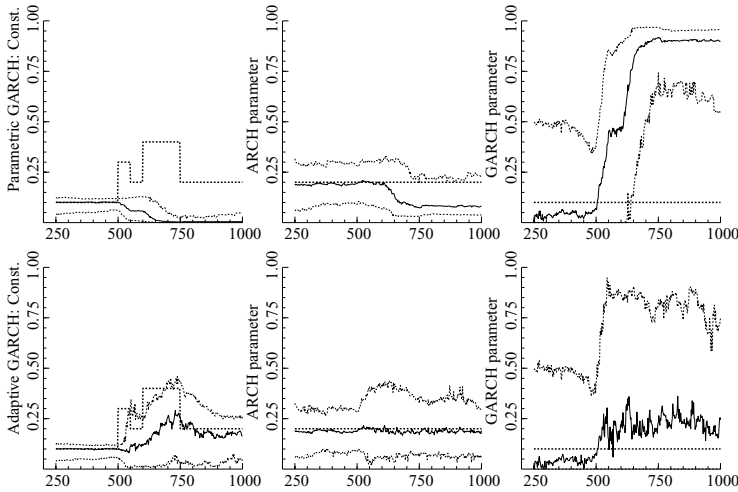


Figure 2. Parameter values estimated by the parametric (top row) and locally adaptive (bottom row) GARCH methods.

we present in what follows PEs averaged over the past month (21 days). The absolute-PE criterion was also used to determine the optimal values of parameters r and ρ (jointly across all simulations and for all $t = 250, \dots, 1000$). The results differ for different models: $r = 0.5$, $\rho = 0.5$ for local constant, $r = 0.5$, $\rho = 1.0$ for local ARCH, and $r = 0.5$, $\rho = 1.5$ for local GARCH.

Let us now compare the adaptively estimated local constant, local ARCH and local GARCH models with the parametric GARCH, which is the best performing parametric model in this set-up. Forecasting one period ahead, the average PEs for all methods and the median lengths of the selected time-homogeneous intervals for adaptive methods are presented on Figure 3 for $t \in [250, 1000]$. First of all, let us observe in the case of the simplest local constant model that even the (median) estimated interval of homogeneity at the end of the first homogeneous period, $1 \leq t < 500$, can actually be shorter than the true one. The reason is that the probability of some 5 or 10 subsequent observations used as I_0 having their sample variance very different from the underlying one increases with the length of the series.

Next, one can notice that all methods are sensitive to jumps in volatility, especially to the first one at $t = 500$: the parametric ones because they ignore a structural break, the adaptive ones because they use a small amount of data after a structural change. In general, the local GARCH performs rather similarly to the parametric GARCH for $t < 650$ because it uses all historical data. After initial volatility jumps, the local GARCH, however, outperforms the parametric one, $650 < t < 775$. Following the last jump at $t = 750$, where the volatility level returns closer to the initial one, the parametric GARCH is best of all methods for some time, $775 < t < 850$, until the adaptive estimation procedure detects the (last) break, and after it, ‘collects’ enough observations for estimation. Then the local GARCH and local ARCH become preferable to the parametric model again, $850 < t$. Interestingly, the local ARCH approximation performs almost as well as both GARCH methods and even outperforms them shortly after structural breaks (except for break at $t = 750$), $600 < t < 775$ and $850 < t < 1000$. Finally, the local constant

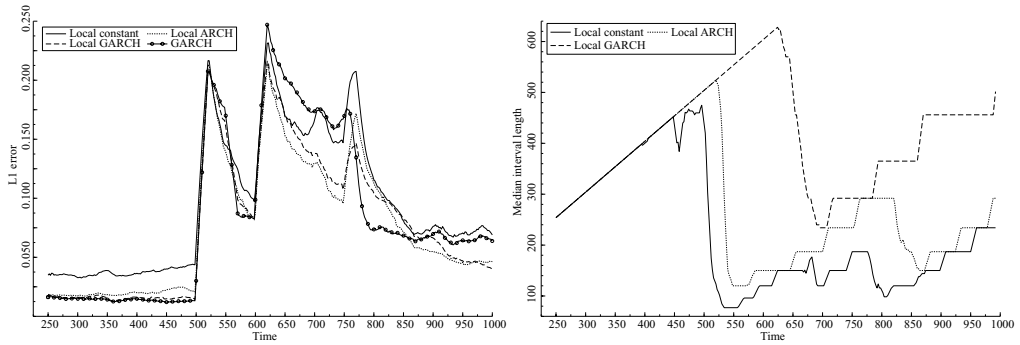


Figure 3. Left-hand panel: Low GARCH-effect simulations—absolute prediction errors one period ahead. Right-hand panel: The median lengths of the adaptively selected intervals.

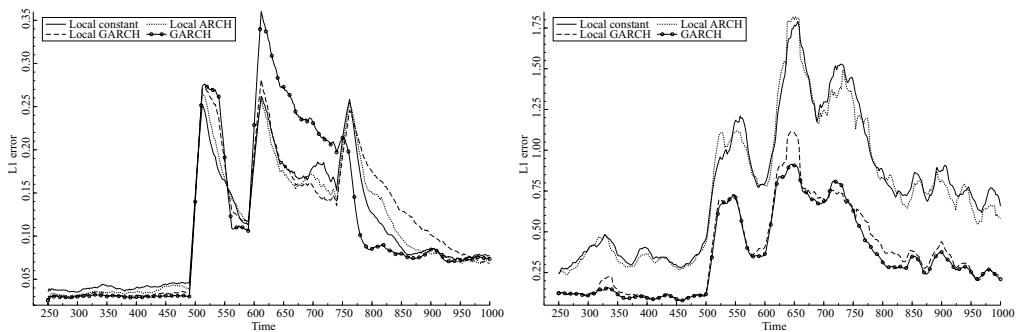


Figure 4. Left-hand panel: Low GARCH-effect simulations—absolute prediction errors 10 periods ahead. Right-hand panel: High GARCH-effect simulations—absolute prediction errors one period ahead.

volatility is lacking behind the other two adaptive methods whenever there is a longer time period without a structural break, but keeps up with them in periods with frequent volatility changes, $500 < t < 650$. All these observations can be documented also by the absolute PE averaged over the whole period $250 \leq t \leq 1000$ (we refer to it as the global PE from now on): the smallest PE is achieved by local ARCH (0.075), then by local GARCH (0.079) and the worst result is from local constant (0.094).

Additionally, all models are compared using the forecasting horizon of 10 days. Most of the results are the same (e.g. parameter estimates) or similar (e.g. absolute PE) to forecasting one period ahead due to the fact that all models rely on at most one past observation. The absolute PEs averaged over one month are summarized for $t \in [250, 1000]$ on Figure 4, which reveals that the difference between local constant volatility, local ARCH and local GARCH models are smaller in this case. As a result, it is interesting to note that: (i) the local constant model becomes a viable alternative to the other methods (it has in fact the smallest global PE 0.107 from all adaptive methods) and (ii) the local ARCH model still outperforms the local GARCH (global

PEs are 0.108 and 0.116, respectively) even though the underlying model is GARCH (with a small value of $\beta = 0.1$ however).

5.2.2. High GARCH effect. Let us now discuss the high GARCH-effect model. One would expect much more prevalent behaviour of both GARCH models, since the underlying GARCH parameter is higher and the changes in the volatility level ω are likely to be small compared to overall volatility fluctuations. Note that the optimal values of tuning constant r and ρ differ from the low GARCH-effect simulations: $r = 0.5$, $\rho = 1.5$ for local constant; $r = 0.5$, $\rho = 1.5$ for local ARCH; and $r = 1.0$, $\rho = 0.5$ for local GARCH.

Comparing the absolute PEs for the one-period-ahead forecast at each time point (Figure 4) indicates that the adaptive and parametric GARCH estimations perform approximately equally well. On the other hand, both the parametric and adaptively estimated ARCH and constant volatility models are lacking significantly. Unreported results confirm, similarly to the low GARCH-effect simulations, that the differences among method are much smaller once a longer prediction horizon of 10 days is used.

6. APPLICATIONS

The proposed adaptive pointwise estimation method will be now applied to real time series consisting of the log-returns of the DAX and S&P 500 stock indices (Sections 6.1 and 6.2). We will again summarize the results concerning both parametric and adaptive methods by the absolute PEs one day ahead averaged over one month. As a benchmark, we employ the parametric GARCH estimated using the last two years of data (500 observations). Since we however do not have the underlying volatility process now, it is approximated by squared returns. Despite being noisy, this approximation is unbiased and provides usually the correct ranking of methods (Andersen and Bollerslev, 1998).

6.1. DAX analysis

Let us now analyse the log-returns of the German stock index DAX from January 1990 till December 2002 depicted at the top of Figure 5. Several periods interesting for comparing the performance of parametric and adaptive pointwise estimates are selected since results for the whole period might be hard to decipher at once.

First, consider the estimation results for years 1991 to 1996. Contrary to later periods, there are structural breaks practically immediately detected by all adaptive methods (July 1991 and June 1992; cf. Stapf and Werner, 2003). For the local GARCH, this differs from less pronounced structural changes discussed later, which are typically detected only with delays of several months. One additional break detected by all methods occurs in October 1994. Note that parameters r and ρ were $r = 0.5$, $\rho = 1.5$ for local constant, $r = 1.0$, $\rho = 1.0$ for local ARCH, and $r = 0.5$, $\rho = 1.5$ for local GARCH.

The results for the period 1991–96 are summarized in the left bottom panel of Figure 5, which depicts the PEs of each adaptive method relative to the PEs of parametric GARCH. First, one can notice that the local constant and local ARCH approximations are preferable till July 1991, where we have less than 500 observations. After the detection of the structural change in June 1991, all adaptive methods are shortly worse than the parametric GARCH due to the limited amount of data used, but then outperform the parametric GARCH till the next structural break in the second half of 1992. A similar behaviour can be observed after the break detected in October 1994,

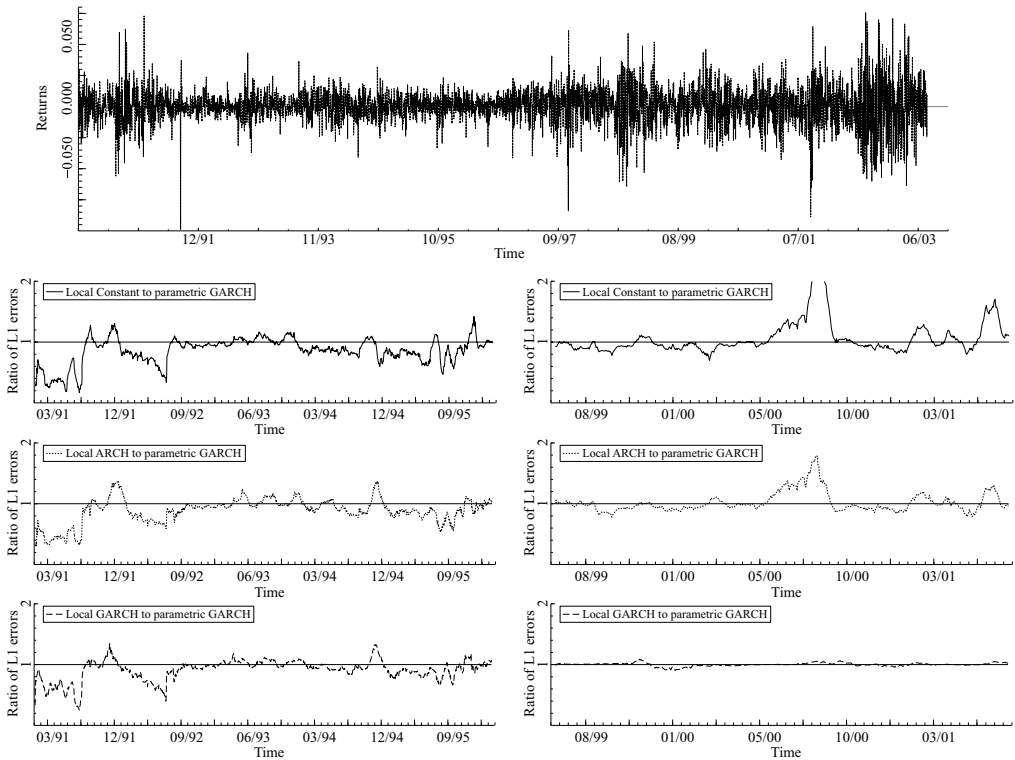


Figure 5. Top panel: The log-returns of DAX series. Bottom panels: The absolute prediction errors of the pointwise adaptive methods relative to the parametric GARCH errors for predictions one period ahead.

where the local constant and local ARCH models actually outperform both the parametric and adaptive GARCH. In the other parts of the data, the performance of all methods is approximately the same, and even though the adaptive GARCH is overall better than the parametric one, the most interesting fact is that the adaptively estimated local constant and local ARCH models perform equally well. In terms of the global PE, the local constant is best (0.829), followed by the local ARCH (0.844) and local GARCH (0.869). This closely corresponds to our findings in simulation study with low GARCH effect in Section 5.2. Note that for other choices of r and ρ , the global PEs are at most 0.835 and 0.851 for the local constant and local ARCH, respectively. This indicates low sensitivity to the choice of these parameters.

Next, we discuss the estimation results for years 1999 to 2001 ($r = 1.0$ for all methods now). After the financial markets were hit by the Asian crisis in 1997 and the Russian crisis in 1998, the market headed to a more stable state in year 1999. The adaptive methods detected the structural breaks in the autumn of 1997 and 1998. The local GARCH detected them, however, with more than a one-year delay—only during 1999. The results in Figure 5 (right bottom panel) confirm that the benefits of the adaptive GARCH are practically negligible compared to the parametric GARCH in such a case. On the other hand, the local constant and ARCH methods perform slightly better than both GARCH methods during the first presented year (July 1999 to June 2000). From July 2000, the situation becomes just the opposite and the performance

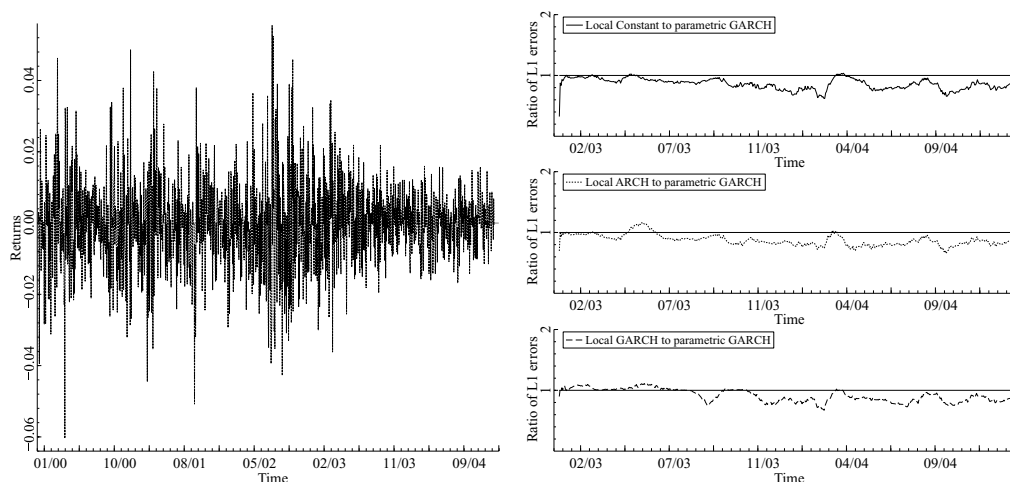


Figure 6. Left-hand panel: The log-returns of S&P 500. Right-hand panel: The absolute prediction errors of the pointwise adaptive methods relative to the parametric GARCH errors for predictions one period ahead.

of the GARCH models is better (parametric and adaptive GARCH estimates are practically the same in this period since the last detected structural change occurred approximately two years ago). Together with previous results, this opens the question of model selection among adaptive procedures as different parametric approximations might be preferred in different time periods. Judging by the global PE, the local ARCH provides slightly better predictions on average than the local constant and local GARCH—despite the ‘peak’ of the PE ratio in the second half of year 2000 (see Figure 5). This, however, depends on the specific choice of loss Λ in (3.6).

Finally, let us mention that the relatively similar behaviour of the local constant and local ARCH methods is probably due to the use of ARCH(1) model, which is not sufficient to capture more complex time developments. Hence, ARCH(p) might be a more appropriate interim step between the local constant and GARCH models.

6.2. S&P 500

Now we turn our attention to more recent data regarding the S&P 500 stock index considered from January 2000 to December 2004; see Figure 6. This period is marked by many substantial events affecting the financial markets, ranging from September 11, 2001, terrorist attacks and the war in Iraq (2003) to the crash of the technology stock-market bubble (2000–02). For the sake of simplicity, a particular time period is again selected: year 2003 representing a more volatile period (the war in Iraq) and year 2004 being a less volatile period. All adaptive methods detected rather quickly a structural break at the beginning of 2003, and additionally they detected a structural break in the second half of 2003, although the adaptive GARCH did so with a delay of more than eight months. The ratios of monthly PE of all adaptive methods to those of the parametric GARCH from January 2003 to December 2004 are summarized on Figure 6 ($r = 0.5$ and $\rho = 1.5$ for all methods).

In the beginning of year 2003, corresponding with 2002 to a more volatile period (see Figure 6), all adaptive methods perform as well as the parametric GARCH. In the middle of year 2003, the local constant and local ARCH models are able to detect another structural change (possibly less pronounced than the one at the beginning of 2003 because of its late detection by the adaptive GARCH). Around this period, the local ARCH shortly performs worse than the parametric GARCH. From the end of 2003 and in year 2004, all adaptive methods starts to outperform the parametric GARCH, where the reduction of the PEs due to the adaptive estimation amounts to 20% on average. All adaptive pointwise estimates exhibit a short period of instability in the first months of 2004, where their performance temporarily worsens to the level of parametric GARCH. This corresponds to ‘uncertainty’ of the adaptive methods about the length of the interval of homogeneity. After this short period, the performance of all adaptive methods is comparable, although the local constant performs overall best of all methods (closely followed by local ARCH) judged by the global PE.

Similarly to the low GARCH-effect simulations and to the analysis of DAX in Section 6.1, it seems that the benefit of pointwise adaptive estimation is most pronounced during periods of stability that follow an unstable period (i.e. year 2004) rather than during a presumably rapidly changing environment. The reason is that, despite possible inconsistency of parametric methods under change points, the adaptive methods tend to have a rather large variance when the intervals of time homogeneity become very short.

7. CONCLUSION

We extend the idea of adaptive pointwise estimation to parametric CH models. In the specific case of ARCH and GARCH, which represent particularly difficult cases due to high data demands and dependence of critical values on underlying parameters, we demonstrate the use and feasibility of the proposed procedure: on the one hand, the adaptive procedure, which itself depends on a number of auxiliary parameters, is shown to be rather insensitive to their choice, and on the other hand, it facilitates the global selection of these parameters by means of fit or forecasting criteria. The real-data applications highlight the flexibility of the proposed time-inhomogeneous models since even simple varying-coefficients models such as constant volatility and ARCH(1) can outperform standard parametric methods such as GARCH(1,1). Finally, the relatively small differences among the adaptive estimates based on different parametric approximations indicate that, in the context of adaptive pointwise estimation, it is sufficient to concentrate on simpler and less data-intensive models such as ARCH(p), $0 \leq p \leq 3$, to achieve good forecasts.

ACKNOWLEDGMENTS

This research was supported by the Deutsche Forschungsgemeinschaft through the SFB 649 ‘Economic Risk’.

REFERENCES

Andersen, T. G. and T. Bollerslev (1998). Answering the skeptics: yes, standard volatility models do provide accurate forecasts. *International Economic Review* 39, 885–905.

- Andreou, E. and E. Ghysels (2002). Detecting multiple breaks in financial market volatility dynamics. *Journal of Applied Econometrics* 17, 579–600.
- Andreou, E. and E. Ghysels (2006). Monitoring disruptions in financial markets. *Journal of Econometrics* 135, 77–124.
- Andrews, D. W. K. (1993). Tests for parameter instability and structural change with unknown change point. *Econometrica* 61, 821–56.
- Bai, J. and P. Perron (1998). Estimating and testing linear models with multiple structural changes. *Econometrica* 66, 47–78.
- Beltratti, A. and C. Morana (2004). Structural change and long-range dependence in volatility of exchange rates: either, neither or both? *Journal of Empirical Finance* 11, 629–58.
- Bollerslev, T. (1986). Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics* 31, 307–27.
- Cai, Z., J. Fan and Q. Yao (2000). Functional coefficient regression models for nonlinear time series. *Journal of the American Statistical Association* 95, 941–56.
- Chen, J. and A. K. Gupta (1997). Testing and locating variance changepoints with application to stock prices. *Journal of the American Statistical Association* 92, 739–47.
- Chen, R. and R. J. Tsay (1993). Functional-coefficient autoregressive models. *Journal of the American Statistical Association* 88, 298–308.
- Cheng, M.-Y., J. Fan and V. Spokoiny (2003). Dynamic nonparametric filtering with application to volatility estimation. In M. G. Akritas and D. N. Politis (Eds.), *Recent Advances and Trends in Nonparametric Statistics*, 315–33. Amsterdam: Elsevier.
- Diebold, F. X. and A. Inoue (2001). Long memory and regime switching. *Journal of Econometrics* 105, 131–59.
- Doornik, J. A. (2002). Object-oriented programming in econometrics and statistics using Ox: a comparison with C++, Java and C#. In S. S. Nielsen (Ed.), *Programming Languages and Systems in Computational Economics and Finance*, 115–47. Dordrecht: Kluwer.
- Engle, R. F. (1982). Autoregressive conditional heteroscedasticity with estimates of the variance of United Kingdom inflation. *Econometrica* 50, 987–1008.
- Fan, J. and W. Zhang (2008). Statistical models with varying coefficient models. *Statistics and Its Interface* 1, 179–95.
- Franco, C. and J.-M. Zakoian (2007). Quasi-maximum likelihood estimation in GARCH processes when some coefficients are equal to zero. *Stochastic Processes and their Applications* 117, 1265–84.
- Glosten, L. R., R. Jagannathan and D. E. Runkle (1993). On the relation between the expected value and the volatility of the nominal excess return on stocks. *Journal of Finance* 48, 1779–801.
- Hansen, B. and S.-W. Lee (1994). Asymptotic theory for the GARCH(1,1) quasi-maximum likelihood estimator. *Econometric Theory* 10, 29–53.
- Härdle, W., H. Herwatz and V. Spokoiny (2003). Time inhomogeneous multiple volatility modelling. *Journal of Financial Econometrics* 1, 55–99.
- Herwatz, H. and H. E. Reimers (2001). Empirical modeling of the DEM/USD and DEM/JPY foreign exchange rate: structural shifts in GARCH-models and their implications. 2001–83, Discussion Paper SFB 373, Humboldt-Universität zu Berlin, Germany.
- Hillebrand, E. (2005). Neglecting parameter changes in GARCH models. *Journal of Econometrics* 129, 121–38.
- Kokoszka, P. and R. Leipus (2000). Change-point estimation in ARCH models. *Bernoulli* 6, 513–39.
- Laurent, S. and J.-P. Peters (2006). *G@RCH 4.2, Estimating and Forecasting ARCH Models*. London: Timberlake Consultants Press.

Mercurio, D. and V. Spokoiny (2004). Statistical inference for time-inhomogeneous volatility models. *Annals of Statistics* 32, 577–602.

Mikosch, T. and C. Starica (1999). Change of structure in financial time series, long range dependence and the GARCH model. Working Paper, Department of Statistics, University of Pennsylvania. See <http://citeseer.ist.psu.edu/mikosch99change.html>.

Mikosch, T. and C. Starica (2004). Changes of structure in financial time series and the GARCH model. *Revstat Statistical Journal* 2, 41–73.

Nelson, D. B. (1991). Conditional heteroskedasticity in asset returns: a new approach. *Econometrica* 59, 347–70.

Pesaran, M. H. and A. Timmermann (2004). How costly is it to ignore breaks when forecasting the direction of a time series? *International Journal of Forecasting* 20, 411–25.

Sentana, E. (1995). Quadratic ARCH models. *Review of Economic Studies* 62, 639–61.

Spokoiny, V. (1998). Estimation of a function with discontinuities via local polynomial fit with an adaptive window choice. *Annals of Statistics* 26, 1356–78.

Spokoiny, V. (2009a). Multiscale local change-point detection with applications to value-at-risk. *Annals of Statistics* 37, 1405–36.

Spokoiny, V. (2009b). Parameter estimation in time series analysis. WIAS Preprint No. 1404, Weierstrass Institute for Applied Analysis and Stochastics, Berlin, Germany.

Stapf, J. and T. Werner (2003). How wacky is DAX? The changing structure of German stock market volatility. Discussion Paper 2003/18, Deutsche Bundesbank, Germany.

Taylor, S. J. (1986). *Modeling Financial Time Series*. Chichester: Wiley.

APPENDIX: PROOFS

Proof of Corollary 2.1: Given the choice of \mathfrak{z}_α , it directly follows from (2.5). □

Proof of Theorem 3.1: Consider the event $\mathcal{B}_k = \{\hat{I} = I_{k-1}\}$ for some $k \leq K$. This particularly means that I_{k-1} is accepted while $I_k = [T - m_k + 1, T]$ is rejected; i.e. there is $I' = [t', T] \subseteq I_k$ and $\tau \in \mathcal{I}(I_k)$ such that $T_{I_k, \tau} > \mathfrak{z}_k = \mathfrak{z}_{I_k, \mathcal{I}(I_k)}$. For every fixed $\tau \in \mathcal{I}(I_k)$ and $J = I_k \setminus [\tau + 1, T]$, $J^c = [\tau + 1, T]$, it holds by definition of $T_{I_k, \tau}$ that

$$T_{I_k, \tau} \leq L_J(\tilde{\theta}_J) + L_{J^c}(\tilde{\theta}_{J^c}) - L_I(\theta_0) = L_J(\tilde{\theta}_J, \theta_0) + L_{J^c}(\tilde{\theta}_{J^c}, \theta_0).$$

This implies by Theorem 2.1 that $\mathbf{P}_{\theta_0}(T_{I_k, \tau} > 2\mathfrak{z}) \leq \exp\{\epsilon(\lambda, \theta_0) - \lambda\mathfrak{z}\}$. Now,

$$\mathbf{P}_{\theta_0}(\mathcal{B}_k) \leq \sum_{t'=T-m_k+1}^{T-m_0} \sum_{\tau=t'+1}^{T-m_0+1} 2 \exp\{\epsilon(\lambda, \theta_0) - \lambda\mathfrak{z}_k/2\} \leq 2 \frac{m_k^2}{2} \exp\{\epsilon(\lambda, \theta_0) - \lambda\mathfrak{z}_k/2\}.$$

Next, by the Cauchy–Schwartz inequality

$$\begin{aligned} \mathbf{E}_{\theta_0} |L_{I_K}(\tilde{\theta}_{I_K}, \hat{\theta})|^r &= \sum_{k=1}^K \mathbf{E}_{\theta_0} [|L_{I_K}(\tilde{\theta}_{I_K}, \tilde{\theta}_{k-1})|^r \mathbf{1}(\mathcal{B}_k)] \\ &\leq \sum_{k=1}^K \mathbf{E}_{\theta_0}^{1/2} |L_{I_K}(\tilde{\theta}_{I_K}, \tilde{\theta}_{k-1})|^{2r} \mathbf{P}_{\theta_0}^{1/2}(\mathcal{B}_k). \end{aligned}$$

Under the conditions of Theorem 2.1, it follows similarly to (2.6) that

$$\mathbf{E}_{\theta_0} |L_{I_K}(\tilde{\theta}_{I_K}, \tilde{\theta}_{k-1})|^{2r} \leq (m_K/m_{k-1})^{2r} \mathfrak{R}_{2r}^*(\theta_0)$$

for some constant $\mathfrak{R}_{2r}^*(\theta_0)$ and $k = 1, \dots, K$, and therefore,

$$\mathbf{E}_{\theta_0} |L_{I_K}(\tilde{\theta}_{I_K}, \hat{\theta})|^r \leq [\mathfrak{R}_{2r}^*(\theta_0)]^{1/2} \sum_{k=1}^K m_k (m_K/m_{k-1})^r \exp\{\epsilon(\lambda, \theta_0)/2 - \lambda \mathfrak{z}_k/4\}$$

and the result follows by simple algebra provided that $a_1\lambda/4 \geq 1$ and $a_2\lambda/4 > 2$. □

LEMMA A.1. *Let \mathbf{P} and \mathbf{P}_0 be two measures such that the Kullback–Leibler divergence $\mathbf{E} \log(d\mathbf{P}/d\mathbf{P}_0)$, satisfies $\mathbf{E} \log(d\mathbf{P}/d\mathbf{P}_0) \leq \Delta < \infty$. Then for any random variable ζ with $\mathbf{E}_0\zeta < \infty$, it holds that $\mathbf{E} \log(1 + \zeta) \leq \Delta + \mathbf{E}_0\zeta$.*

Proof: By simple algebra one can check that for any fixed y the maximum of the function $f(x) = xy - x \log x + x$ is attained at $x = e^y$ leading to the inequality $xy \leq x \log x - x + e^y$. Using this inequality and the representation $\mathbf{E} \log(1 + \zeta) = \mathbf{E}_0\{Z \log(1 + \zeta)\}$ with $Z = d\mathbf{P}/d\mathbf{P}_0$ we obtain

$$\begin{aligned} \mathbf{E} \log(1 + \zeta) &= \mathbf{E}_0\{Z \log(1 + \zeta)\} \leq \mathbf{E}_0(Z \log Z - Z) + \mathbf{E}_0(1 + \zeta) \\ &= \mathbf{E}_0(Z \log Z) + \mathbf{E}_0\zeta - \mathbf{E}_0Z + 1. \end{aligned}$$

It remains to note that $\mathbf{E}_0Z = 1$ and $\mathbf{E}_0(Z \log Z) = \mathbf{E} \log Z$. □

Proof of Theorem 4.1: Lemma A.1 applied with $\zeta = \varrho(\hat{\theta}, \theta)/\mathbf{E}_\theta\varrho(\hat{\theta}, \theta)$ yields the result in the view of

$$\begin{aligned} \mathbf{E}_\theta(Z_{I,\theta} \log Z_{I,\theta}) &= \mathbf{E} \log Z_{I,\theta} = \mathbf{E} \sum_{t \in I} \log \frac{p[Y_t, g(X_t)]}{p[Y_t, g(X_t(\theta))]} \\ &= \mathbf{E} \sum_{t \in I} \mathbf{E} \left\{ \log \frac{p[Y_t, g(X_t)]}{p[Y_t, g(X_t(\theta))]} \middle| \mathcal{F}_{t-1} \right\} = \mathbf{E} \Delta_{I_k}(\theta). \end{aligned} \quad \square$$

Proof of Corollary 4.1: It is Theorem 4.1 formulated for $\varrho(\theta', \theta) = L_I(\theta', \theta)$. □

Proof of Theorem 4.2: The first inequality follows from Corollary 4.1, the second one from condition (3.4) and the property $x \geq \log x$ for $x > 0$. □

Proof of Theorem 4.3: Let $\hat{k} = k > k^*$. This means that I_k is not rejected as homogeneous. Next, we show that for every $k > k^*$ the inequality $T_{I_k, \tau} \leq T_{I_k, \mathcal{S}(I_k)} \leq \mathfrak{z}_k$ with $\tau = T - m_{k^*} = T - |I_{k^*}|$ implies $L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \tilde{\theta}_{I_k}) \leq \mathfrak{z}_{k^*}$. Indeed with $J = I_k \setminus I_{k^*}$, this means that, by construction, $\mathfrak{z}_k \leq \mathfrak{z}_{k^*}$ for $k > k^*$ and

$$\mathfrak{z}_k \geq T_{I_k, \tau} = L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \tilde{\theta}_{I_k}) + L_J(\tilde{\theta}_J, \tilde{\theta}_{I_k}) \geq L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \tilde{\theta}_{I_k}).$$

It remains to note that

$$|L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \hat{\theta})|^r \leq |L_{I_{k^*}}(\tilde{\theta}_{I_{k^*}}, \hat{\theta}_{I_{k^*}})|^r \mathbf{1}(\hat{k} < k^*) + \mathfrak{z}_{k^*}^r \mathbf{1}(\hat{k} > k^*),$$

which obviously yields the assertion. □

Dynamic semiparametric factor models in risk neutral density estimation

Enzo Giacomini · Wolfgang Härdle ·
Volker Krättschmer

Received: 1 March 2009 / Accepted: 31 August 2009 / Published online: 18 September 2009
© Springer-Verlag 2009

Abstract Dynamic semiparametric factor models (DSFM) simultaneously smooth in space and are parametric in time, approximating complex dynamic structures by time invariant basis functions and low dimensional time series. In contrast to traditional dimension reduction techniques, DSFM allows the access of the dynamics embedded in high dimensional data through the lower dimensional time series. In this paper, we study the time behavior of risk assessments from investors facing random financial payoffs. We use DSFM to estimate risk neutral densities from a dataset of option prices on the German stock index DAX. The dynamics and term structure of risk neutral densities are investigated by Vector Autoregressive (VAR) methods applied on the estimated lower dimensional time series.

Keywords Dynamic factor models · Dimension reduction · Risk neutral density

1 Introduction

Large datasets containing various samples of high dimensional observations became common in diverse fields of science with advances in measurement and computational techniques. In many applications the data come in curves, i.e., as observations of discretized values of smooth random functions, presenting evident functional structure. In these cases, it is natural to perform statistical inference using functional data analysis techniques.

E. Giacomini (✉) · W. Härdle · V. Krättschmer
CASE—Center for Applied Statistics and Economics, Humboldt-Universität zu Berlin,
Spandauerstr. 1, 10178 Berlin, Germany
e-mail: enzogiacomini@gmail.com

V. Krättschmer
Institute of Mathematics, Technische Universität Berlin, Straße des 17. Juni 136, 10623 Berlin,
Germany

Consider a dataset $\{(Y_{jt}, X_{jt})\}$, $j = 1, \dots, J_t$, $t = 1, \dots, T$, containing noisy samples of a real valued smooth random function $\mathcal{F} \in L_2(\mathcal{X})$, $\mathcal{X} \subseteq \mathbb{R}^d$, $d \in \mathbb{N}$, evaluated at unbalanced design points as

$$Y_{jt} = \mathcal{F}_t(X_{jt}) + \varepsilon_{jt}, \quad (1.1)$$

where ε_{jt} denote unknown zero-mean error terms and $\{\mathcal{F}_t\}$ are realizations of \mathcal{F} . Each sample $S_t = \{(Y_{jt}, X_{jt}) : j = 1, \dots, J_t\}$, $t = 1, \dots, T$, may correspond to observations on, e.g., different individuals, time periods or experimental conditions. Examples in biomedicine are measurements of growth curves and brain potentials across individuals, see Kneip and Gasser (1992) and Gasser and Kneip (1995), in econometrics such are expenditures across households and implied volatilities across trading days, see Kneip (1994) and Fengler et al. (2007).

A large branch of functional data analysis concentrates on approximating \mathcal{F} by lower dimensional objects. Distributions on function spaces are highly complex objects and dimension reduction techniques present a feasible and interpretable approach for investigating them. Functional principal components analysis (FPCA), based on the Karhunen–Loève expansion of \mathcal{F} is the most prominent and widely used dimension reduction technique, see Rice and Silverman (1991) and Ramsay and Dalzell (1991).

Asymptotic results on FPCA have been obtained by Dauxois et al. (1982) and Hall et al. (2006) for observed functional data $\{\mathcal{F}_t\}$. For non-observable data, the standard approach is to perform FPCA on presmoothed $\{\widehat{\mathcal{F}}_t\}$, see Benko et al. (2009) for recent developments. In practical applications, however, presmoothing may suffer from design-sparseness, see Cont and Fonseca (2002) and Fengler et al. (2007).

In general lines, previous literature combines PCA and dimension reduction with presmoothing for effective dimensional space at fixed time horizon. Various applications, however, involve the dynamics of the unobserved random functions, calling for dimension reduction techniques that smooth in space and are parametric in time.

In this paper, we investigate the dynamics of $\{\mathcal{F}_t\}$ by reducing dimensionality without presmoothing. \mathcal{F}_t is considered as a linear combination of $L + 1 \ll T$ unknown smooth basis functions $m_l \in L_2(\mathcal{X})$, $l = 0, \dots, L$:

$$\mathcal{F}_t(X_{jt}) = \sum_{l=0}^L Z_{lt} m_l(X_{jt}), \quad (1.2)$$

where $Z_t = (Z_{0t}, \dots, Z_{Lt})^\top$ is an unobservable random vector taking values on \mathbb{R}^{L+1} with $Z_{0t} = 1$. Defining the tuple of functions $m = (m_0, \dots, m_L)^\top$, the Dynamic Semiparametric Factor Model (DSFM) reads as

$$Y_{jt} = Z_t^\top m(X_{jt}) + \varepsilon_{jt}. \quad (1.3)$$

The basis functions are estimated nonparametrically avoiding specification issues. Their estimation is performed simultaneously with Z_t , i.e., the smoothing is transferred directly to m_l and design-sparseness issues become secondary. In addition, the random process $\{Z_t\}$ is allowed to be non-stationary. Park et al. (2009) show that

under (1.2) the autocorrelation structures of $\{\widehat{Z}_t\}$ and $\{Z_t\}$ are asymptotically equivalent; therefore, no loss is incurred by inferring the dynamics from the estimated $\{\widehat{Z}_t\}$, and there is no payment for not knowing the true $\{Z_t\}$. This result is essential for investigating cointegration between dynamical systems, see Brüggemann et al. (2008) for an econometric application.

Note that the common regressors model, Kneip (1994), also represents unobservable functions by (1.2). There are, however, crucial differences between the DSFM and common regressors:

1. In DSFM, $\{Z_t\}$ is a (non-stationary) random process with autocovariance structure inferable from $\{\widehat{Z}_t\}$.
2. DSFM is implementable in unbalanced designs.
3. DSFM avoids presmoothing by transferring the smoothing to the basis functions.

Thus DSFM goes beyond traditional dimension reductions techniques (FPCA and common regressors) as it captures structural dynamics embedded in the observations.

In economics, there is substantial interest in the behavior (over time) of investors facing risks and its relation to macroeconomic and financial indicators. The knowledge about the dynamics of risk assessments from investors is essential for many applications ranging from pricing of illiquid instruments to risk management.

Option prices contain information on risk assessments from investors facing future financial payoffs, summarized in the risk neutral densities q , see Ait-Sahalia and Lo (1998). An European call option with price C_t at time $0 \leq t \leq T$, maturity date $T > 0$ and strike $K > 0$ is a financial instrument that delivers the random payoff $(S_T - K)^+$ at time T where S_t is the price of an underlying asset at time t . Breeden and Litzenberger (1978) show that under no arbitrage assumptions the risk neutral density is obtained from the European call price function C_t through the relation

$$q_{t,T}(s_T | s_t) = e^{r(T-t)} \frac{\partial^2 C_t(s_t, r, K, T-t)}{\partial K^2} \Big|_{K=s_T}, \tag{1.4}$$

where $r > 0$ is interest rate, see Sect. 4 for details.

We estimate risk neutral densities based on observed intraday prices of calls on the German stock index (DAX). Each observation consists of a price Y_{jt} on a design point $X_{jt} = (\kappa_{jt}, \tau_{jt})^\top$ where $j = 1, \dots, J_t$, denote the transactions at day $t = 1, \dots, T$, κ is the moneyness, a monotone transformation of strikes K , and $\tau = T - t$ is the time to maturity associated with the option. Stock exchange regulations impose prespecified values for tradable maturities resulting in degenerated designs, see Fig. 1.

Following Ait-Sahalia and Lo (1998) and Fengler et al. (2007), call prices are transformed into log-implied volatilities $\widetilde{Y}_{jt} = \log C_{BS}^{-1}(Y_{jt})$, where C_{BS} is the Black–Scholes call price function defined in Sect. 4. These are assumed as discretized noisy values of the log-implied volatility surface evaluated at $\{X_{jt}\}$:

$$\widetilde{Y}_{jt} = \log \mathcal{V}_t(X_{jt}) + \varepsilon_{jt}, \tag{1.5}$$

where $\mathcal{V} \in L_2(\mathcal{X})$, $\mathcal{X} \subset \mathbb{R}_+^2$, is a smooth random function, called the implied volatility surface, and ε_{jt} is an error term. The realizations $\{\mathcal{V}_t\}$ are filtered out from the data with DSFM and, remarking that C_{BS} is a function of K , the risk neutral densities are

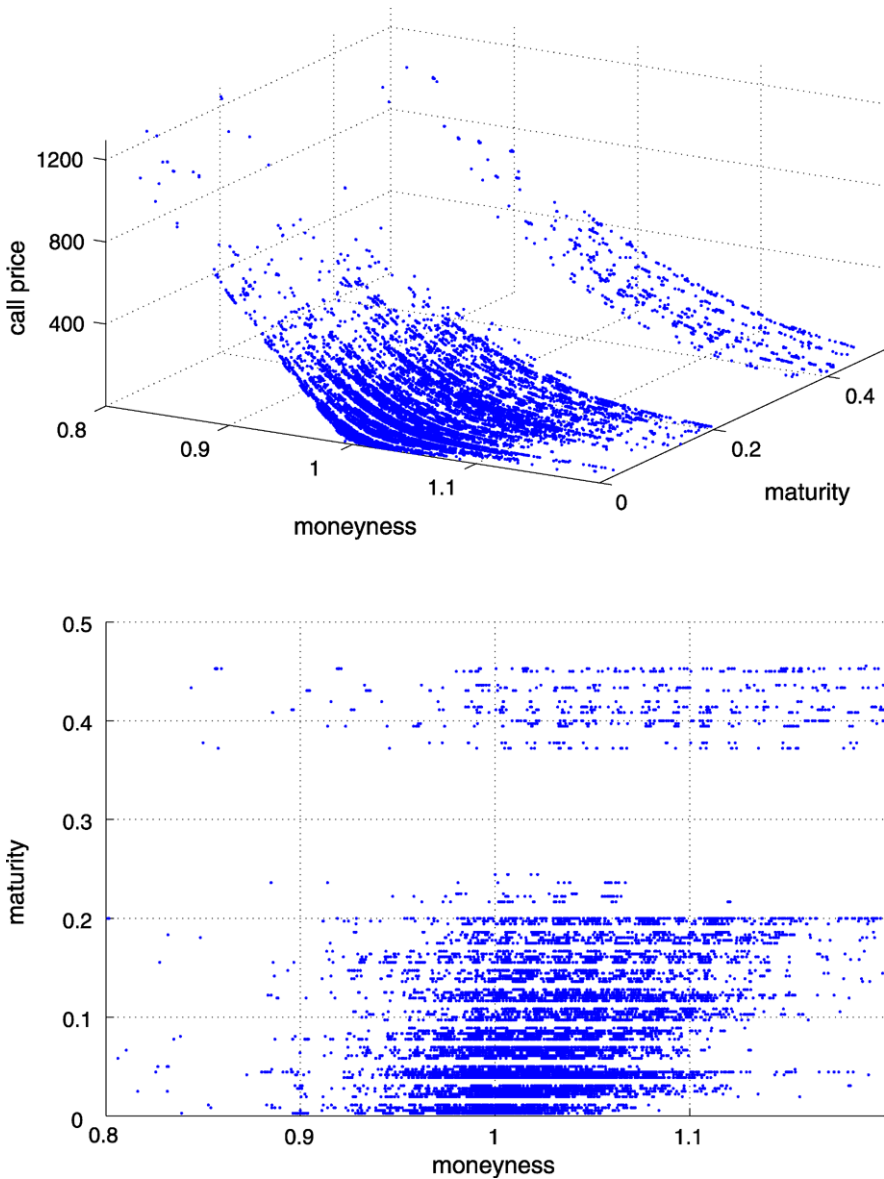


Fig. 1 Samples S_t , $t = 1, \dots, 22$, of DAX call prices traded on January 2001 (left). Corresponding unbalanced design $\{X_{jt}\}$ (right)

obtained by (1.4) with $C_{BS}(\widehat{V})$ as an estimator for C_t . The dynamics of the estimated $\{\widehat{q}_{t,T}\}$ is analyzed based on the autocorrelation structure of $\{\widehat{Z}_t\}$.

In the sequel, the DSFM estimation method and its asymptotic properties are described (Sect. 2). In Sect. 3, the risk neutral densities are defined, and in Sect. 4 they are estimated from observed prices of European call options on the DAX index

(ODAX dataset). Their dynamic structure is then analyzed by vector autoregressive models.

2 Estimation method

Consider a dataset $\{(Y_{jt}, X_{jt})\}$, $j = 1, \dots, J_t, t = 1, \dots, T$, such that

$$Y_{jt} = \sum_{l=0}^L Z_{lt} m_l(X_{jt}) + \varepsilon_{jt}, \tag{2.1}$$

where ε_{jt} are unknown error terms with $E[\varepsilon_{jt}] = 0$ and $E[\varepsilon_{jt}^2] < \infty$. The variables $X_{11}, \dots, X_{T,J_T}, \varepsilon_{1,1}, \dots, \varepsilon_{T,J_T}$ are independent. Here $Z_t = (Z_{0t}, \dots, Z_{Lt})^\top$ is an unobservable random vector taking values on \mathbb{R}^{L+1} with $Z_{0t} = 1$ and $m_l \in L_2(\mathcal{X})$, $l = 0, \dots, L$, are unknown smooth functions, called basis functions, mapping $\mathcal{X} \subseteq \mathbb{R}^d, d \in \mathbb{N}$, into real values.

Following Park et al. (2009), the basis functions are estimated using a series expansion. Defining K normed functions $\psi_k : \mathcal{X} \rightarrow \mathbb{R}, \int_{\mathcal{X}} \psi_k^2(x) dx = 1, k = 1, \dots, K$, and an $((L + 1) \times K)$ matrix of coefficients $\Gamma = (\gamma_{l,k}), \gamma_{l,k} \in \mathbb{R}$, the tuple of functions $m = (m_0, \dots, m_L)^\top$ is approximated by $\Gamma^\top \psi$ where $\psi = (\psi_1, \dots, \psi_K)^\top$. For simplicity of notation, we assume that $J_t = J$ does not depend on t . We define the least squares estimators as

$$(\widehat{\Gamma}, \widehat{Z}) = \arg \min_{\Gamma \in \mathcal{G}, Z \in \mathcal{Z}} \sum_{t=1}^T \sum_{j=1}^J \{Y_{jt} - Z_t^\top \Gamma \psi(X_{jt})\}^2, \tag{2.2}$$

where $\mathcal{G} = \mathcal{M}(L + 1, K), \mathcal{Z} = \{Z \in \mathcal{M}(T, L + 1) : Z_{0t} = 1\}$ and $\mathcal{M}(a, b)$ is the set of all $(a \times b)$ matrices. The basis functions m are estimated by $\widehat{m} = \widehat{\Gamma}^\top \psi$.

Theorem (2.1) gives the asymptotic behavior of the least squares estimators $(\widehat{\Gamma}, \widehat{Z})$. See Park et al. (2009) for the proof.

Theorem 2.1 *Suppose that DSFM holds and that $(\widehat{\Gamma}, \widehat{Z})$ is defined by (2.2). Under Assumptions (A1)–(A8), see Appendix, it holds for $K, J \rightarrow \infty$:*

$$\frac{1}{T} \sum_{1 \leq t \leq T} \|\widehat{Z}_t^\top \widehat{\Gamma} - Z_t^\top \Gamma^*\|^2 = \mathcal{O}_P(\delta_K^2 + \xi^2).$$

See (A5) and (A8) for the definitions of δ_K and ξ . Note that the model (2.1) is only identifiable up to linear transformations. Consider an $((L + 1) \times (L + 1))$ regular matrix $B = (b_{ij})$ with $b_{1j} = \delta_{1j}$ and $b_{i1} = \delta_{i1}$ for $i, j = 1, \dots, L + 1$, where $\delta_{ij} = 1(i = j)$. Define $Z_t^* = B^\top Z_t, m^* = B^{-1}m$. Then from (1.2)

$$\mathcal{F}_t(X) = Z_t^\top m(X) = Z_t^\top B B^{-1} m(X) = Z_t^{*\top} m^*(X)$$

for $X \in \mathcal{X}$. On the other hand, it is always possible to chose orthonormal basis functions by setting $m^* = Hm$ where H is an orthogonal matrix.

Theorem (2.2) states that for any \widehat{Z}_t there exists a random matrix B such that the autocovariances of $\{\widetilde{Z}_t\}$, $\widetilde{Z}_t = B^\top \widehat{Z}_t$, are asymptotically equivalent to the autocovariances of the true unobservable $\{Z_t\}$. This equivalence is transferred to classical estimation and testing procedures in the context of, e.g., vector autoregressive models and, in particular, justifies inference based on $\{\widetilde{Z}_t\}$ when $\{Z_t\}$ is a VAR process. Define for $H_t \in \mathcal{Z}$, $t = 1, \dots, T$: $\overline{H} = T^{-1} \sum_{t=1}^T H_t$, $H_{c,t} = H_t - \overline{H}$ and $H_{n,t} = (T^{-1} \sum_{s=1}^T H_{c,s} H_{c,s}^\top)^{-1/2} H_{c,t}$.

Theorem 2.2 *Suppose that DSFM holds and that $(\widehat{\Gamma}, \widehat{Z})$ is defined by (2.2). Under Assumptions (A1)–(A11), see Appendix, there exists a random matrix B such that for $h \neq 0$, $h_d = \max(1, 1 - h)$, $h_u = \max(T, T - h)$ and $T \rightarrow \infty$:*

$$\frac{1}{T} \sum_{t=h_d}^{h_u} \widetilde{Z}_{c,t} (\widetilde{Z}_{c,t+h} - \widetilde{Z}_{c,t})^\top - \frac{1}{T} \sum_{t=h_d}^{h_u} Z_{c,t} (Z_{c,t+h} - Z_{c,t})^\top = \mathcal{O}_P(T^{-1/2}),$$

where $\widetilde{Z}_t = B^\top \widehat{Z}_t$. Moreover,

$$\frac{1}{T} \sum_{t=h_d}^{h_u} \widetilde{Z}_{n,t} \widetilde{Z}_{n,t+h}^\top - \frac{1}{T} \sum_{t=h_d}^{h_u} Z_{n,t} Z_{n,t+h}^\top = \mathcal{O}_P(T^{-1/2}).$$

See Park et al. (2009) for the proof. Note that, in contrast to FPCA, DSFM does not require stationarity neither for $\{Z_t\}$ nor for $\{\varepsilon_t\}$, but only weak assumptions on the average behavior of Z_t , like being a martingale difference, see Appendix.

3 Risk neutral density estimation

3.1 Risk neutral densities

Consider a financial market with one risky asset and one riskless bond with constant interest rate $r > 0$. Let the price of the asset traded on the market be described by the real valued random process $\{S_t\}$, $t = [0, T]$, $T < \infty$, on a filtered probability space $(\Omega, \{\mathcal{F}_t\}, \mathbb{P})$ with $\mathcal{F}_t = \sigma(S_u, u \leq t)$ and $\mathcal{F}_0 = \{\emptyset, \Omega\}$. Assume further no arbitrage in the financial market in the sense that there exists a (risk neutral) probability measure \mathbb{Q} equivalent to \mathbb{P} under which the discounted price process $\{e^{-rt} S_t\}$ is a martingale.

A European call option at strike $K > 0$ is a financial instrument that pays $\Psi(S_T) = (S_T - K)^+$ at time T . By the risk-neutral valuation principle w.r.t. \mathbb{Q} , the price C_t of a European call option at time t is defined to be

$$C_t = e^{-r(T-t)} E^{\mathbb{Q}}[\Psi(S_T) | \mathcal{F}_t]. \tag{3.1}$$

Assuming that $\{S_t\}$ is a \mathbb{Q} -Markov process and denoting the \mathbb{P} -density of \mathbb{Q} by π , the price can be rewritten as

$$C_t = e^{-r(T-t)} E[\Psi(S_T) \mathcal{K}_\pi^t(S_t, S_T) | S_t],$$

where E denotes the expectation under \mathbb{P} and $\mathcal{K}_\pi^t(S_t, S_T) \stackrel{\text{def.}}{=} \frac{E[\pi|S_t, S_T]}{E[\pi|S_t]}$. The conditional risk neutral distribution of S_T is defined as

$$Q_{S_T|S_t=s_t}([S_T \leq x]) \stackrel{\text{def.}}{=} \int_{-\infty}^x \mathcal{K}_\pi^t(s_t, \cdot) dP_{S_T|S_t=s_t}, \tag{3.2}$$

where $P_{S_T|S_t=s_t}$ is the conditional distribution of S_T under $S_t = s_t$. Specializing to the following two factor model, we assume that the price process has dynamics given by

$$dS_t = S_t \mu(Y_t) dt + S_t \sigma(Y_t) dW_t^1,$$

here W^1 is a standard \mathbb{P} -Brownian motion and Y denotes an external economic factor process modeled by

$$dY_t = g(Y_t) + \rho dW_t^1 + \bar{\rho} dW_t^2,$$

where $\rho \in [-1, 1]$ is some correlation factor, $\bar{\rho} \stackrel{\text{def.}}{=} \sqrt{1 - \rho^2}$ and W^2 is a standard \mathbb{P} -Brownian motion independent of W^1 under \mathbb{P} . Market models of this type are popular in mathematical finance and economics, in particular, if Y follows an Ornstein-Uhlenbeck dynamics with mean reversion term $g(y) = \iota(\theta - y)$ for constants $\theta \geq 0$ and $\iota > 0$. Moreover, $\{S_t\}$ is a \mathbb{Q} -Markov process for any \mathbb{Q} , see Hernández-Hernández and Schied (2007) and the conditional risk neutral distribution $Q_{S_T|S_t=s_t}$ has a density function denoted by $q_{t,T}(\cdot|s_t)$. Hence, recalling (3.1), the call prices can be expressed as

$$C_t(s_t, r, K, T - t) = e^{-r(T-t)} \int (s_T - K)^+ q_{t,T}(s_T|s_t) ds_T.$$

We assume that the observed prices in the financial market are built based on the risk neutral valuation principle w.r.t. an unknown risk neutral measure \mathbb{Q} . Our interest lies in estimating the conditional risk neutral distribution $Q_{S_T|S_t=s_t}$, or equivalently the risk neutral density function $q_{t,T}(\cdot|s_t)$, implied by \mathbb{Q} through (3.2).

3.2 Estimation

Adapting Breeden and Litzenberger (1978), one can show that the risk neutral density function $q_{t,T}(\cdot|s_t)$ is obtained as the second derivative of the call price function C_t with respect to strike K

$$q_{t,T}(s_T|s_t) = e^{r\tau} \left. \frac{\partial^2 C_t(s_t, r, K, \tau)}{\partial K^2} \right|_{K=s_T}, \tag{3.3}$$

where $\tau = T - t$ is the time to maturity.

The unknown price function C_t might be smoothed out of price observations and used in (3.3) to recover risk neutral densities. Here we follow the semiparametric approach from Ait-Sahalia and Lo (1998) where the smoothing is carried out in the space of implied volatilities.

The implied volatility surface is the function $v_t : \mathbb{R}_+^2 \rightarrow \mathbb{R}_+$ satisfying for all $(K, \tau) \in \mathbb{R}_+^2$

$$C_t(s_t, r, K, \tau) = C_{BS}\{s_t, r, K, \tau, v_t(K, \tau)\}, \tag{3.4}$$

where $C_{BS}(s, r, K, \tau, v) = s\Phi(d_1) - Ke^{-r\tau}\Phi(d_2)$ is the Black–Scholes price of Ψ with strike K and maturity τ , $\Phi(\cdot)$ is the cumulative distribution function of the standard normal distribution, $d_1 = \{\log(\frac{s}{K}) + (r + \frac{1}{2}v^2)\tau\}/(v\sqrt{\tau})$ and $d_2 = d_1 - v\sqrt{\tau}$.

More generally, the implied volatility surface is considered a smooth random function $\mathcal{V} \in L_2(\mathcal{X})$ on the space $\mathcal{X} \subset \mathbb{R}^2$ of strikes K and maturities τ . Combining (3.3) and (3.4), the functional random variable $\mathcal{H} \in L_2(\mathcal{X})$, called the risk neutral (RN) surface, is defined as

$$\begin{aligned} \mathcal{H}(s, r, K, \tau, \mathcal{V}) &= e^{r\tau} D^2 C_{BS}(s, r, K, \tau, \mathcal{V}) \\ &= \varphi(d_2) \left\{ \frac{1}{K\sqrt{\tau}\mathcal{V}} + \frac{2d_1}{\mathcal{V}} D\mathcal{V} + K\sqrt{\tau} \frac{d_1 d_2}{\mathcal{V}} (D\mathcal{V})^2 + K\sqrt{\tau} D^2 \mathcal{V} \right\}, \end{aligned} \tag{3.5}$$

where D^m denotes the m th partial derivative with respect to K and $\varphi(\cdot)$ the probability density function of the standard normal distribution. The explicit derivation of (3.5) and a detailed treatment of implied volatilities can be found in Hafner (2004) and Fengler (2005). Clearly, lower dimension objects describing \mathcal{V} may be used to analyze the RN surface \mathcal{H} .

A functional dataset containing realizations of the implied volatility surface \mathcal{V} is, however, not available, as in an exchange only discretized values of \mathcal{V}_t corrupted by noise are registered from trades. On each day $t = 1, \dots, T$ there are J_t options traded, each intraday trade $j = 1, \dots, J_t$ corresponds to an observed option price Y_{jt} at a pair of moneyness κ and maturities τ , $X_{jt} = (\kappa_{jt}, \tau_{jt})^\top$ where $\kappa = e^{r\tau} K/s_t$. Let $C_{BS}(v) = C_{BS}(v; s, r, K, \tau)$ denote the Black–Scholes price as a function of v with all other arguments held constant. As $C_{BS}(v)$ is continuous and monotone in v with inverse C_{BS}^{-1} , the observed implied volatility associated with trade j at day t is then $v_{jt} = C_{BS}^{-1}(Y_{jt})$. Figure 2 shows the implied volatilities from options on the German Stock Index DAX traded on 2 May 2000, the sparse and degenerated design is caused by regulation imposed by stock exchanges on the tradable maturities from call options.

For numerical tractability, see Fengler et al. (2007), observations v_{jt} are transformed into log-implied volatilities $\tilde{Y}_{jt} = \log v_{jt}$ and based on $\{(\tilde{Y}_{jt}, X_{jt})\}$, we use DSFM to model

$$\tilde{Y}_{jt} = Z_t^\top m(X_{jt}) + \varepsilon_{jt}. \tag{3.6}$$

The implied volatility surface at t is estimated by $\hat{\mathcal{V}}_t = \exp(\hat{Z}_t^\top \hat{\Gamma} \psi)$, recall (2.2). The RN surface is estimated using (3.5) by $\hat{\mathcal{H}}_t = \mathcal{H}(s_t, r, K, \tau, \hat{\mathcal{V}}_t)$. The dynamics of the unobservable sequence of RN surfaces $\{\mathcal{H}_t\}$ implied in the observations may be investigated by analyzing the lower dimensional $\{\hat{Z}_t\}$.

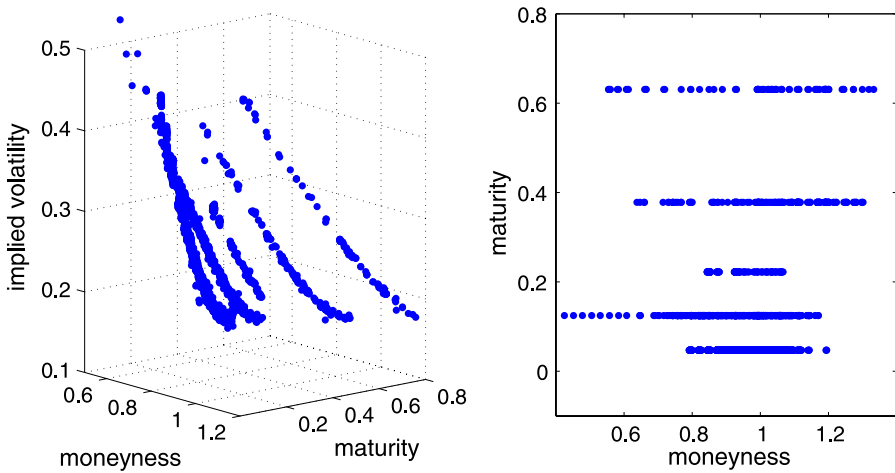


Fig. 2 Implied volatilities (left) and data design (right), ODAx on 2 May 2000

Table 1 Descriptive statistics, number of intraday observations $J_t, t = 1, \dots, 253$

Mean	Std. dev.	Max	Min
2845.92	1589.90	11298	616

4 Application

In this section, the implied volatility and risk neutral surfaces are estimated with DSFM from intraday prices of calls on the DAX index, i.e., S_t represents the value of the DAX index at time t . The dataset contains prices observed from 1 Jan. 2001 to 1 Jan. 2002 corresponding to $T = 253$ trading days. The descriptive statistics of the number of intraday observations J_t are in Table 1, the total number of intraday observations across days is $\sum_{t=1}^T J_t = 720017$.

Tensor B-splines, quadratic in τ and cubic in κ directions placed on 8×6 knots, are used for the series estimators of m . The number of basis functions is chosen based on

$$EV(L) = 1 - \frac{\sum_{t=1}^T \sum_{j=1}^{J_t} \{\tilde{Y}_{jt} - \hat{Z}_t^\top \hat{m}(X_{jt})\}^2}{\sum_{t=1}^T \sum_{j=1}^{J_t} (\tilde{Y}_{jt} - \bar{Y})^2},$$

where $\bar{Y} = (\sum_{t=1}^T \sum_{j=1}^{J_t} \tilde{Y}_{jt}) / \sum_{t=1}^T J_t$. The value $EV(L)$ may be interpreted as the ratio of variation explained by the model to total variation. As established by numerous simulations in Park et al. (2009), the order of the splines and number of knots have negligible influence on $EV(L)$.

4.1 Simulation

The choice of the number of basis functions based on the explained variation criteria is validated by a small simulation study. Datasets $\{(Y_{jt}, X_{jt})\}$ are generated following

$$\begin{aligned}
 Y_{jt} &= \sum_{l=0}^{L^*} Z_{lt} m_l(X_{jt}) + \varepsilon_{jt}, \quad j = 1, \dots, J, \quad t = 1, \dots, T, \\
 \varepsilon_{jt} &\sim N(0, \sigma_\varepsilon^2), \\
 X_{jt} &\sim U([0, 1]^2),
 \end{aligned}
 \tag{4.1}$$

where ε_{jt} and X_{jt} are i.i.d. For $\zeta_t = (Z_{1t}, \dots, Z_{L^*t})^\top$, with 0_d denoting the $(d \times 1)$ vector of zeros and I_d the d identity matrix we define

$$\begin{aligned}
 Z_t &= (1, \zeta_t)^\top, \\
 \zeta_t &= A_{L^*} \zeta_{t-1} + u_t, \\
 u_t &\sim N(0_{L^*}, \sigma_u^2 I_{L^*}),
 \end{aligned}$$

where u_t is i.i.d. and A_{L^*} is a square matrix containing the first L^* rows and L^* columns from A ,

$$A = \begin{pmatrix} 0.95 & -0.2 & 0 & 0.1 \\ 0 & 0.8 & 0.1 & 0.2 \\ 0.1 & 0 & 0.6 & -0.1 \\ 0 & 0.1 & -0.2 & 0.5 \end{pmatrix}.$$

The basis functions are defined as

$$\begin{aligned}
 m_0(\kappa, \tau) &= 1, \\
 m_1(\kappa, \tau) &= 3.46(\kappa - 0.5), \\
 m_2(\kappa, \tau) &= 9.45\{(\kappa - 0.5)^2 + (\tau - 0.5)^2\} - 1.6, \\
 m_3(\kappa, \tau) &= 1.41 \sin(2\pi\tau), \\
 m_4(\kappa, \tau) &= 1.41 \cos(2\pi\kappa),
 \end{aligned}$$

and are close to orthogonal, enhancing similar choice from Park et al. (2009). The value L^* denotes the true number of dynamic basis functions.

Setting $T = 500$, $J = 100$, $\sigma_\varepsilon = 0.05$, and $\sigma_u = 0.1$, $i = 1, \dots, 100$ samples following (4.1) are generated with $L^* = 2, 3$ and 4. Each of them is estimated by DSFM with $L = 1, \dots, 6$, and the corresponding $EV_i(L)$ is computed. The average explained variation under the true L^* , defined as $\mathcal{E}\mathcal{V}(L; L^*) = \frac{1}{100} \sum_i EV_i(L)$, is also calculated.

Table 2 shows $\mathcal{E}\mathcal{V}(L; L^*)$ and indicates that the increase in the average explained variation between estimation with L^* and $L^* + 1$ dynamic basis functions, $\mathcal{E}\mathcal{V}(L^* + 1; L^*) - \mathcal{E}\mathcal{V}(L^*; L^*)$, is close to zero across values of L^* . Therefore,

Table 2 Average explained variation $\mathcal{E}^{\mathcal{V}}(L; L^*)$ based on 100 samples from (4.1), across number of dynamic basis functions used in the estimation L and the true L^*

$\mathcal{E}^{\mathcal{V}}(L; L^*)$		L^*		
		2	3	4
L	1	0.86	0.75	0.71
	2	0.99	0.90	0.89
	3	0.99	0.99	0.97
	4	0.99	0.99	0.99
	5	0.99	0.99	0.99

Table 3 Number of basis functions and explained variation

L	1	2	3	4	5
$EV(L)$	0.77	0.97	0.98	0.98	0.98

for DSFM estimation, we select the smallest L such that $EV(L - 1) < EV(L) \approx EV(L + 1)$.

4.2 Results

The implied volatility and RN surfaces are estimated with DSFM as in (3.6) with $L = 3$. Table 3 shows that the addition of the fourth or fifth dynamic basis function results in negligible increase in $EV(L)$.

Following Fengler et al. (2007) and Park et al. (2009), the estimated \widehat{Z}_t and \widehat{m} are respectively transformed and orthonormalized so that $\{\widehat{Z}_t^\top \widehat{m}_l\}$ has a larger contribution than $\{\widehat{Z}_{(l+1)t}^\top \widehat{m}_{l+1}\}$, $l = 1, \dots, L - 1$, to the total variation $\sum_{t=1}^T \int \widehat{Z}_t^\top \widehat{m}$. This transformation aims to improve the interpretation of the basis functions in the analysis of the dynamics of implied volatility surfaces. In the analysis of risk neutral surfaces dynamics, however, it does not present a clear advantage. The covariance structures from $\{\widehat{Z}_t\}$ and $\{Z_t\}$ are then asymptotically equivalent up to orthogonal transformations.

Figures 3 and 4 depict the estimated loading factors series $\{\widehat{Z}_t\}$ and basis functions \widehat{m}_l . The upward and downward peaks observed in \widehat{Z}_{2t} occur on days 6 Feb. 2001 and 5 Nov. 2001 and are caused respectively by extremely unbalanced design and low price levels. The first day has $J_t = 1697$ observations concentrated on short maturities, while the latter has $J_t = 3268$ with very low prices at high maturities.

From (3.5), we obtain a sequence of RN surfaces $\{\widehat{\mathcal{H}}_t\}$, $t = 1, \dots, 253$. We define $\widehat{\mathcal{H}}_t(\kappa, \tau)$ as $\mathcal{H}(\kappa, \tau; s_t, r, \widehat{V}_t)$ where $\kappa = e^{r\tau} K / s_t$. Figure 5 shows $\widehat{\mathcal{H}}_t(\kappa, \tau)$ across moneyness κ and maturity τ at t corresponding to 10 Jul. 2001.

In a first step, we investigate the covariance structure of $\{\widehat{Z}_t\}$ by means of VAR analysis. Table 4 presents the parameters from the VAR(2) model fitted on $\{\widehat{Z}_t\}$. The order 2 is selected based on Akaike (AIC), Schwarz (SC) and Hannan–Quinn (HQ) criteria, see Table 5. Moreover, the VAR(2) model is stationary as the roots of the characteristic polynomial lie inside of the unit circle.

A natural issue is to analyze the dependences between $\{Z_t\}$ and the shape of the RN surfaces $\{\widehat{\mathcal{H}}_t\}$. In order to investigate this relation, we compute the skewness

Fig. 3 Estimated $\{\widehat{Z}_{lt}\}$, $l = 1, 2, 3$ (top to bottom)

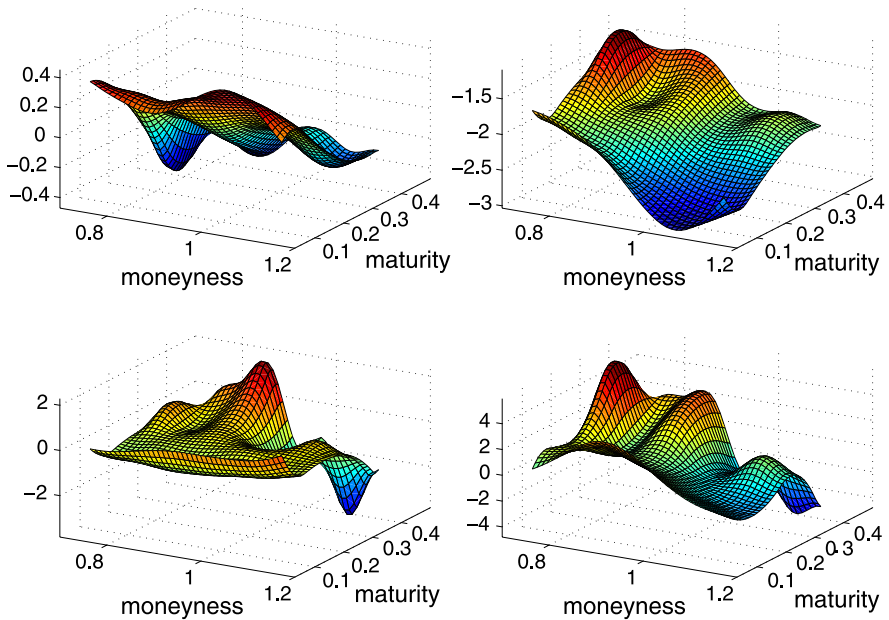
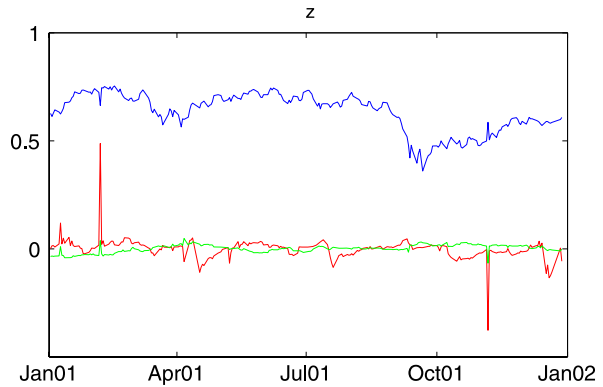


Fig. 4 Estimated basis functions $\widehat{m}_l, l = 0, \dots, 3$, clockwise

γ and excess kurtosis η of $\widehat{q}_{t,T}(\cdot|s_t)$ across t for a maturity τ where $\widehat{q}_{t,T}(\cdot|s_t) = \widehat{\mathcal{H}}_t(\cdot, \tau)$. Figure 6 displays the skewness $\{\gamma_t\}$ and excess kurtosis $\{\eta_t\}$ associated with $\widehat{q}_{t,T}$ for maturity $\tau = 18$ days together with $\{\widehat{Z}_{1t}\}$ and $\{\widehat{Z}_{3t}\}$, motivating the investigation of their joint autocovariance structure.

The dynamic structure of the pairs $\{(\widehat{Z}_{1t}, \eta_t)\}$ and $\{(\widehat{Z}_{3t}, \gamma_t)\}$ for $\tau = 18$ is modeled by VAR(2) models. The choice of the VAR order is again based on AIC, SC, and HQ selection criteria. Portmanteau and LM tests on VAR residuals reject autocorrelations up to lag 12 and the roots of the characteristic polynomial lie inside of the unit circle.

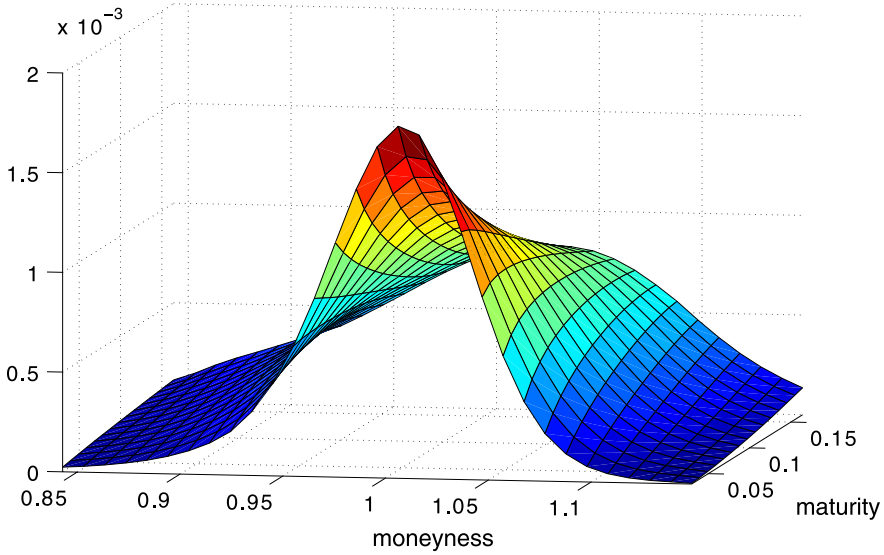


Fig. 5 Estimated RN surface, $\widehat{\mathcal{H}}_t$ at t corresponding to 10 Jul. 2001

Table 4 Estimated parameters for the VAR(2) model on $\{\widehat{Z}_t\}$

	VAR(2)						
	Const	$\widehat{Z}_{1,t-1}$	$\widehat{Z}_{1,t-2}$	$\widehat{Z}_{2,t-1}$	$\widehat{Z}_{2,t-2}$	$\widehat{Z}_{3,t-1}$	$\widehat{Z}_{3,t-2}$
\widehat{Z}_{1t}	0.01	1.09	-0.16	0.10	-0.36	0.32	-0.23
\widehat{Z}_{2t}	0.01	-0.27	0.26	0.31	0.12	-1.14	0.33
\widehat{Z}_{3t}	0.01	-0.08	0.62	-0.05	-0.04	0.41	0.35

Table 5 Lag selection criteria for VAR models on $\{\widehat{Z}_t\}$. The asterisks denote the smallest value for each criterion

Order	AIC	SC	HQ
1	-11.03	-10.99	-11.01
2	-15.71	-15.54*	-15.64*
3	-15.77*	-15.46	-15.64
4	-15.76	-15.32	-15.58
5	-15.72	-15.16	-15.45

Modeling the dynamics of risk neutral densities using DSFM allows quantifying the mechanisms governing risk perceptions from agents acting in a market. Insights are obtained in two directions, concerning the autocovariance structure of $\{\widehat{Z}_t\}$, i.e., the time behavior of the RN surfaces and their cross-correlation with the skewness and excess kurtosis from the estimated risk neutral densities, i.e., the relation between the dynamics and shape of the obtained RN surfaces. As seen in Tables 6 and 7 the excess kurtosis and skewness from $\widehat{q}_{t,T}$ at maturity $\tau = 18$ are determined by the corresponding lagged values of \widehat{Z}_t .

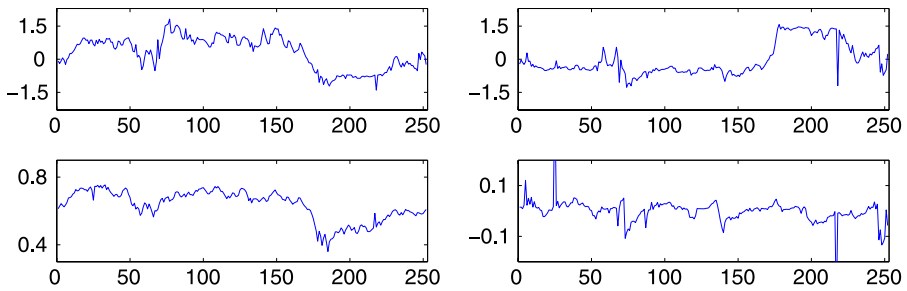


Fig. 6 Left: RN excess kurtosis $\{\eta_t\}$, $\tau = 18$ (top), $\{\widehat{Z}_{1t}\}$ (bottom). Right: RN skewness $\{\gamma_t\}$, $\tau = 18$ (top), $\{\widehat{Z}_{2t}\}$ (bottom)

Table 6 Estimated parameters for the VAR(2) model on $\{\widehat{Z}_{1t}, \eta_t\}$

	VAR(2)				
	Const	$\widehat{Z}_{1,t-1}$	$\widehat{Z}_{1,t-2}$	η_{t-1}	η_{t-2}
\widehat{Z}_{1t}	0.04	0.86	0.08	0.01	0.00
η_t	-0.51	2.63	-1.75	0.67	0.19

Table 7 Estimated parameters for the VAR(2) model on $\{\widehat{Z}_{3t}, \gamma_t\}$

	VAR(2)				
	Const	$\widehat{Z}_{3,t-1}$	$\widehat{Z}_{3,t-2}$	γ_{t-1}	γ_{t-2}
\widehat{Z}_{3t}	0.00	0.20	0.27	0.01	-0.02
γ_t	0.00	-1.69	0.68	0.81	0.24

The presented methodology allows the investigation of the dynamics from risk neutral skewness and excess kurtosis based on statistical inference on $\{\widehat{Z}_t\}$. A natural further step is to perform econometric analysis on the cointegration between the lower dimensional time series and macroeconomic and financial indicators. This could provide deeper insights into the relation between risk assessments from investors acting in a market and the flow of economic information at which they are exposed.

Acknowledgements Financial support from the Deutsche Forschungsgemeinschaft via SFB 649 “Economic Risk” is gratefully acknowledged. The authors also thank the editor, an associate editor and two referees for their helpful comments.

Appendix: Assumptions

The results from Theorems 2.1 and 2.2, see Park et al. (2009), rely on the following assumptions:

- (A1) The variables X_{11}, \dots, X_{JT} , $\varepsilon_{11}, \dots, \varepsilon_{JT}$ and Z_1, \dots, Z_T are independent. The process Z_t is allowed to be nonrandom.
- (A2) For $t = 1, \dots, T$, the variables X_{1t}, \dots, X_{Jt} are identically distributed, have support $[0, 1]^d$ and a density f_t that is bounded from below and above on $[0, 1]^d$, uniformly over $t = 1, \dots, T$.

(A3) We assume that $E[\varepsilon_{jt}] = 0$ for $t = 1, \dots, T$ and $j = 1, \dots, J$, and

$$\sup_{t=1, \dots, T, j=1, \dots, J} E \exp[c\varepsilon_{jt}^2] < \infty$$

for $c > 0$ small enough.

(A4) The functions ψ_k may depend on the increasing indices T and J and are normed so that $\int_{[0,1]^d} \psi_k^2(x) dx = 1$ for $k = 1, \dots, K$. Furthermore, $\sup_{x \in [0,1]^d} \|\psi(x)\| = \mathcal{O}(K^{1/2})$.

(A5) The components m_0, \dots, m_L can be approximated by ψ_1, \dots, ψ_K , i.e.,

$$\delta_K = \sup_{x \in [0,1]^d} \inf_{\Gamma \in \mathcal{G}} |m(x) - \Gamma \psi(x)| \rightarrow 0 \tag{A.1}$$

for $l = 0, \dots, L$ and $K \rightarrow \infty$. We denote by Γ^* the matrix that fulfills

$$\sup_{x \in [0,1]^d} |m(x) - \Gamma^* \psi(x)| \leq 2\delta_K.$$

(A6) There exist constants $0 < C_L < C_U < \infty$ such that all eigenvalues of the random matrix $T^{-1} \sum_{t=1}^T Z_t Z_t^\top$ lie in the interval $[C_L, C_U]$ with probability tending to one.

(A7) The minimization (2.2) runs over all values of (Γ, z) with

$$\sup_{x \in [0,1]^d} \max_{1 \leq t \leq T} \|Z_t^\top \Gamma \psi(x)\| \leq M_T,$$

where M_T fulfills $\max_{1 \leq t \leq T} \|Z_t\| \leq M_T/C_m$ (with probability tending to one) for a constant $C_m > \sup_{x \in [0,1]^d} \|m(x)\|$.

(A8) It holds that

$$\xi^2 = (K + T)M_T^2 \log(JTM_T)(JT)^{-1} \rightarrow 0, \tag{A.2}$$

where the dimension L is fixed.

(A9) Z_t is a martingale difference with $E[Z_t|Z_1, \dots, Z_{t-1}] = 0$ and for some $C > 0$ $E[\|Z_t\|^2|Z_1, \dots, Z_{t-1}] < C$ (a.s.). The matrix $E[Z_t Z_t^\top]$ has full rank. The process Z_t is independent of X_{11}, \dots, X_{TJ} and $\varepsilon_{11}, \dots, \varepsilon_{TJ}$.

(A10) The functions m_0, \dots, m_L are linearly independent. In particular, no function is equal to 0.

(A11) It holds that $(K^{1/2}M_T + T^{1/4})(\xi + \delta_K) = \mathcal{O}(1)$.

References

Ait-Sahalia, Y., Lo, A.: Nonparametric estimation of state-price densities implicit in financial asset prices. *J. Finance* **53**, 499–547 (1998)

Benko, M., Kneip, A., Härdle, W.: Common functional principal components. *Ann. Stat.* **37**(1), 1–34 (2009)

Breeden, D., Litzenberger, R.: Prices of state-contingent claims implicit in options prices. *J. Bus.* **51**, 621–651 (1978)

- Brüggemann, R., Härdle, W., Mungo, J., Trenkler, C.: VAR modeling for dynamic loadings driving volatility strings. *J. Financ. Econ.* **6**, 361–381 (2008)
- Cont, R., da Fonseca, J.: The dynamics of implied volatility surfaces. *Quant. Finance* **2**, 45–60 (2002)
- Dauxois, J., Pousse, A., Romain, Y.: Asymptotic theory for the principal component analysis of a vector random function: some applications to statistical inference. *J. Multivar. Anal.* **12**, 136–154 (1982)
- Fengler, M.: *Semiparametric Modeling of Implied Volatility*. Springer, Heidelberg (2005)
- Fengler, M., Härdle, W., Mammen, E.: A semiparametric factor model for implied volatility surface dynamics. *J. Financ. Econ.* **5**, 189–218 (2007)
- Gasser, T., Kneip, A.: Searching for structure in curve samples. *J. Am. Stat. Assoc.* **90**(432), 1179–1188 (1995)
- Hafner, R.: *Stochastic Implied Volatility*. Springer, Heidelberg (2004)
- Hall, P., Müller, H., Wang, J.: Properties of principal component methods for functional and longitudinal data analysis. *Ann. Stat.* **34**(3), 1493–1517 (2006)
- Hernández-Hernández, D., Schied, A.: A control approach to robust maximization with logarithmic utility and time-consistent penalties. *Stoch. Process. Appl.* **117**(8), 980–1000 (2007)
- Kneip, A.: Nonparametric estimation of common regressors for similar curve data. *Ann. Stat.* **22**(3), 1386–1427 (1994)
- Kneip, A., Gasser, T.: Statistical tools to analyse data representing a sample of curves. *Ann. Stat.* **20**(3), 1266–1305 (1992)
- Park, B., Mammen, E., Härdle, W., Borak, S.: Time series modelling with semiparametric factor dynamics. *J. Am. Stat. Assoc.* **104**(485), 284–298 (2009)
- Ramsay, J.O., Dalzell, C.T.: Some tools for functional data analysis. *J. R. Stat. Soc. B* **53**(3), 539–572 (1991)
- Rice, J., Silverman, B.W.: Estimating the mean and covariance structure nonparametrically when the data are curves. *J. R. Stat. Soc. B* **53**, 233–243 (1991)

Inhomogeneous Dependence Modeling With Time-Varying Copulae

Enzo GIACOMINI

Center for Applied Statistics and Economics, Humboldt-Universität zu Berlin, 10178 Berlin, Germany; and Weierstrass Institute for Applied Analysis and Stochastics, 10117 Berlin, Germany (*giacomini@wiwi.hu-berlin.de*)

Wolfgang HÄRDLE

AU1 Weierstrass Institute for Applied Analysis and Stochastics, 10117 Berlin, Germany

Vladimir SPOKOINY

Center for Applied Statistics and Economics, Humboldt-Universität zu Berlin, 10178 Berlin, Germany; and Weierstrass Institute for Applied Analysis and Stochastics, 10117 Berlin, Germany

Measuring dependence in multivariate time series is tantamount to modeling its dynamic structure in space and time. In risk management, the nonnormal behavior of most financial time series calls for non-Gaussian dependences. The correct modeling of non-Gaussian dependences is, therefore, a key issue in the analysis of multivariate time series. In this article we use copula functions with adaptively estimated time-varying parameters for modeling the distribution of returns. Furthermore, we apply copulae to the estimation of value-at-risk of portfolios and show their better performance over the *RiskMetrics* approach.

KEY WORDS: Adaptive estimation; Nonparametric estimation; Value-at-risk.

1. INTRODUCTION

Time series of financial data are high dimensional and typically have a non-Gaussian behavior. The standard modeling approach based on properties of the multivariate normal distribution therefore often fails to reproduce the stylized facts (i.e., fat tails, asymmetry) observed in returns from financial assets.

A correct understanding of the time-varying multivariate (conditional) distribution of returns is vital to many standard applications in finance such as portfolio selection, asset pricing, and value-at-risk (var) calculation. Empirical evidence from asymmetric return distributions have been reported in the recent literature. Longin and Solnik (2001) investigate the distribution of joint extremes from international equity returns and reject multivariate normality in their lower orthant; Ang and Chen (2002) test for conditional correlation asymmetries in U.S. equity data, rejecting multivariate normality at daily, weekly, and monthly frequencies; and Hu (2006) models the distribution of index returns with mixtures of copulae, finding asymmetries in the dependence structure across markets. For a concise survey on stylized empirical facts from financial returns see Cont (2001) and Granger (2003).

Modeling distributions with copulae has drawn attention from many researchers because it avoids the “procrustean bed” of normality assumptions, producing better fits of the empirical characteristics of financial returns. A natural extension is to apply copulae in a dynamic framework with conditional distributions modeled by copulae with time-varying parameters. The question, though, is how to steer the time-varying copulae parameters. This question is the focus of this article.

A possible approach is to estimate the parameter from structurally invariant periods. There is a broad field of econometric literature on structural breaks. Tests for unit root in macroeconomic series against stationarity with a structural

break at a known change point have been investigated by Perron (1989), and for an unknown change point by Zivot and Andrews (1992), Stock (1994) and Hansen (2001); Andrews (1993) tests for parameter instability in nonlinear models; Andrews and Ploberger (1994) construct asymptotic optimal tests for multiple structural breaks. In a different set up, Quintos, Fan, and Philips (2001) test for a constant tail index coefficient in Asian equity data against a break at an unknown point.

Time-varying copulae and structural breaks are combined in Patton (2006). The dependence structure across exchange rates is modeled with time-varying copulae with a parameter specified to evolve as an ARMA-type process. Tests for a structural break in the ARMA coefficients at a known change point have been performed, and strong evidence of a break was found. In a similar fashion, Rodriguez (2007) models the dependence across sets of Asian and Latin American stock indexes using time-varying copula where the parameter follows regime-switching dynamics. Common to these articles is that they use a fixed (parametric) structure for the pattern of changes in the copula parameter. **AU2**

In this article we follow a semiparametric approach, because we are not specifying the parameter changing scheme. Rather, we select locally the time-varying copula parameter. The choice is performed via an adaptive estimation under the assumption of local homogeneity: For every time point there exists an interval of time homogeneity in which the copula parameter can be well approximated by a constant. This interval is recovered from the data using local change point analysis. This does not imply that the model follows a change

point structure. The adaptive estimation also applies when the parameter varies smoothly from one value to another (see Spokoiny 2008).

[F1] Figure 1 shows the time-varying copula parameter determined by our procedure for a portfolio composed of daily prices of six German equities and the “global” copula parameter, shown by a constant horizontal line. The absence of parametric specification for time variations in the dependence structure (its dynamics is obtained adaptively from the data) allows for flexibility in estimating dependence shifts across time.

The obtained time-varying dependence structure can be used in financial engineering applications, the most prominent being the calculation of the var of a portfolio. Using copulae with adaptively estimated dependence parameters we estimate the var from DAX portfolios over time. As a benchmark procedure **[AU3]** we choose *RiskMetrics*, a widely used methodology based on conditional normal distributions with a GARCH specification for the covariance matrix. Backtesting underlines the improved performance of the proposed adaptive time-varying copulae fitting.

This article is organized as follows: Section 2 presents the basic copulae definitions, Section 3 discusses the var and its estimation procedure. The adaptive copula estimation is exposed in Section 4 and is applied to simulated data in Section 5. In Section 6, the var from DAX portfolios is estimated based on adaptive time-varying copulae. The estimation performance is compared with the *RiskMetrics* approach by means of back-testing. Section 7 concludes.

2. COPULAE

Copulae merge marginally into joint distributions, providing a natural way for measuring the dependence structure between random variables. Copulae are present in the literature since Sklar (1959), although related concepts originate in Hoeffding (1940) and Fréchet (1951), and have been widely studied in the statistical literature (see Joe 1997, Nelsen 1998, and Mari and Kotz 2001). Applications of copulae in finance, insurance, and econometrics have been investigated in Embrechts, McNeil, and Straumann (2002); Embrechts, Hoeing, and Juri (2003a); Franke, Härdle, and Hafner (2004); and Patton (2004) among others. Cherubini, Luciano, and Vecchiato (2004) and McNeil, Frey, and Embrechts (2005) provide an overview of copulae for practical problems in finance and insurance.

Assuming absolutely continuous distributions and continuous marginals throughout this article, we have from Sklar’s

theorem that for a d -dimensional distribution function F with marginal cdf’s F_1, \dots, F_d there exists a unique copula $C : [0, 1]^d \rightarrow [0, 1]$ satisfying

$$F(x_1, \dots, x_d) = C\{F_1(x_1), \dots, F_d(x_d)\} \quad (2.1)$$

for every $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathbb{R}^d$. Conversely, for a random vector $\mathbf{X} = (X_1, \dots, X_d)^T$ with cdf $F_{\mathbf{X}}$, the copula of \mathbf{X} may be written as $C_{\mathbf{X}}(u_1, \dots, u_d) = F_{\mathbf{X}}\{F_1^{-1}(u_1), \dots, F_d^{-1}(u_d)\}$, where $u_j = F_j(x_j)$, F_j is the cdf of X_j , and $F_j^{-1}(\alpha) = \inf\{x_j : F_j(x_j) \geq \alpha\}$ its generalized inverse, $j = 1, \dots, d$. A prominent copula is the Gaussian

$$C_{\Psi}^{Ga}(u_1, \dots, u_d) = F_Y\{\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_d)\} \quad (2.2)$$

where $\Phi(s)$, $s \in \mathbb{R}$ stands for the one-dimensional standard normal cdf, F_Y is the cdf of $\mathbf{Y} = (Y_1, \dots, Y_d)^T \sim N_d(\mathbf{0}, \Psi)$, $\mathbf{0}$ is the $(d \times 1)$ vector of zeros, and Ψ is a correlation matrix. The Gaussian copula represents the dependence structure of the multivariate normal distribution. In contrast, the Clayton copula given by

$$C_{\theta}(u_1, \dots, u_d) = \left\{ \left(\sum_{j=1}^d u_j^{-\theta} \right) - d + 1 \right\}^{-\theta^{-1}} \quad (2.3)$$

for $\theta > 0$, expresses asymmetric dependence structures.

The dependence at upper and lower orthants of a copula C may be expressed by the upper and lower tail dependence coefficients $\lambda_U = \lim_{u \rightarrow 0} \hat{C}(u, \dots, u)/u$ and $\lambda_L = \lim_{u \rightarrow 0} C(u, \dots, u)/u$, where $u \in (0, 1]$ and \hat{C} is the survival copula of C (see Joe 1997 and Embrechts, Lindskog, and McNeil 2003b). Although Gaussian copulae are asymptotically independent at the tails ($\lambda_L = \lambda_U = 0$), the d -dimensional Clayton copulae exhibit lower tail dependence ($\lambda_L = d^{-1/\theta}$) but are asymptotically independent at the upper tail ($\lambda_U = 0$). Joe (1997) provides a summary of diverse copula families and detailed description of their properties.

For estimating the copula parameter, consider a sample $\{\mathbf{x}_t\}_{t=1}^T$ of realizations from \mathbf{X} where the copula of \mathbf{X} belongs to a parametric family $C = \{C_{\theta}, \theta \in \Theta\}$. Using Equation (2.1), the log-likelihood reads as $L(\theta; \mathbf{x}_1, \dots, \mathbf{x}_T) = \sum_{t=1}^T [\log c(F_1(x_{t,1}), \dots, F_d(x_{t,d}); \theta) + \sum_{j=1}^d \log f_j(x_{t,j})]$, where $c(u_1, \dots, u_d) = \partial^d C(u_1, \dots, u_d) / \partial u_1 \dots \partial u_d$ is the density of the copula C and f_j is the probability density function of F_j . The canonical maximum likelihood estimator $\hat{\theta}$ maximizes the pseudo log-likelihood with empirical marginal cdf’s $\hat{L}(\theta) = \sum_{t=1}^T \log c\{\hat{F}_1(x_{t,1}), \dots, \hat{F}_d(x_{t,d}); \theta$, where

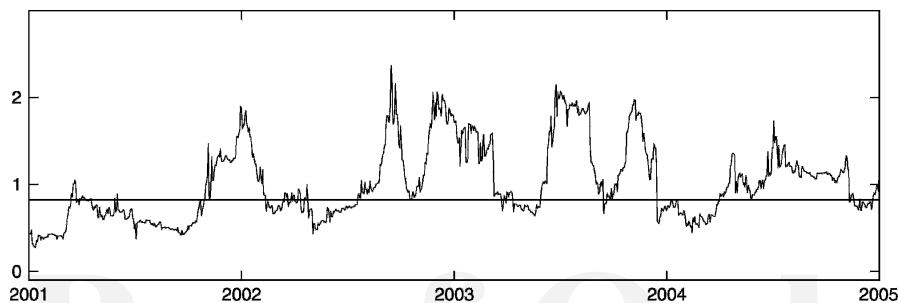


Figure 1. Time-varying dependence. Time-varying dependence parameter and global parameter (horizontal line) estimated with Clayton copula, stock returns from Allianz, Münchener Rückversicherung, BASF, Bayer, DaimlerChrysler, and Volkswagen.

$$\widehat{F}_j(s) = \frac{1}{T+1} \sum_{k=1}^T 1_{\{x_{kj} \leq s\}} \quad (2.4)$$

for $j = 1, \dots, d$. Note that \widehat{F}_j differs from the usual empirical cdf by the denominator $T + 1$. This ensures that $\{\widehat{F}_1(x_{t,1}), \dots, \widehat{F}_d(x_{t,d})\}^\top \in (0, 1)^d$ and avoids infinite values the copula density may take on the boundary of the unit cube (see McNeil, Frey, and Embrechts 2005). Joe (1997); Cherubini, Luciano, and Vecchiato (2004); and Chen and Fan (2006) provide a detailed exposition of inference methods for copulae.

3. VALUE-AT-RISK AND COPULAE

The dependence (over time) between asset returns is especially important in risk management, because the profit and loss (P&L) function determines the var. More precisely, the var of a portfolio is determined by the multivariate distribution of risk factor increments. If $\mathbf{w} = (w_1, \dots, w_d)^\top \in \mathbb{R}^d$ denotes a portfolio of positions on d assets and $\mathbf{S}_t = (S_{t,1}, \dots, S_{t,d})^\top$ a nonnegative random vector representing the prices of the assets at time t , the value V_t of the portfolio \mathbf{w} is given by $V_t = \sum_{j=1}^d w_j S_{t,j}$. The random variable

$$L_t = (V_t - V_{t-1}), \quad (3.1)$$

called the profit and loss (P&L) function, expresses the change in the portfolio value between two subsequent time points. Defining the log-returns $\mathbf{X}_t = (X_{t,1}, \dots, X_{t,d})^\top$, where $X_{t,j} = \log S_{t,j} - \log S_{t-1,j}$ and $\log S_{0,j} = 0, j = 1, \dots, d$, Equation (3.1) can be written as

$$L_t = \sum_{j=1}^d w_j S_{t-1,j} \{\exp(X_{t,j}) - 1\}. \quad (3.2)$$

The cdf of L_t is given by $F_{t,L_t}(x) = P_t(L_t \leq x)$. The var at level α from a portfolio \mathbf{w} is defined as the α quantile from F_{t,L_t} :

$$\text{var}_t(\alpha) = F_{t,L_t}^{-1}(\alpha). \quad (3.3)$$

It follows from Equation (3.2) that F_{t,L_t} depends on the specification of the d -dimensional distribution of the risk factors \mathbf{X}_t . Thus, modeling their distribution over time is essential for obtaining the quantiles (Eq. 3.3).

The *RiskMetrics* technique, a widely used methodology for var estimation, assumes that risk factors \mathbf{X}_t follow a conditional multivariate normal distribution $\mathcal{L}(\mathbf{X}_t | \mathcal{F}_{t-1}) = N(\mathbf{0}, \boldsymbol{\Sigma}_t)$, where $\mathcal{F}_{t-1} = \sigma(\mathbf{X}_1, \dots, \mathbf{X}_{t-1})$ is the σ field generated by the first $t - 1$ observations, and estimates the covariance matrix $\boldsymbol{\Sigma}_t$ for one period return as

$$\widehat{\boldsymbol{\Sigma}}_t = \lambda \widehat{\boldsymbol{\Sigma}}_{t-1} + (1 - \lambda) \mathbf{X}_{t-1} \mathbf{X}_{t-1}^\top, \quad (3.4)$$

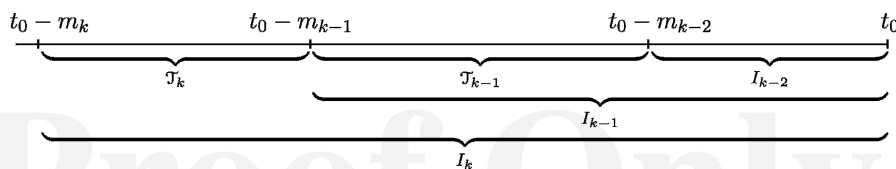


Figure 2. Local change point procedure. Choice of intervals I_k and I_{k-1} .

where the parameter λ is the so-called decay factor. $\lambda = 0.94$ provides the best backtesting results for daily returns according to Morgan (1996). Using the copulae-based approach, one first corrects the contemporaneous mean and volatility in the log-returns process:

$$X_{t,j} = \mu_{t,j} + \sigma_{t,j} \varepsilon_{t,j}, \quad (3.5)$$

where $\mu_{t,j} = E[X_{t,j} | \mathcal{F}_{t-1}]$ is the conditional mean and $\sigma_{t,j}^2 = E[(X_{t,j} - \mu_{t,j})^2 | \mathcal{F}_{t-1}]$ is the conditional variance of $X_{t,j}$. The standardized innovations $\boldsymbol{\varepsilon}_t = (\varepsilon_{t,1}, \dots, \varepsilon_{t,d})^\top$ have joint cdf $F_{\boldsymbol{\varepsilon}_t}$ given by

$$F_{\boldsymbol{\varepsilon}_t}(x_1, \dots, x_d) = C_\theta \{F_{t,1}(x_1), \dots, F_{t,d}(x_d)\}, \quad (3.6)$$

where $F_{t,j}$ is the cdf of $\varepsilon_{t,j}$ and C_θ is a copula belonging to a parametric family $C = C_\theta, \theta \in \Theta$. For details on the previous model specification, see Chen and Fan (2006) and Chen, Fan, and Tsyrennikov (2006). For the Gaussian copula with Gaussian marginals, we recover the conditional Gaussian *RiskMetrics* framework.

To obtain the var in this setup, the dependence parameter and cdf's from residuals are estimated from a sample of log-returns and are used to generate P&L Monte Carlo samples. Their quantiles at different levels are the estimators for the var (see Embrechts, McNeil, and Straumann 2002).

The whole procedure can be summarized as follows (see Härdle, Kleinow, and Stahl 2002; and Giacomini and Härdle 2005): For a portfolio $\mathbf{w} \in \mathbb{R}^d$ and a sample $\{x_{t,j}\}_{t=1}^T, j = 1, \dots, d$ of log-returns, the var at level α is estimated according to the following steps:

1. Determination of innovations $\{\hat{\varepsilon}_t\}_{t=1}^T$ by, for example, “deGARCHing”
2. Specification and estimation of marginal cdf's $F_j(\hat{\varepsilon}_j)$
3. Specification of a parametric copula family C and estimation of the dependence parameter θ
4. Generation of Monte Carlo sample of innovations ε and losses L
5. Estimation of $\widehat{\text{var}}(\alpha)$, the empirical α quantile of F_L

4. MODELING WITH TIME-VARYING COPULAE

Similar to the *RiskMetrics* procedure, one can perform a moving (fixed-length) window estimation of the copula parameter. This procedure, though, does not fine-tune local changes in dependences. In fact, the cdf $F_{\boldsymbol{\varepsilon}_t}$ from Equation (3.6) is modeled as $F_{t,\boldsymbol{\varepsilon}_t} = C_\theta \{F_{t,1}(\cdot), \dots, F_{t,d}(\cdot)\}$ with probability measure P_θ . The moving window of fixed width will estimate a θ_t for each t , but it has clear limitations. The choice of a small window results in a high pass filtering and, hence, in a very unstable estimate with huge variability. The choice of a large window leads to a poor sensitivity of the estimation procedure

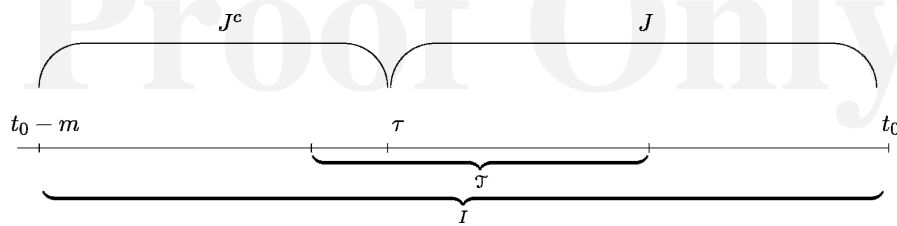


Figure 3. Homogeneity test. Testing interval I , tested interval I , and subintervals J and J^c for a point $\tau \in I$.

and to a high delay in the reaction to changes in dependence measured by the parameter θ_r .

To choose an interval of homogeneity, we use a local parametric fitting approach as introduced by Polzehl and Spokoiny (2006), Belomestny and Spokoiny (2007) and Spokoiny (2008). The basic idea is to select for each time point t_0 an interval $I_{t_0} = [t_0 - m_{t_0}, t_0]$ of length m_{t_0} in such a way that the time-varying copula parameter θ_t can be well approximated by a constant value θ . The question is, of course, how to select m_{t_0} in an online situation from historical data. The aim should be to select I_{t_0} as close as possible to the so-called ‘‘oracle’’ choice interval. The oracle choice is defined as the largest interval $I = [t_0 - m_{t_0}^*, t_0]$, for which the small modeling bias condition

$$\Delta_I(\theta) = \sum_{t \in I} \mathcal{K}(P_{\theta_t}, P_{\theta}) \leq \Delta \quad (4.1)$$

for some $\Delta \geq 0$ holds. Here, θ is constant and $\mathcal{K}(P_{\vartheta}, P_{\vartheta'}) = E_{\vartheta} \log\{p(y, \vartheta)/p(y, \vartheta')\}$ denotes the Kullback-Leibler divergence. In such an oracle choice interval, the parameter $\theta_{t_0} = \theta_t|_{t=t_0}$ can be ‘‘optimally’’ estimated from $I = [t_0 - m_{t_0}^*, t_0]$. The error and risk bounds are calculated in Spokoiny (2008). It is important to mention that the concept of local parametric approximation allows one to treat in a unified way the case of ‘‘switching regime’’ models with spontaneous changes of parameters and the ‘‘smooth transition’’ case when the parameter varies smoothly in time.

The oracle choice of the interval of homogeneity depends on the unknown time-varying copula parameter θ_r . The next section presents an adaptive (data-driven) procedure that mimics the oracle in the sense that it delivers the same accuracy of estimation as the oracle one. The trick is to find the largest interval in which the hypothesis of a local constant copula

parameter is supported. The local change point (LCP) detection procedure originates from Mercurio and Spokoiny (2004) and sequentially tests the hypothesis: θ_t is constant (i.e., $\theta_t = \theta$) within some interval I (local parametric assumption).

The LCP procedure for a given point t_0 starts with a family of nested intervals $I_0 \subset I_1 \subset I_2 \subset \dots \subset I_K = I_{K+1}$ of the form $I_k = [t_0 - m_k, t_0]$. The sequence m_k determines the length of these interval ‘‘candidates’’ (see Section 4.2). Every interval I_k leads to an estimate $\hat{\theta}_k$ of the copula parameter θ_{t_0} . The procedure selects one interval \hat{I} out of the given family and, therefore, the corresponding estimate $\hat{\theta} = \hat{\theta}_{\hat{I}}$.

The idea of the procedure is to screen each interval $J_k = [t_0 - m_k, t_0 - m_{k-1}]$ sequentially and check each point $\tau \in J_k$ as a possible change point location (see Section 4.1 for more details). The family of intervals I_k and J_k are illustrated in Figure 2. The interval I_k is accepted if no change point is detected within J_1, \dots, J_k . If the hypothesis of homogeneity is rejected for an interval candidate I_k , the procedure stops and selects the latest accepted interval. The formal description reads as follows:

Start the procedure with $k = 1$ and test the hypothesis $H_{0,k}$ of no structural changes within J_k using the larger testing interval I_{k+1} . If no change points were found in J_k , then I_k is accepted. Take the next interval J_{k+1} and repeat the previous step until homogeneity is rejected or the largest possible interval $I_K = [t_0 - m_K, t_0]$ is accepted. If $H_{0,k}$ is rejected for J_k , the estimated interval of homogeneity is the last accepted interval $\hat{I} = I_{k-1}$. If the largest possible interval I_K is accepted, we take $\hat{I} = I_K$. We estimate the copula dependence parameter θ at time instant t_0 from observations in \hat{I} , assuming the homogeneous model within \hat{I} (i.e., we define $\hat{\theta}_{t_0} = \hat{\theta}_{\hat{I}}$). We also denote by \hat{I}_k the largest accepted interval after k steps of

Table 1. Critical values $z_k(\rho, \theta^*)$

k	$\theta^* = 0.5$			$\theta^* = 1.0$			$\theta^* = 1.5$		
	$\rho = 0.2$	$\rho = 0.5$	$\rho = 1.0$	$\rho = 0.2$	$\rho = 0.5$	$\rho = 1.0$	$\rho = 0.2$	$\rho = 0.5$	$\rho = 1.0$
1	3.64	3.29	2.88	3.69	3.29	2.84	3.95	3.49	2.96
2	3.61	3.14	2.56	3.43	2.91	2.35	3.69	3.02	2.78
3	3.31	2.86	2.29	3.32	2.76	2.21	3.34	2.80	2.09
4	3.19	2.69	2.07	3.04	2.57	1.80	3.14	2.55	1.86
5	3.05	2.53	1.89	2.92	2.22	1.53	2.95	2.65	1.49
6	2.87	2.26	1.48	2.92	2.17	1.19	2.83	2.04	0.94
7	2.51	1.88	1.02	2.64	1.82	0.56	2.62	1.79	0.31
8	2.49	1.72	0.35	2.33	1.39	0.00	2.35	1.33	0.00
9	2.18	1.23	0.00	2.03	0.81	0.00	2.10	0.60	0.00
10	0.92	0.00	0.00	0.82	0.00	0.00	0.79	0.00	0.00

NOTE: Critical values are obtained according to Equation (4.2), based on 5,000 simulations. Clayton copula, $m_0 = 20$ and $c = 1.25$.

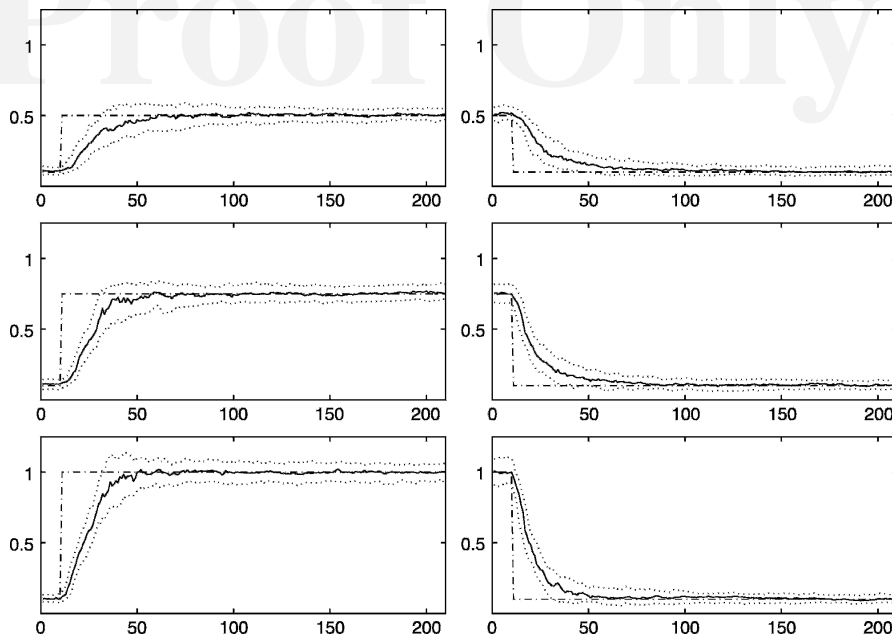


Figure 4. LCP and sudden jump in copula parameter. Pointwise median (full), and 0.25 and 0.75 quantiles (dotted) from $\hat{\theta}_t$. True parameter θ_t (dashed) with $\vartheta_a = 0.10$, $\vartheta_b = 0.50, 0.75$, and 1.00 (left, top to bottom); and $\vartheta_b = 0.10$, $\vartheta_a = 0.50, 0.75$, and 1.00 (right, top to bottom). Based on 100 simulations from Clayton copula, estimated with LCP, $m_0 = 20$, $c = 1.25$, and $\rho = 0.5$.

the algorithm and, by $\hat{\theta}_k$ the corresponding estimate of the copula parameter.

It is worth mentioning that the objective of the described estimation algorithm is not to detect the points of change for the copula parameter, but rather to determine the current dependence structure from historical data by selecting an interval of time homogeneity. This distinguishes our approach from other procedures for estimating a time-varying parameter by change point detection. A visible advantage of our approach is that it equally applies to the case of spontaneous changes in the dependence structure and in the case of smooth transition in the copula parameter. The obtained dependence structure can be used for different purposes in financial engineering, the most prominent being the calculation of the var (see also Section 6).

The theoretical results from Spokoiny and Chen (2007) and Spokoiny (2008) indicate that the proposed procedure provides the rate optimal estimation of the underlying parameter when this varies smoothly with time. It has also been shown that the procedure is very sensitive to structural breaks and provides the minimal possible delay in detection of changes, where the delay depends on the size of change in terms of Kullback-Leibler divergence.

4.1 Test of Homogeneity Against a Change Point Alternative

In the homogeneity test against a change point alternative we want to check every point of an interval I (recall Fig. 2), here called the “tested interval,” on a possible change in the dependence structure at this moment. To perform this check, we assume a larger testing interval I of form $I = [t_0 - m, t_0]$, so that I is an internal subset within I . The null hypothesis H_0 means that $\forall t \in I, \theta_t = \theta$ (i.e., the observations in I follow the

model with dependence parameter θ). The alternative hypothesis H_1 claims that $\exists \tau \in I$ such that $\theta_t = \theta_1$ for $t \in J = [\tau, t_0]$ and $\theta_t = \theta_2 \neq \theta_1$ for $t \in J^c = [t_0 - m, \tau]$ (i.e., the parameter θ changes spontaneously in some point $\tau \in I$). Figure 3 depicts I, I , and the subintervals J and J^c determined by the point $\tau \in I$.

Let $L_I(\theta)$ be the log-likelihood and $\tilde{\theta}_I$ the maximum likelihood estimate for the interval I . The log-likelihood functions corresponding to H_0 and H_1 are $L_I(\theta)$ and $L_J(\theta_1) + L_{J^c}(\theta_2)$, respectively. The likelihood ratio test for the single change point with known fixed location τ can be written as

Table 2. Detection delay statistics

$(\vartheta_a, \vartheta_b)$	r	Mean	SD	Max	Min
(0.50, 0.10)	0.25	9.06	7.28	56	0
	0.50	13.64	9.80	60	0
	0.75	21.87	14.52	89	3
(0.75, 0.10)	0.25	5.16	4.24	21	0
	0.50	8.85	5.55	25	0
	0.75	16.72	10.37	64	3
(1.00, 0.10)	0.25	4.47	2.94	12	0
	0.50	7.94	4.28	22	0
	0.75	14.79	7.38	62	5
(0.10, 0.50)	0.25	8.94	6.65	36	0
	0.50	14.21	9.06	53	0
	0.75	21.43	12.15	68	0
(0.10, 0.75)	0.25	9.00	4.80	25	0
	0.50	14.30	5.96	40	3
	0.75	21.00	10.97	75	6
(0.10, 1.00)	0.25	7.39	3.67	19	0
	0.50	13.10	4.13	22	2
	0.75	20.13	7.34	55	10

NOTE: The detection delays δ are calculated as in Equation (5.1), with the statistics based on 100 simulations. Clayton copula, $m_0 = 20$, $c = 1.25$, and $\rho = .5$. SD, standard deviation.

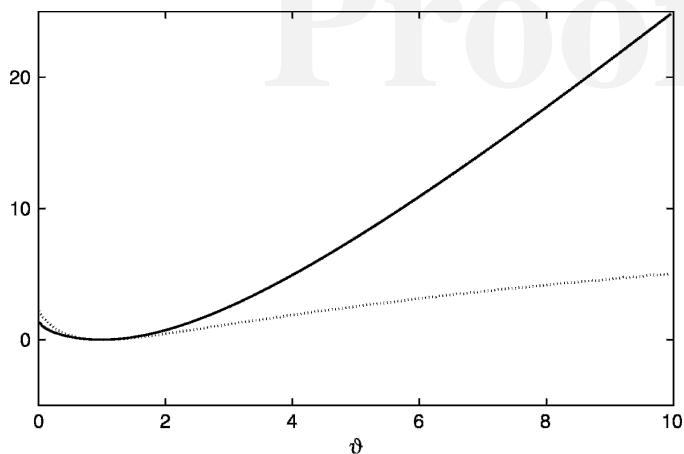


Figure 5. Divergences for upward and downward jumps. Kullback-Leibler divergences $\mathcal{K}(0.10, \vartheta)$ (full) and $\mathcal{K}(\vartheta, 0.10)$ (dashed) for Clayton copula.

$$T_{I,\tau} = \max_{\theta_1, \theta_2} \{L_J(\theta_1) + L_{J^c}(\theta_2)\} - \max_{\theta} L_I(\theta) \\ = L_J(\tilde{\theta}_J) + L_{J^c}(\tilde{\theta}_{J^c}) - L_I(\tilde{\theta}_I).$$

The test statistic for an unknown change point location is defined as $T_I = \max_{\tau \in I} T_{I,\tau}$. The change point test compares this test statistic with a critical value \mathfrak{z}_I , which may depend on the interval I . One rejects the hypothesis of homogeneity if $T_I > \mathfrak{z}_I$.

4.2 Parameters of the LCP Procedure

To apply the LCP testing procedure for local homogeneity, we have to specify some parameters. This includes selecting interval candidates I_k or, equivalently, of the tested intervals \mathcal{J}_k and choosing respective critical values \mathfrak{z}_k . One possible parameter set that has been used successfully in simulations is presented in the following section.

4.2.1 Selection of interval candidates \mathcal{J}_k and internal points I_k . It is useful to take the set of numbers m_k defining the length of I_k and \mathcal{J}_k in the form of a geometric grid. We fix the

value m_0 and define $m_k = \lfloor m_0 c^k \rfloor$ for $k = 1, 2, \dots, K$ and $c > 1$ where $\lfloor x \rfloor$ means the integer part of x . We set $I_k = [t_0 - m_k, t_0]$ and $\mathcal{J}_k = [t_0 - m_k, t_0 - m_{k-1}]$ for $k = 1, 2, \dots, K$ (see Fig. 2).

4.2.2 Choice of the critical values \mathfrak{z}_k . The algorithm is in fact a multiple testing procedure. Mercurio and Spokoiny (2004) suggested selecting the critical value z_k to provide the overall first type error probability of rejecting the hypothesis of homogeneity in the homogeneous situation. Here we follow another proposal from Spokoiny and Chen (2007), which focuses on estimation losses caused by the “false alarm”—in our case obtaining a homogeneity interval that is too small—rather than on its probability.

In the homogeneous situation with $\theta_t \equiv \theta^*$ for all $t \in I_{k+1}$, the desirable behavior of the procedure is that after the first k steps the selected interval \hat{I}_k coincides with I_k and the corresponding estimate $\hat{\theta}_k$ coincides with $\tilde{\theta}_k$, which means there is no false alarm. On the contrary, in the case of a false alarm, the selected interval \hat{I}_k is smaller than I_k and, hence, the corresponding estimate $\hat{\theta}_k$ has larger variability than $\tilde{\theta}_k$. This means that the false alarm during the early steps of the procedure is more critical than during the final steps, because it may lead to selecting an estimate with very high variance. The difference between $\hat{\theta}_k$ and $\tilde{\theta}_k$ can naturally be measured by the value $L_{I_k}(\tilde{\theta}_k, \hat{\theta}_k) = L_{I_k}(\tilde{\theta}_k) - L_{I_k}(\hat{\theta}_k)$ normalized by the risk of the nonadaptive estimate $\tilde{\theta}_k$, $\mathfrak{R}(\theta^*) = \max_{k \geq 1} E_{\theta^*} |L_{I_k}(\tilde{\theta}_k, \theta^*)|^{1/2}$. The conditions we impose read as

$$E_{\theta^*} |L_{I_k}(\tilde{\theta}_k, \hat{\theta}_k)|^{1/2} \leq \rho \mathfrak{R}(\theta^*), \quad k = 1, \dots, K, \quad \theta^* \in \Theta. \quad (4.2)$$

The critical values \mathfrak{z}_k are selected as minimal values providing these constraints. In total we have K conditions to select K critical values $\mathfrak{z}_1, \dots, \mathfrak{z}_K$. The values \mathfrak{z}_k can be selected sequentially by Monte Carlo simulation, where one simulates under $H_0 : \theta_t = \theta^*, \forall t \in I_k$. The parameter ρ defines how conservative the procedure is. A small ρ value leads to larger critical values and hence to a conservative and nonsensitive procedure, whereas an increase in ρ results in more sensitivity at cost of stability. For details, see Spokoiny and Chen (2007) or Spokoiny (2008).

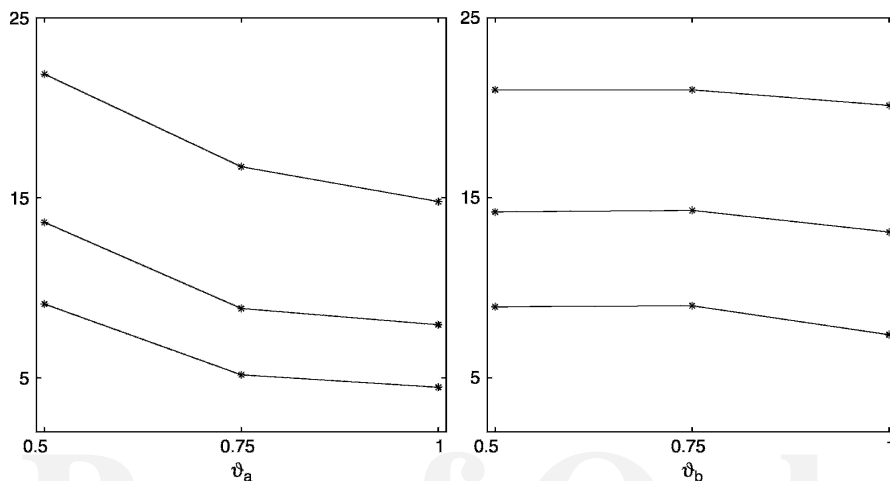


Figure 6. Mean detection delay and parameter jumps. Mean detection delays (dots) at rule $r = 0.75, 0.50$, and 0.25 from top to bottom. Left: $\vartheta_b = 0.10$ (upward jump). Right: $\vartheta_a = 0.10$ (downward jump), based on 100 simulations from Clayton copula, $m_0 = 20$, $c = 1.25$, and $\rho = 0.5$.

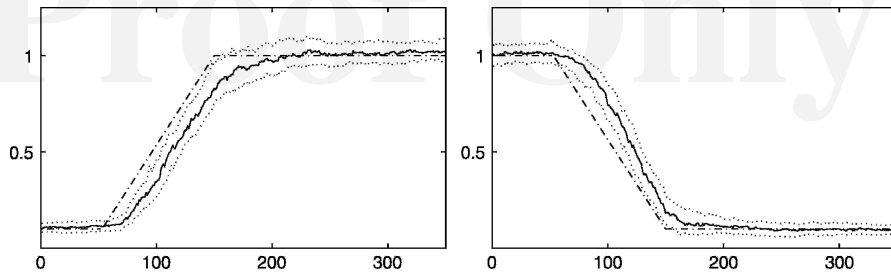


Figure 7. LCP and smooth change in copula parameter. Pointwise median (full), 0.25 and 0.75 quantiles (dotted) from $\hat{\theta}_t$ and true parameter θ_t (dashed) with $\vartheta_a = 0.10$ and $\vartheta_b = 1.00$ (left), and $\vartheta_a = 1.00$ and $\vartheta_b = 0.10$ (right). Based on 100 simulations from Clayton copula, estimated with LCP, $m_0 = 20$, $c = 1.25$, and $\rho = 0.5$.

5. SIMULATED EXAMPLES

In this section we apply the LCP procedure on simulated data with a dependence structure given by the Clayton copula. We generate sets of six-dimensional data with a sudden jump in the dependence parameter given by

$$\theta_t = \begin{cases} \vartheta_a & \text{if } -390 \leq t \leq 10 \\ \vartheta_b & \text{if } 10 < t \leq 210 \end{cases}$$

for different values of $(\vartheta_a, \vartheta_b)$: One of them is fixed at .1 (close to independence) and the other is set to larger values.

The LCP procedure is implemented with the family of interval candidates in form of a geometric grid defined by $m_0 = 20$ and $c = 1.25$. The critical values, selected according to Equation (4.2) for different ρ and θ^* , are displayed in Table 1. The choice of θ^* has negligible influence in the critical values for fixed ρ , therefore we use $\mathfrak{z}_1, \dots, \mathfrak{z}_K$ obtained with $\theta^* = 1.0$. Based on our experience, see Spokoiny and Chen (2007) and Spokoiny (2008), the default choice for ρ is 0.5.

Figure 4 shows the pointwise median and quantiles of the estimated parameter $\hat{\theta}_t$ for distinct values of $(\vartheta_a, \vartheta_b)$ based on 100 simulations. The detection delay δ at rule $r \in [0, 1]$ to jump of size $\gamma = \theta_t - \theta_{t-1}$ at t is expressed by

$$\delta(t, \gamma, r) = \min\{u \geq t : \hat{\theta}_u = \theta_{t-1} + r\gamma\} - t \quad (5.1)$$

and represents the number of steps necessary for the estimated parameter to reach the r fraction of a jump in the true parameter.

Detection delays are proportional to the probability of error of type II (i.e., the probability of accepting homogeneity in case of a jump). Thus, tests with higher power correspond to lower delays δ . Moreover, because the Kullback-Leibler divergences for upward and downward jumps are proportional to the power of the respective homogeneity tests, larger divergences result in faster jump detections.

The descriptive statistics for detection delays to jumps at $t = 11$ for different values of $(\vartheta_a, \vartheta_b)$ are in Table 2. The mean detection delay decreases with $\gamma = \vartheta_b - \vartheta_a$ and are higher for downward jumps than for upward jumps. Figure 5 shows that for Clayton copulae the Kullback-Leibler divergence is higher for upward jumps than for downward jumps. Figure 6 displays the mean detection delays against jump size for upward and downward jumps.

The LCP procedure is also applied on simulated data with smooth transition in the dependence parameter given by

$$\theta_t = \begin{cases} \vartheta_a & \text{if } -350 \leq t \leq 50 \\ \vartheta_a + \frac{t-50}{100}(\vartheta_b - \vartheta_a) & \text{if } 50 < t \leq 150 \\ \vartheta_b & \text{if } 150 < t \leq 350. \end{cases}$$

Figure 7 depicts the pointwise median and quantiles of the estimated parameter $\hat{\theta}_t$ and the true parameter θ_t for $(\vartheta_a, \vartheta_b)$ set to $(0.10, 1.00)$ and $(1.00, 0.10)$.

6. EMPIRICAL RESULTS

In this section the var from German stock portfolios is estimated based on time-varying copulae and *RiskMetrics* approaches. The time-varying copula parameters are selected by local change point (LCP) and moving window procedures. Backtesting is used to evaluate the performances of the three methods in var estimation.

Two groups of six stocks listed on DAX are used to compose the portfolios. Stocks from group 1 belong to three different industries: automotive (Volkswagen and DaimlerChrysler), insurance (Allianz and Münchener Rückversicherung), and chemical (Bayer and BASF). Group 2 is composed of stocks from six industries: electrical (Siemens), energy (E.ON), metallurgical (ThyssenKrupp), airlines (Lufthansa), pharmaceutical (Schering), and chemical (Henkel). The portfolio values are calculated using 1,270 observations, from January 1, 2000 to December 31, 2004, of the daily stock prices (data available at <http://sfb649.wiwi.hu-berlin.de/fedc>).

The selected copula belongs to the Clayton family (Eq. 2.3). Clayton copulae have a natural interpretation and are well advocated in risk management applications. In line with the stylized facts for financial returns, Clayton copulae are asymmetric and present lower tail dependence, modeling joint

Table 3. p Values from tests on residuals $\hat{\varepsilon}_{t,j}$

j	Ljung-Box		ARCH	
	Group 1	Group 2	Group 1	Group 2
1	0.33	0.52	0.15	0.04
2	0.13	0.35	0.15	0.98
3	0.21	0.08	0.34	0.72
4	0.99	0.05	0.10	0.18
5	0.90	0.07	0.91	0.77
6	0.28	0.81	0.28	0.94

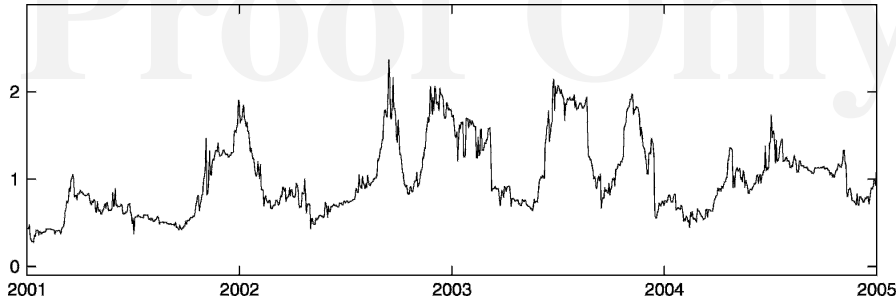


Figure 8. Time-varying dependence, group 1. Copula parameter $\hat{\theta}_t$ estimated with LCP method, Clayton copula, $m_0 = 20$, $c = 1.25$, and $\rho = 0.5$.

extreme events at lower orthants with higher probability than Gaussian copulae for the same correlation, see McNeil, Frey, and Embrechts (2005). This fact is essential for var calculations and is illustrated by the ratio between Equations (2.2) and (2.3) for off-diagonal elements of Ψ set to 0.25 and $\theta = 0.5$. For the quantiles $u_i = 0.05$, $i = 1, \dots, 6$ the ratio $C_{\Psi}^{Ga}(u_1, \dots, u_6) / C_{\theta}(u_1, \dots, u_6)$ equals 2.3×10^{-2} , whereas for the 0.01 quantiles it equals 1.3×10^{-3} .

The var estimation follows the steps described in Section 3. Using the *RiskMetrics* approach, the log-returns X_t are assumed conditionally normal distributed with zero mean and covariance matrix following a GARCH specification with fixed decay factor $\lambda = 0.94$ as in Equation (3.4).

In the time-varying copulae estimation, the log-returns are modeled as in Equation (3.5), where the innovations ε_t have cdf $F_{t,\varepsilon_i}(x_1, \dots, x_d) = C_{\theta_i}\{F_{t,1}(x_1), \dots, F_{t,d}(x_d)\}$ and C_{θ} is the Clayton copula. The univariate log-returns $X_{t,j}$ corresponding to stock j are devolatilized according to *RiskMetrics* (i.e., with zero conditional means and conditional variances $\sigma_{t,j}^2$ estimated by the univariate version of Equation (3.4) with a decay factor equal to 0.94). We note that this choice sets the same specification for the dynamics of the univariate returns across all methods (*RiskMetrics*, moving windows, and LCP), making their performances in var estimation comparable. Moreover, as the means from daily returns are clearly dominated by the variances and are approximately independent on the available information sets (see Jorion 1995; Fleming, Kirby, and Ostdiek 2001; and Christoffersen and Diebold 2006), their specification is very unlikely to cause a perceptible bias in the estimated variances and dependence parameters. Therefore, the zero mean assumption is, as pointed out by Kim, Malz, and Mina (1999), as good as any other choice. Daily returns are also modeled with zero conditional means in Fan and Gu (2003) and Härdle, Herwartz, and Spokoiny (2003) among others.

The GARCH specification (Eq. 3.4) with $\lambda = .94$ optimizes variance forecasts across a large number of assets (Morgan 1996), and is widely used in the financial industry. Different choices for the decay factor (like 0.85 or 0.98) result in negligible changes (about 3%) in the estimated dependence parameter.

The p values from the Ljung-Box test for serial correlation and from ARCH test for heteroscedasticity effects in the obtained residuals $\hat{\varepsilon}_{t,j}$ are in Table 3. Normality is rejected by Jarque-Bera test, with p values approximately 0.00 for all residuals in both groups. The empirical cdf's of residuals as defined in Equation (2.4) are used for the copula estimation.

With the moving windows approach, the size of the estimating window is fixed as 250 days corresponding to 1 business year (the same size is used in, for example, Fan and Gu (2003)). For the LCP procedure, following Section 4.2, we set the family of interval candidates as a geometric grid with $m_0 = 20$, $c = 1.25$, and $\rho = 0.5$. We have chosen these parameters from our experience in simulations (for details on robustness of the reported results with respect to the choice of m_0 and c , refer to Spokoiny (2008)).

The performance of the var estimation is evaluated based on backtesting. At each time t , the estimated var at level α for a portfolio w is compared with the realization l_t of the corresponding P&L function (see Eq. 3.2), with an exceedance occurring for each l_t less than $\widehat{\text{var}}_t(\alpha)$. The ratio of the number of exceedances to the number of observations gives the exceedance ratio

$$\hat{\alpha}_w(\alpha) = \frac{1}{T} \sum_{t=1}^T \mathbf{1}_{\{l_t < \widehat{\text{var}}_t(\alpha)\}}$$

Because the first 250 observations are used for estimation, $T = 1,020$. The difference between $\hat{\alpha}$ and the desired level α is expressed by the relative exceedance error

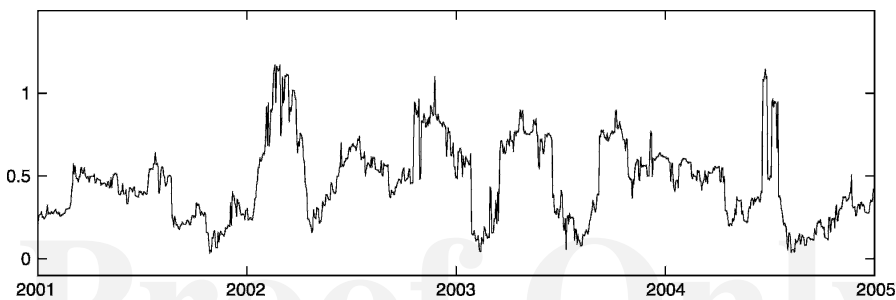


Figure 9. Time-varying dependence, group 2. Copula parameter $\hat{\theta}_t$ estimated with LCP method, Clayton copula, $m_0 = 20$, $c = 1.25$, and $\rho = 0.5$.

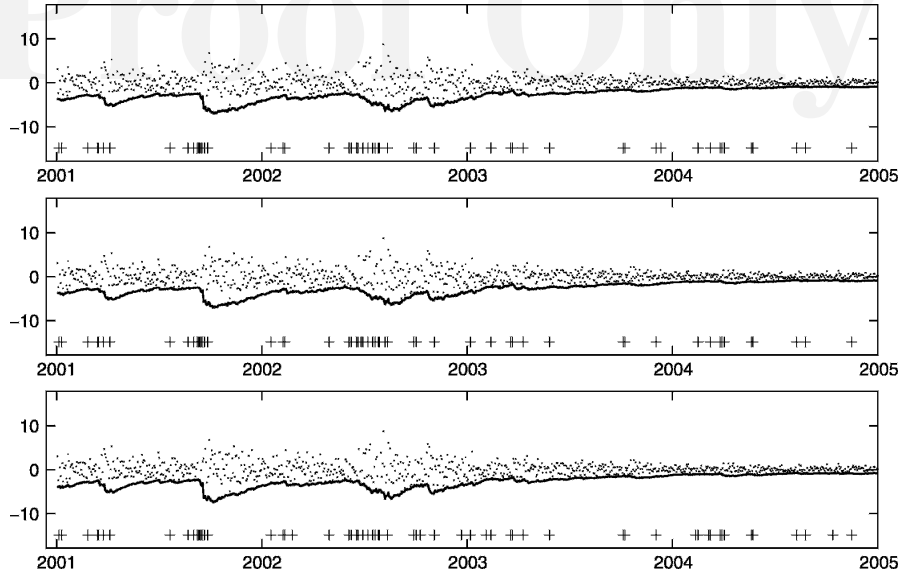


Figure 10. Estimated var across methods, group 1. P&L realizations l_t (dots), $\widehat{\text{var}}_t(\alpha)$ (line), and exceedance times (crosses). Estimated with LCP (top), moving windows (middle), and *RiskMetrics* (bottom) for equally weighted portfolio w^* at level $\alpha = 0.05$.

$$e_w = (\hat{\alpha} - \alpha)/\alpha.$$

We compute exceedance ratios and relative exceedance errors to levels $\alpha = 0.05$ and 0.01 for a set $W = \{w^*, w_n; n = 1, \dots, 100\}$ of portfolios, where each $w_n = (w_{n,1}, \dots, w_{n,6})^\top$ is a realization of a random vector uniformly distributed on $S = \{(x_1, \dots, x_6) \in \mathbb{R}^6 : \sum_{i=1}^6 x_i = 1, x_i \geq .1\}$, and $w^* = 1/6 \mathbf{I}_6$, with \mathbf{I}_d denoting the $(d \times 1)$ vector of ones, is the equally weighted portfolio. The degree of diversification of a portfolio can be measured based on the majorization preordering on S (see Marshall and Olkin 1979). In other words, a portfolio w_a is more diversified than portfolio w_b if $w_a \prec w_b$. Under the majorization preordering the vector w^* satisfies $w^* \prec w$ for all $w \in S$; therefore, the equally weighted portfolio is the most diversified portfolio from W , see Ibragimov and Walden (2007).

The average relative exceedance error over portfolios and the corresponding standard deviation

$$A_W = \frac{1}{|W|} \sum_{w \in W} e_w$$

$$D_W = \left\{ \frac{1}{|W|} \sum_{w \in W} (e_w - A_W)^2 \right\}^{\frac{1}{2}}$$

are used to evaluate the performances of the time-varying copulae and *RiskMetrics* methods in var estimation.

The dependence parameter estimated with LCP for stocks from groups 1 and 2 are shown in Figures 8 and 9. The different industry concentrations in each group are reflected in the higher parameter values obtained for group 1. The P&L and the var at level 0.05 estimated with LCP, moving windows, and

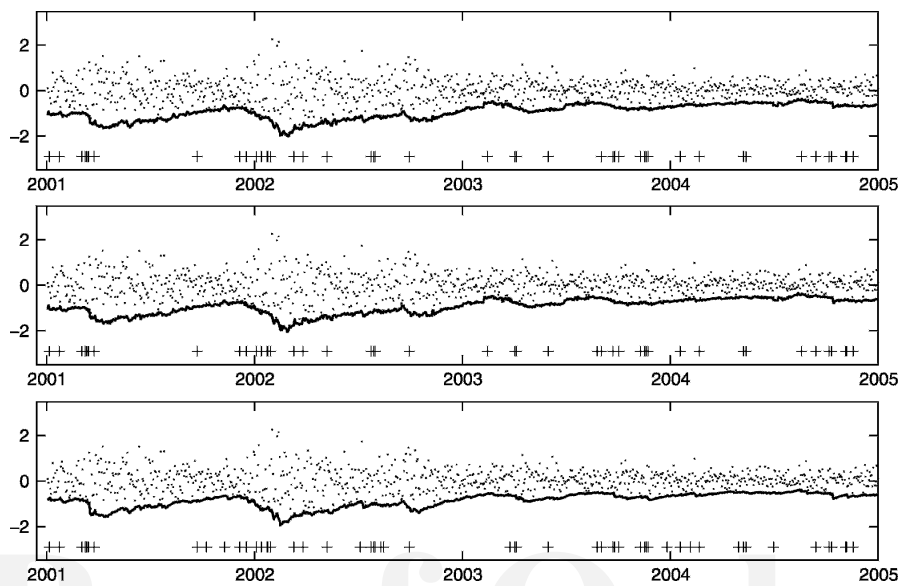


Figure 11. Estimated var across methods, group 2. P&L realizations l_t (dots), $\widehat{\text{var}}_t(\alpha)$ (line), and exceedance times (crosses). Estimated with LCP (top), moving windows (middle), and *RiskMetrics* (bottom) for equally weighted portfolio w^* at level $\alpha = 0.05$.

Table 4. Exceedance ratios and errors, group 1

	RiskMetrics		Moving windows		LCP	
	$\alpha = 5.00$	$\alpha = 1.00$	$\alpha = 5.00$	$\alpha = 1.00$	$\alpha = 5.00$	$\alpha = 1.00$
	$\hat{\alpha}_{w^*}$	6.11	1.48	5.62	0.59	5.52
$\hat{\alpha}_{w_1}$	5.91	1.38	5.42	0.49	5.42	0.69
$\hat{\alpha}_{w_2}$	6.40	1.28	5.91	0.49	5.71	0.59
A_W	0.23	0.45	0.11	-0.49	0.11	-0.36
D_W	0.04	0.14	0.06	0.08	0.06	0.10

NOTE: Exceedance ratios for portfolios w^* , w_1 , and w_2 , and average and standard deviation from relative exceedance errors. Across levels and methods, ratios and levels are expressed as a percentage.

RiskMetrics methods for the equally weighted portfolio w^* are in Figures 10 (group 1) and 11 (group 2). Exceedance ratios for portfolios w^* , w_1 , and w_2 ; average relative exceedance errors; and corresponding standard deviations across methods and levels are shown in Tables 4 (group 1) and 5 (group 2).

Based on the exceedance errors, the LCP procedure outperforms the moving windows (second best) and *RiskMetrics* methods in var estimation in group 1. At level 0.05, the average error associated with copula methods is about half the error from *RiskMetrics* estimation for nearly the same standard deviation. At level 0.01, the LCP average error is the smallest in absolute value, and copula methods present less standard deviations. At this level, copula methods overestimate var, and *RiskMetrics* underestimates it. Although overestimation of var means that a financial institution would be requested to keep more capital aside than necessary to guarantee the desired confidence level, underestimation means that less capital is reserved and the desired level is not guaranteed. Therefore, from the regulatory point of view, overestimation is preferred to underestimation. In the less concentrated group 2, LCP outperforms moving windows and *RiskMetrics* at the level 0.05, presenting the smallest average error in magnitude for nearly the same value of D_W . At level 0.01, copula methods overestimate and *RiskMetrics* underestimates the var by about 60%.

It is interesting to note the effect of portfolio diversification on the exceedance errors for group 1 and level 0.01. The errors decrease with increasing portfolio diversification for copulae methods but become larger under the *RiskMetrics* estimation. For other groups and levels, the diversification effects are not clear. Refer to Ibragimov (2007) and Ibragimov and Walden

(2007) for details on the effects of portfolio diversification under heavy-tailed distributions in risk management.

7. CONCLUSION

In this article we modeled the dependence structure from German equity returns using time-varying copulae with adaptively estimated parameters. In contrast to Patton (2006) and Rodriguez (2007), we neither specified the dynamics nor assumed regime switching models for the copula parameter. The parameter choice was performed under the local homogeneity assumption with homogeneity intervals recovered from the data through local change point analysis.

We used time-varying Clayton copulae, which are asymmetric and present lower tail dependence, to estimate the var from portfolios of two groups of German securities, presenting different levels of industry concentration. *RiskMetrics*, a widely used methodology based on multivariate normal distributions, was chosen as a benchmark for comparison. Based on back-testing, the adaptive copula achieved the best var estimation performance in both groups, with average exceedance errors mostly small in magnitude and corresponding to sufficient capital reserve for covering losses at the desired levels.

The better var estimates provided by Clayton copulae indicate that the dependence structure from German equities may contain nonlinearities and asymmetries, such as stronger dependence at lower tails than at upper tails, that cannot be captured by the multivariate normal distribution. This asymmetry translates into extremely negative returns being more correlated than extremely positive returns. Thus, our results for the German equities resemble those from Longin and Solnik (2001), Ang and Chen (2002) and Patton (2006) for international markets, U.S. equities, and Deutsch mark/Japanese yen exchange rates, where empirical evidence for asymmetric dependences with increasing correlations in market downturns were found.

Furthermore, in the non-Gaussian framework, with nonlinearities and asymmetries taken into consideration through the use of Clayton copulae, the adaptive estimation produces better var fits than the moving window estimation. The high sensitive adaptive procedure can capture local changes in the dependence parameter that are not detected by the estimation with a scrolling window of fixed size.

The main advantage of using time-varying copulae to model dependence dynamics is that the normality assumption is not needed. With the proposed adaptively estimated time-varying copulae, neither normality assumption nor specification for the dependence dynamics are necessary. Hence, the method provides more flexibility in modeling dependences between markets and economies over time.

ACKNOWLEDGMENTS

Financial support from the *Deutsche Forschungsgemeinschaft* via *SFB 649 "Ökonomisches Risiko,"* Humboldt-Universität zu Berlin is gratefully acknowledged. The authors also thank the editor, an associate editor, and two referees for their helpful comments.

[Received October 2006. Revised November 2007.]

Table 5. Exceedance ratios and errors, group 2

	RiskMetrics		Moving windows		LCP	
	$\alpha = 5.00$	$\alpha = 1.00$	$\alpha = 5.00$	$\alpha = 1.00$	$\alpha = 5.00$	$\alpha = 1.00$
	$\hat{\alpha}_{w^*}$	5.42	1.58	4.53	0.39	4.53
$\hat{\alpha}_{w_1}$	5.81	1.77	5.02	0.39	5.02	0.39
$\hat{\alpha}_{w_2}$	5.62	1.58	5.12	0.39	5.22	0.30
A_W	0.16	0.57	-0.10	-0.65	-0.09	-0.65
D_W	0.04	0.16	0.06	0.09	0.06	0.08

NOTE: Exceedance ratios for portfolios w^* , w_1 , and w_2 , and average and standard deviation from relative exceedance errors. Across levels and methods, ratios and levels are expressed as a percentage.

AU8

F10, T1

T4, T5

AU9

REFERENCES

- Andrews, D. W. K. (1993), "Tests for Parameter Instability and Structural Change With Unknown Change Point," *Econometrica*, 61, 821–856.
- Andrews, D. W. K., and Ploberger, W. (1994), "Optimal Tests When a Nuisance Parameter Is Present Only Under the Alternative," *Econometrica*, 62, 1383–1414.
- Ang, A., and Chen, J. (2002), "Asymmetric Correlations of Equity Portfolios," *Journal of Financial Economics*, 63, 443–494.
- Belomestny, D., and Spokoiny, V. (2007), "Spatial Aggregation of Local Likelihood Estimates With Applications to Classification," *The Annals of Statistics*, 35, 2287–2311.
- Chen, X., and Fan, Y. (2006), "Estimation and Model Selection of Semiparametric Copula-Based Multivariate Dynamic Models Under Copula Misspecification," *Journal of Econometrics*, 135, 125–154.
- Chen, X., Fan, Y., and Tsyrennikov, V. (2006), "Efficient Estimation of Semiparametric Multivariate Copula Models," *Journal of the American Statistical Association*, 101, 1228–1240.
- Cherubini, U., Luciano, E., and Vecchiato, W. (2004), *Copula Methods in Finance*, Chichester: Wiley.
- Christoffersen, P., and Diebold, F. (2006), "Financial Asset Returns, Direction-of-Change Forecasting, and Volatility Dynamics," *Management Science*, 52, 1273–1287.
- Cont, R. (2001), "Empirical Properties of Asset Returns: Stylized Facts and Statistical Issues," *Quantitative Finance*, 1, 223–236.
- Embrechts, P., Hoeing, A., and Juri, A. (2003a), "Using Copulae to Bound the Value-at-Risk for Functions of Dependent Risks," *Finance and Stochastics*, 7, 145–167.
- Embrechts, P., Lidskog, F., and McNeil, A. (2003b), "Modelling Dependence with Copulas and Applications to Risk Management," in *Handbook of Heavy Tailed Distributions in Finance*, ed. S. Rachev, Amsterdam: North-Holland, pp. 329–384.
- Embrechts, P., McNeil, A., and Straumann, D. (2002), "Correlation and Dependence in Risk Management: Properties and Pitfalls," in *Risk Management: Value at Risk and Beyond*, ed. M. Dempster, Cambridge, UK: Cambridge University Press.
- Fan, J., and Gu, J. (2003), "Semiparametric Estimation of Value-at-Risk," *The Econometrics Journal*, 6, 261–290.
- Fleming, J., Kirby, C., and Ostdiek, B. (2001), "The Economic Value of Volatility Timing," *The Journal of Finance*, 56, 239–354.
- Franke, J., Härdle, W., and Hafner, C. (2004), *Statistics of Financial Markets*, Heidelberg: Springer-Verlag.
- Fréchet, M. (1951), "Sur les Tableaux de Correlation Dont les Marges Sont Données," *Annales de l'Université de Lyon, Sciences Mathématiques et Astronomie*, 14, 5–77.
- Giacomini, E., and Härdle, W. (2005), "Value-at-Risk Calculations With Time Varying Copulae," in *Bulletin of the International Statistical Institute, Proceedings of the 55th Session*.
- Granger, C. (2003), "Time Series Concept for Conditional Distributions," *Oxford Bulletin of Economics and Statistics*, 65, 689–701.
- Hansen, B. E. (2001), "The New Econometrics of Structural Change: Dating Breaks in U.S. Labor Productivity," *The Journal of Economic Perspectives*, 15, 117–128.
- Härdle, W., Herwartz, H., and Spokoiny, V. (2003), "Time Inhomogeneous Multiple Volatility Modelling," *Journal of Financial Econometrics*, 1, 55–95.
- Härdle, W., Kleinow, T., and Stahl, G. (2002), *Applied Quantitative Finance*, Springer-Verlag, Heidelberg.
- Hoeffding, W. (1940), "Maßstabvariante Korrelationstheorie," *Schriften des mathematischen Seminars und des Instituts für angewandte Mathematik der Universität Berlin*, 5, 181–233.
- Hu, L. (2006), "Dependence Patterns Across Financial Markets: A Mixed Copula Approach," *Applied Financial Economics*, 16, 717–729.
- Ibragimov, R. (2007), "Efficiency of Linear Estimators Under Heavy-Tailedness: Convolutions of α -Symmetric Distributions," *Econometric Theory*, 23, 501–517.
- Ibragimov, R., and Walden, J. (2007), "The Limits of Diversification When Losses May be Large," *Journal of Banking and Finance*, 31, 2551–2569.
- Joe, H. (1997), *Multivariate Models and Dependence Concepts*, London: Chapman & Hall.
- Jorion, P. (1995), "Predicting Volatility in the Foreign Exchange Market," *The Journal of Finance*, 50, 507–528.
- Morgan, J. P. (1996), *RiskMetrics Technical Document*, New York: RiskMetrics Group.
- Kim, J., Malz, A. M., and Mina, J. (1999), *Long Run Technical Document*, New York: RiskMetrics Group.
- Longin, F., and Solnik, B. (2001), "Extreme Correlation on International Equity Markets," *The Journal of Finance*, 56, 649–676.
- Mari, D., and Kotz, S. (2001), *Correlation and Dependence*, London: Imperial College Press.
- Marshall, A., and Olkin, I. (1979), *Inequalities: Theory of Majorizations and Its Applications*, New York: Academic Press.
- McNeil, A. J., Frey, R., and Embrechts, P. (2005), *Quantitative Risk Management: Concepts, Techniques and Tools*, Princeton, NJ: Princeton University Press.
- Mercurio, D., and Spokoiny, V. (2004), "Estimation of Time Dependent Volatility via Local Change Point Analysis With Applications to Value-at-Risk," *Annals of Statistics*, 32, 577–602.
- Nelsen, R. (1998), *An Introduction to Copulas*, New York: Springer-Verlag.
- Patton, A. (2004), "On the Out-of-Sample Importance of Skewness and Asymmetric Dependence for Asset Allocation," *Journal of Financial Econometrics*, 2, 130–168.
- (2006), "Modelling Asymmetric Exchange Rate Dependence," *International Economic Review*, 47, 527–556.
- Perron, P. (1989), "The Great Crash, the Oil Price Shock and the Unit Root Hypothesis," *Econometrica*, 57, 1361–1401.
- Polzehl, J., and Spokoiny, V. (2006), "Propagation–Separation Approach for Likelihood Estimation," *Probability Theory and Related Fields*, 135, 335–362.
- Quintos, C., Fan, Z., and Philips, P. C. B. (2001), "Structural Change Tests in Tail Behaviour and the Asian Crisis," *The Review of Economic Studies*, 68, 633–663.
- Rodriguez, J. C. (2007), "Measuring Financial Contagion: A Copula Approach," *Journal of Empirical Finance*, 14, 401–423.
- Sklar, A. (1959), "Fonctions de Répartition à n Dimensions et Leurs Marges," *Publications de l'Institut de Statistique de l'Université de Paris*, 8, 229–231.
- Spokoiny, V. (2008), *Local Parametric Methods in Nonparametric Estimation*, Berlin, Heidelberg: Springer-Verlag.
- Spokoiny, V., and Chen, Y. (2007), *Multiscale Local Change Point Detection with Applications to Value-at-Risk*, Preprint 904, Berlin: Weierstrass Institute Berlin.
- Stock, J.H. (1994), "Unit Roots, Structural Breaks and Trends," in *Handbook of Econometrics*, Vol. 4, ed. R. F. Engle and D. McFadden, Amsterdam: North-Holland, pp. 2739–2841.
- Zivot, E., and Andrews, D. W. K. (1992), "Further Evidence on the Great Crash, the Oil Price Shock and the Unit Root Hypothesis," *Journal of Business & Economic Statistics*, 10, 251–270.



Dynamics of state price densities

Wolfgang Härdle^a, Zdeněk Hlávka^{b,*}

^a CASE—Center for Applied Statistics and Economics, Wirtschaftswissenschaftliche Fakultät, Humboldt-Universität zu Berlin, Spandauer Str. 1, 10178 Berlin, Germany

^b Charles University in Prague, Department of Statistics, Sokolovská 83, 18675 Praha, Czech Republic

ARTICLE INFO

Article history:

Received 11 January 2009

Accepted 12 January 2009

Available online 15 January 2009

JEL classification:

C13

C14

G13

Keywords:

Option pricing

State price density

Nonlinear least squares

Constrained estimation

ABSTRACT

State price densities (SPDs) are an important element in applied quantitative finance. In a Black–Scholes world they are lognormal distributions, but in practice volatility changes and the distribution deviates from log-normality. In order to study the degree of this deviation, we estimate SPDs using EUREX option data on the DAX index via a nonparametric estimator of the second derivative of the (European) call pricing function. The estimator is constrained so as to satisfy no-arbitrage constraints and corrects for the intraday covariance structure in option prices. In contrast to existing methods, we do not use any parametric or smoothness assumptions.

© 2009 Elsevier B.V. All rights reserved.

1. Introduction

The dynamics of option prices carries information on changes in state price densities (SPDs). The SPD contains important information on the behavior and expectations of the market and is used for pricing and hedging. The most important application of an SPD is that it allows one to price options with complicated payoff functions simply by (numerical) integration of the payoff with respect to this density.

Prices $C_t(K, T)$ of European options with strike price K observed at time t and expiring at time T allow one to deduce the state price density $f(\cdot)$ using the relationship (Breedon and Litzenberger, 1978)

$$f(K) = \exp\{r(T-t)\} \frac{\partial^2 C_t(K, T)}{\partial K^2}. \quad (1)$$

Eq. (1) can be used to estimate the SPD $f(K)$ from the observed option prices. An extensive overview of parametric and other estimation techniques can be found, for example, in Jackwerth (1999). An application to option pricing is given in Buehler (2006).

Kernel smoothers were in this framework proposed and successfully applied by, for example, Ait-Sahalia and Lo (1998), Ait-Sahalia and Lo (2000), Ait-Sahalia et al. (2000), or Huynh et al. (2002). Ait-Sahalia and Duarte (2003) proposed a method for

nonparametric estimation of the SPD under constraints like positivity, convexity, and boundedness of the first derivative. Bondarenko (2003) calculates arbitrage-free SPD estimates using positive convolution approximation (PCA) methodology and demonstrates its properties in a Monte Carlo study based on closing prices of the S&P 500 options. Another sophisticated approach based on smoothing splines allowing one to include these constraints is described and applied on simulated data in Yatchew and Härdle (2006). In the majority of these papers, the focus was more on the smoothing techniques rather than on a no-arbitrage argument, although a crucial element of local volatility models is the absence of arbitrage (Dupire, 1994). Highly numerically efficient pricing algorithms, for example, by Andersen and Brotherton-Ratcliffe (1997), rely heavily on no-arbitrage properties. Kahalé (2004) proposed a procedure that requires solving a set of nonlinear equations with no guarantee of a unique solution. Moreover, for that algorithm the data feed is already (unrealistically) expected to be arbitrage free (Fengler, 2005; Fengler et al., 2007). In addition, the covariance structure of the quoted option prices (Renault, 1997) is rarely incorporated into the estimation procedure.

In Table 1, we give an overview of selected properties of different estimation techniques. The parametric approach may be used to estimate parameters of a probability density lying in some preselected family. The parametric models may be further extended by considering more flexible probability densities or mixtures of distributions. Approaches based on nonparametric smoothing techniques are more flexible since the shape of a nonparametric SPD estimate is not fixed in advance and the method controls only the smoothness of the estimate. For example,

* Corresponding author. Tel.: +420 221 913 284; fax: +420 283 073 341.

E-mail addresses: haerdle@wiwi.hu-berlin.de (W. Härdle), hlavka@karlin.mff.cuni.cz (Z. Hlávka).

Table 1
Summary of properties of parametric and nonparametric estimators.

	Methods			
	Parametric	Standard smoothing method	Nonparametric under constraints	This paper
Shape	Fixed	Flexible	Flexible	Flexible
Control	Choice of family	Smoothness	Smoothness	None
SPD support	Infinite	Restricted	Restricted	Restricted
Constraints	By design	Local	Yes	Yes

the smoothness of a kernel regression estimator depends mostly on the choice of the bandwidth parameter, the smoothness of the PCA estimator (Bondarenko, 2003) depends on the choice of the kernel, and the smoothness of the NNLS estimator (Yatchew and Härdle, 2006) is controlled by constraining the Sobolev norm of the SPD; using these nonparametric estimators, systematic bias may typically occur in the case of oversmoothing. Constraints on estimators are more easily implemented for globally valid parametric models than for local (nonparametric) models. The use of a standard smoothing technique which does not account for the constraints is not advisable. The value of the nonparametric estimate cannot be calculated in regions without any data and, therefore, the support of nonparametrically estimated SPDs is limited by the range of the observed strike prices even for nonparametric-under-constraints techniques.

Most of the commonly used estimation techniques do not specify explicitly the source of random error in the observed option prices; see Renault (1997) for an extensive review of this subject. A common approach in SPD estimation is to use either the closing option prices or to correct the intraday option prices by the current value of the underlying asset. Both approaches lack interpretation if the shape of the SPD changes rapidly. This can be made clear by a gedankenexperiment: if the shape of the SPD changes dramatically during the day, correcting the observed option prices by the value of the underlying asset and then estimating the SPD would lead to an estimate of some (nonexisting) daily average of the true SPDs. We try to circumvent this problem by introducing a simple model for the intraday covariance structure of option prices which allows us to estimate the value of the true SPD at an arbitrarily chosen fixed time; see also Hlávka and Svojík (2008). Most often, we are interested in the estimation of the current SPD.

We develop a simple estimation technique in order to construct constrained SPD estimates from the observed intraday option prices which are treated as repeated observations collected during a certain time period. The proposed technique involves constrained LS-estimation, it enables us to construct confidence intervals for the current value of the SPD and prediction intervals for its future development, and it does not depend on any tuning (smoothness) parameter. The construction of a simple approximation of the covariance structure of the observed option prices follows naturally from the derivation of our nonparametric constrained estimator. This covariance structure is interesting in itself; it separates two sources of random errors, and it is applicable to other SPD estimators.

We study the development of the estimated SPDs in Germany over 8 years. A no-arbitrage argument is imposed at each time point, leading (mathematically) to the above-mentioned no-arbitrage constraints. This, of course, is a vital feature for trading purposes where the derived (implied) volatility surfaces for different strikes and maturities are needed for proper judgment of risk and return.

The resulting SPDs and implied volatility surfaces are not smooth per se. In most applications, this is not a disadvantage though, since, first, we may smooth the resulting SPD estimates (Hlávka and Svojík, 2008) and, second, we are mostly interested in functionals of the estimated SPD like, for example, the expected payoff or the forward price. Another important feature that can be

easily estimated from the nonsmooth SPDs are the quantiles; see Section 6.2 for an application.

In Section 2, we introduce the notation, discuss constraints that are necessary for estimating SPDs, and we construct a very simple unconstrained SPD estimator using simple linear regression. In Section 3, this estimator is modified so that it satisfies the shape constraints given in Section 2.1. We demonstrate that the covariance structure of the option prices exhibits correlations depending both on the strike price and time of the trade in Section 4. In Section 5, we apply our estimation technique on option prices observed in the year 1995, and we show that the proposed approximation of the covariance structure removes the dependency and heteroscedasticity of the residuals. The dynamics of the estimated SPDs in years 1995–2003 is studied in Section 6.

2. Construction of the estimate

The fair price of a European call option with payoff $(S_T - K)_+ = \max(S_T - K, 0)$, with S_T denoting the price of the stock at time T , t the current time, K the strike price, and r the risk-free interest rate, can be written as

$$C_t(K, T) = \exp\{-r(T - t)\} \int_0^\infty (S_T - K)_+ f(S_T) dS_T, \tag{2}$$

i.e., as the discounted expected value of the payoff with respect to the SPD $f(\cdot)$. For the sake of simplicity of the following presentation, we assume in the rest of the paper that the discount factor $\exp\{-r(T - t)\} = 1$. In applications, this is achieved by correcting the observed option prices by the known risk-free interest rate r and the time to maturity $(T - t)$ in (2). At the time of the trade, the current index price and volatility are common to all options and, hence, do not appear explicitly in Eq. (2).

Let us denote the i -th observation of the strike price by K_i and the corresponding option price, divided by the discount factor $\exp\{-r(T - t)\}$ from (2), by $C_i = C_{t,i}(K_i, T)$. In practice, on any given day t , one observes option prices repeatedly for a small number of distinct strike prices. Therefore, it is useful to adopt the following notation. Let $\mathcal{C} = (C_1, \dots, C_n)^\top$ be the vector of the observed option prices on day t sorted by strike price. Then, the vector of strike prices has the following structure:

$$\mathcal{K} = \begin{pmatrix} K_1 \\ K_2 \\ \vdots \\ K_n \end{pmatrix} = \begin{pmatrix} k_1 \mathbf{1}_{n_1} \\ k_2 \mathbf{1}_{n_2} \\ \vdots \\ k_p \mathbf{1}_{n_p} \end{pmatrix},$$

where $k_1 < k_2 < \dots < k_p$, $n_j = \sum_{i=1}^n \mathbf{I}(K_i = k_j)$, with $\mathbf{I}(\cdot)$ denoting the indicator function and $\mathbf{1}_n$ a vector of ones of length n .

2.1. Assumptions and constraints

Let us now concentrate on options corresponding to a single maturity T observed at fixed time t . Let us assume that the i -th observed option price (corresponding to strike price K_i) follows the model

$$C_{t,i}(K_i, T) = \mu(K_i) + \varepsilon_i, \tag{3}$$

where ε_i are iid random variables with zero mean and variance σ^2 . In practice, one might expect that the errors exhibit correlations depending on the strike price and time. Heteroscedasticity can

be incorporated in model (3) if we assume that the random errors ε_i have variance $\text{Var } \varepsilon_i = \sigma_{k_i}^2$, leading to weighted least squares. The assumptions on the distribution of random errors will be investigated in more detail in Section 5.3. Following Renault (1997), we interpret the observed option price as the price given by a pricing formula plus an error term, and in Section 4 we suggest a covariance structure for the observed option prices taking into account the dependencies across strike prices and times of trade.

Harrison and Pliska (1981) characterized the absence of arbitrage by the existence of a unique risk neutral SPD $f(\cdot)$. From formula (2) and the properties of a probability density it follows that, in a continuous setting, the function $\mu(\cdot)$, defined on \mathbb{R}^+ , has to satisfy the following no-arbitrage constraints:

- 1': it is positive,
- 2': it is decreasing in K ,
- 3': it is convex,
- 4': its second derivative exists and it is a density (i.e., nonnegative and it integrates to one).

Let us now have a look at functions satisfying Constraints 1'–4'.

Lemma 1. Suppose that $\mu : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ satisfies Constraints 1'–4'. Then the first derivative, $\mu^{(1)}(\cdot)$, is nondecreasing and such that $\lim_{x \rightarrow 0} \mu^{(1)}(x) = -1$ and $\lim_{x \rightarrow +\infty} \mu^{(1)}(x) = 0$.

Proof. Constraint 4' implies that the first derivative, $\mu^{(1)}$, exists and that it is differentiable. $\lim_{x \rightarrow +\infty} \mu^{(1)}(x)$ exists since the function $\mu^{(1)}$ is nondecreasing (Constraint 3') and bounded (Constraint 2'). Next, $\lim_{x \rightarrow \infty} \mu^{(1)}(x) = 0$ since a negative limit would violate Constraint 1' for large x ($\mu^{(1)}(x)$ cannot be positive since $\mu(x)$ is decreasing). Finally, Constraint 4', $1 = \int_0^\infty \mu^{(2)}(x) dx = \lim_{x \rightarrow +\infty} \mu^{(1)}(x) - \lim_{x \rightarrow 0} \mu^{(1)}(x)$, implies that $\lim_{x \rightarrow 0} \mu^{(1)}(x) = -1$. \square

Remark 1. Lemma 1 allows us to restate Constraints 3' and 4' in terms of $\mu^{(1)}(\cdot)$ by assuming that $\mu^{(1)}(\cdot)$ is differentiable, nondecreasing, and such that $\lim_{x \rightarrow 0} \mu^{(1)}(x) = -1$ and $\lim_{x \rightarrow +\infty} \mu^{(1)}(x) = 0$.

In this section, we stated only constraints guaranteeing that the SPD estimate will be a probability density. Constraints for the expected value of the SPD estimate are discussed in Section 3.6.

2.2. Existence and uniqueness

In this subsection we address the issue of existence and uniqueness of a regression function, $\hat{C}(\cdot)$, satisfying the required assumptions and constraints. In practice, we do not deal with a continuous function. Hence, we restate Constraints 1'–4' for discrete functions, defined only on a finite set of distinct points, say $k_1 < \dots < k_p$, in terms of their function values, $C(k_i)$, and their scaled first differences, $C_{k_i, k_j}^{(1)} = \{C(k_i) - C(k_j)\} / \{k_i - k_j\}$.

- 1: $C(k_i) \geq 0, i = 1, \dots, p$,
- 2: $k_i < k_j$ implies that $C(k_i) \geq C(k_j)$,
- 3: $k_i < k_j < k_l$ implies that $-1 \leq C_{k_i, k_j}^{(1)} \leq C_{k_j, k_l}^{(1)} \leq 0$.

It is easy to see that Constraints 1–2 are discrete versions of Constraints 1' and 2'. Constraint 3 is a discrete version of Constraints 3' and 4'; see Remark 1.

From now on, similarly as in Robertson et al. (1988), we think of the collection, \mathcal{C} , of functions satisfying Constraints 1–3 as a subset of a p -dimensional Euclidean space, where p is the number of distinct k_i 's. The constrained regression, \hat{C} , is in this setting the closest point of \mathcal{C} to the vector C of the observed option prices with distances measured by the usual Euclidean distance

$$d(f, C) = (f - C)^\top (f - C) = \sum_{i=1}^n \{f(K_i) - C(K_i)\}^2. \quad (4)$$

From this point of view, the regression function, \hat{C} , consists only of the values of the function in the points k_1, \dots, k_p . The first and second differences are used to approximate the first and the second derivatives, respectively.

We claim that the set, \mathcal{C} , of functions satisfying Constraints 1–3 is closed in the topology induced by the metric given by Euclidean distance and it is convex, i.e., if $f, g \in \mathcal{C}$ and $0 \leq a \leq 1$, then $af + (1 - a)g \in \mathcal{C}$.

Lemma 2. If $\hat{C} \in \mathcal{C}$ is the regression of $C(K_i), i = 1, \dots, n$, on $k_1 < \dots < k_p$ under Constraints 1–3 and if a and b are constants such that $a \leq C(K_i) \leq b, \forall i$, then $a \leq \hat{C}(k_i) \leq b + (k_p - k_i)$.

Proof. It is not possible that $\hat{C}(k_i)$ lies above b for all k_i 's (otherwise we would get a better fit only by shifting $\hat{C}(k_i)$). The upper bound now follows from Constraint 3.

The validity of the lower bound may be demonstrated similarly. Clearly, it is not possible that $\hat{C}(k_i)$ lie below a for all k_i 's. Moreover, it is not possible that $\hat{C}(k_1) \geq \dots \geq \hat{C}(k_i) \geq a > \hat{C}(k_{i+1}) \geq \dots \geq \hat{C}(k_p)$ for any i , since in such a situation the fit could be trivially improved by increasing $\hat{C}(k_{i+1}), \dots, \hat{C}(k_p)$ by some small amount, for example, by $a - \hat{C}(k_{i+1})$, without violating any of the Constraints 1–3. \square

Theorem 1. A regression, $\hat{C} = \arg \min_{f \in \mathcal{C}} d(f, C)$, satisfying Constraints 1–3, exists and it is unique.

Proof. Lemma 2 implies that \hat{C} belongs to a subset, \mathcal{B} , of \mathcal{C} bounded below by a and above by $b + (k_p - k_1)$. Thinking of the functions as points in Euclidean space, it is clear that the continuous function $d(f, C)$ attains its minimum on the closed and bounded set \mathcal{B} . The uniqueness of \hat{C} follows from the convexity of \mathcal{B} using, for example, Robertson et al. (1988, Theorem 1.3.1). \square

2.3. Linear model

With the given option data, Constraints 1–3 of Section 2.2 can be reformulated using linear regression models with constraints.

In the following, we fix the time t and the expiry date T and we omit these symbols from the notation. In Section 2.2 we have noted that the option prices are repeatedly observed for a small number p of distinct strike prices. Defining the expected values of the option prices for a given strike price, $\mu_j = \mu(k_j) = E\{C(k_j)\}$, we can write

$$\begin{aligned} \mu_p &= \beta_0, \\ \mu_{p-1} &= \beta_0 + \beta_1, \\ \mu_{p-2} &= \beta_0 + 2\beta_1 + \beta_2, \\ \mu_{p-3} &= \beta_0 + 3\beta_1 + 2\beta_2 + \beta_3, \\ &\vdots \\ \mu_1 &= \beta_0 + (p-1)\beta_1 + (p-2)\beta_2 + \dots + \beta_{p-1}. \end{aligned}$$

Thus, we fit our data using coefficients $\beta_j, j = 1, \dots, p$. The conditional means $\mu_i, i = 1, \dots, p$ are replaced by the same number of parameters $\beta_j, j = 0, \dots, p-1$, which allow us to impose the shape constraints in a more natural way.

The interpretation of the coefficients β_j can be seen in Fig. 1, which shows a simple situation with only four distinct strike prices ($p = 4$). β_0 is the mean option price at point 4. Constraint 1', Section 2.1, implies that it has to be positive. β_1 is the difference between the mean option prices at point 4 and point 3; Constraint 2' implies that it has to be positive. The next coefficient, β_2 , approximates the change in first derivative in point 3 and it can be interpreted as an approximation of the second derivative in point 3. Constraint 3' implies that β_2 has to be positive. Similarly, β_3 is an estimate of the (positive) second derivative in point 2. Constraint 4' can be rewritten as $\beta_2 + \beta_3 \leq 1$.

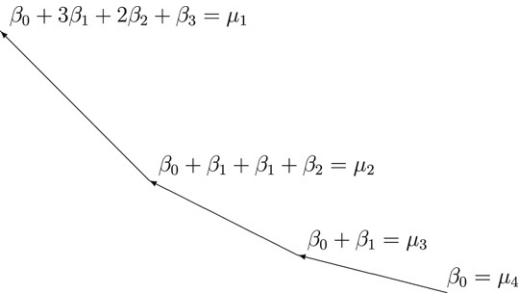


Fig. 1. Illustration of the dummy variables for call options.

In practice, we start with the construction of a design matrix which allows us to write the above model in the following linear form. For simplicity of presentation, we again set $p = 4$:

$$\begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \end{pmatrix} = \begin{pmatrix} 1 & 3 & 2 & 1 \\ 1 & 2 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}. \tag{5}$$

Ignoring the constraints on the coefficients would lead to a simple linear regression problem. Unfortunately, this approach does not have to lead, and usually does not, to interpretable and stable results.

Model (5) in the above form can be reasonably interpreted only if the observed strike prices are equidistant and if the distances between the neighboring observed strike prices are equal to one. If we want to keep the interpretation of the parameters β_j as the derivatives of the estimated function, we should use the design matrix

$$\Delta = \begin{pmatrix} 1 & \Delta_p^1 & \Delta_{p-1}^1 & \Delta_{p-2}^1 & \cdots & \Delta_3^1 & \Delta_2^1 \\ 1 & \Delta_p^2 & \Delta_{p-1}^2 & \Delta_{p-2}^2 & \cdots & \Delta_3^2 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \Delta_p^{p-2} & \Delta_{p-1}^{p-2} & 0 & \cdots & 0 & 0 \\ 1 & \Delta_p^{p-1} & 0 & 0 & \cdots & 0 & 0 \\ 1 & 0 & 0 & 0 & \cdots & 0 & 0 \end{pmatrix}, \tag{6}$$

where $\Delta_j^i = \max(k_j - k_i, 0)$ denotes the positive part of the distance between k_i and k_j , the i -th and the j -th ($1 \leq i \leq j \leq p$) sorted distinct observed values of the strike price.

The vector of conditional means μ can be written in terms of the parameters β as follows:

$$\begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix} = \mu = \Delta\beta = \Delta \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{pmatrix}. \tag{7}$$

The constraints on the conditional means μ_j can now be expressed as conditions on the parameters of the model (7). Namely, it suffices to request that $\beta_i > 0$, $i = 0, \dots, p - 1$ and that $\sum_{j=2}^{p-1} \beta_j \leq 1$.

The model for the option prices can now be written as

$$C(\mathcal{K}) = \mathcal{X}_\Delta \beta + \varepsilon, \tag{8}$$

where \mathcal{X}_Δ is the design matrix obtained by repeating each row of matrix Δ n_i times, $i = 1, \dots, p$.

3. Implementing the constraints

In order to impose Constraints 1–3 on parameters β_i , $i = 0, \dots, p - 1$, we propose the following reparameterization of the

model in terms of parameters $\theta = (\theta_0, \dots, \theta_{p-1})^\top$:

$$\beta_0(\theta) = \exp(\theta_0),$$

$$\beta_1(\theta) = \exp(\theta_1),$$

⋮

$$\beta_{p-1}(\theta) = \exp(\theta_{p-1}),$$

under the constraint that $\sum_{j=2}^{p-1} \exp(\theta_j) < 1$. Clearly, the parameters $\beta_i(\theta)$ satisfy the constraints

$$\beta_i(\theta) > 0, \quad i = 0, \dots, p - 1,$$

$$\sum_{j=2}^{p-1} \beta_j(\theta) < 1.$$

This means that the parameters $\beta_2(\theta), \dots, \beta_{p-1}(\theta)$ can be considered as point estimates of the state price density (the estimates have to be positive and integrate to less than one). Furthermore, in view of Lemma 1, it is worthwhile to note that the parameters also satisfy

$$-\sum_{j=1}^k \beta_j \in (-1, 0), \quad \text{for } k = 1, \dots, p - 1.$$

The model (8) rewritten in terms of parameters θ_i , $i = 0, \dots, p$, is a nonlinear regression model which can be estimated using standard nonlinear least squares or maximum likelihood methods (Seber and Wild, 2003). The main advantage of these methods is that the asymptotic distribution is well known and that the asymptotic variance of the estimator can be approximated using numerical methods implemented in many statistical packages.

3.1. Reparameterization

The following reparameterization of the model in terms of parameters $\xi = (\xi_0, \dots, \xi_p)^\top$ simplifies the calculation of the estimates because it guarantees that all constraints are automatically satisfied:

$$\beta_0(\xi) = \exp(\xi_0),$$

$$\beta_1(\xi) = \frac{\exp(\xi_1)}{\sum_{j=1}^p \exp(\xi_j)},$$

⋮

$$\beta_{p-1}(\xi) = \frac{\exp(\xi_{p-1})}{\sum_{j=1}^p \exp(\xi_j)}.$$

This property simplifies the numerical minimization algorithm needed for the calculation of the estimates.

The equality

$$\frac{1}{\sum_{j=1}^{p-1} \beta_j(\xi)} = 1 + \frac{\exp(\xi_p)}{\sum_{j=1}^{p-1} \exp(\xi_j)}$$

shows the meaning of the additional parameter ξ_p . Setting this parameter to $-\infty$ would be the same as requiring that $\sum_{j=1}^{p-1} \beta_j(\xi) = 1$. Large values of the parameter ξ_p indicate that the estimated coefficients sum to less than one or, in other words, the observed strike prices do not cover the support of the estimated SPD. Notice that, by setting $\xi_p = -\infty$, we could easily modify our procedure and impose the equality constraint $\sum_{j=1}^{p-1} \beta_j(\xi) = 1$.

3.2. Inverse transformation of model parameters

For the numerical algorithm, it is useful to know how to calculate ξ 's from given β 's. This is needed, for example, to obtain reasonable starting points for the iterative procedure maximizing the likelihood.

Lemma 3. Given $\beta = (\beta_1, \dots, \beta_p)^\top$, where $\beta_p = 1 - \sum_{i=1}^{p-1} \beta_i$, the parameters $\xi = (\xi_1, \dots, \xi_p)^\top$ satisfy the system of equations

$$(\beta \mathbf{1}_p^\top - \mathbf{I}_p) \exp \xi^\top = \mathcal{A} \exp \xi^\top = 0, \tag{9}$$

where \mathbf{I}_p is the $(p \times p)$ identity matrix. Furthermore,

$$\text{rank} \mathcal{A} = p - 1. \tag{10}$$

The system of Eq. (9) has infinitely many solutions, which can be expressed as

$$\exp(\xi) = (\mathcal{A}^- \mathcal{A} - \mathbf{I}_p) z, \tag{11}$$

where \mathcal{A}^- denotes a generalized inverse of \mathcal{A} and where z is an arbitrary vector in \mathbb{R}^p such that the right-hand side of (11) is positive.

Proof. Parts (9) and (10) follow from the definition of $\beta(\xi)$ and from simple algebra (notice that the sum of rows of \mathcal{A} is equal to zero). Part (11) follows, for example, from Anděl (1985, Theorem IV.18). \square

It remains to choose the vector z in (11) so that the solution of the system of Eq. (9) is positive.

Proposition 1. The rank of the matrix $\mathcal{A}^- \mathcal{A} - \mathbf{I}_p$ is 1. Hence, any solution of the system of Eq. (9) is a multiple of the first column of the matrix $\mathcal{A}^- \mathcal{A} - \mathbf{I}_p$. The vector z in (11) can be chosen, for example, as $z = \pm \mathbf{1}_p$, where the sign is chosen so that the resulting solution is positive.

Proof. The definition of a generalized inverse is

$$\mathcal{A} \mathcal{A}^- \mathcal{A} - \mathcal{A} = \mathcal{A} (\mathcal{A}^- \mathcal{A} - \mathbf{I}_p) = 0. \tag{12}$$

Lemma 3 says that $\text{rank} \mathcal{A} = p - 1$ and, hence, Eq. (12) implies that $\text{rank}(\mathcal{A}^- \mathcal{A} - \mathbf{I}_p) \leq 1$. Noticing that $\mathcal{A}^- \mathcal{A} \neq \mathbf{I}_p$ means that $\text{rank}(\mathcal{A}^- \mathcal{A} - \mathbf{I}_p) > 0$, and concludes the proof. \square

3.3. The algorithm

The proposed algorithm consists of the following steps:

- 1: obtain a reasonable initial estimate $\hat{\beta}$, for example, by running the Pool-Adjacent-Violators algorithm (Robertson et al., 1988, Chapter 1) on the unconstrained least squares estimates of the first derivative of the curve,
- 2: transform the initial estimate $\hat{\beta}$ into the estimate $\hat{\xi}$ using the method described in Section 3.2,
- 3: estimate the parameters of the model (8) by minimizing the sum of squares $\{C(\mathcal{K}) - \mathcal{X}_\Delta \beta(\xi)\}^\top \{C(\mathcal{K}) - \mathcal{X}_\Delta \beta(\xi)\}$ in terms of ξ (see Section 3.1) using numerical methods.

An application of this simple algorithm on real data is given in Section 5.1.

3.4. Asymptotic confidence intervals

We construct confidence intervals based on the parameterization $\beta(\theta)$ introduced at the beginning of this section. The confidence limits for parameters θ_i are exponentiated in order to obtain valid pointwise confidence bounds for the true SPD. The main advantage of this approach is that such confidence bounds are always positive.

An alternative approach would be to construct confidence intervals based on the parameterizations in terms of β_i (Section 2.3) or ξ_i (Section 3.1). However, the limits of confidence intervals for β_i may be negative and confidence intervals for the SPD based on parameters ξ_i would have very complicated shapes in high-dimensional space and could not be easily calculated and interpreted.

Another approach to the construction of the asymptotic confidence intervals can be based on the maximum likelihood theory. Assuming normality, the log-likelihood for the model (8) can be written as

$$l(C, \mathcal{X}_\Delta, \theta, \sigma) = -n \log \sigma - \frac{1}{2\sigma^2} \{C - \mathcal{X}_\Delta \beta(\theta)\}^\top \times \{C - \mathcal{X}_\Delta \beta(\theta)\}, \tag{13}$$

where \mathcal{X}_Δ is the design matrix given in (8). This normality assumption is justified later by a residual analysis. The maximum likelihood estimator is defined as

$$\hat{\theta} = \arg \max_{\theta} l(C, \mathcal{X}_\Delta, \theta, \sigma), \tag{14}$$

and it has asymptotically a p -dimensional normal distribution with mean θ and the variance given by the inverse of the Fisher information matrix:

$$\mathcal{F}_n^{-1} = \left\{ -E \left(\frac{\partial^2}{\partial \theta \partial \theta^\top} l(C, \mathcal{X}_\Delta, \theta, \sigma) \right) \right\}^{-1}. \tag{15}$$

More precisely, $n^{1/2}(\hat{\theta} - \theta) \xrightarrow{\mathcal{L}} N_p(0, \mathcal{F}_n^{-1})$. In this framework, the Fisher information matrix can be estimated by using the numerically differentiated Hessian matrix of the log-likelihood. For details we refer, for example, to Serfling (1980, Chapter 4). The confidence intervals calculated for parameters θ may be transformed (exponentiated) to a confidence intervals for the SPD (β). We have not pursued the maximum likelihood approach since it was numerically less stable in this situation.

Note that, under the assumptions of normality, the maximum likelihood estimate is equal to the nonlinear least squares estimate (Seber and Wild, 2003, Section 2.2), and the asymptotic variance of $\hat{\theta} = \exp(\beta)$ may be approximated by $\text{Var} \hat{\theta} = \{\text{diag}(\exp \hat{\theta}) \mathcal{X}_\Delta^\top \mathcal{X}_\Delta \text{diag}(\exp \hat{\theta})\}^{-1} \hat{\sigma}^2$. Hence, asymptotic confidence intervals for θ_i may be calculated as $(\hat{\theta}_i \pm u_{1-\alpha/2} \hat{s}_{ii})$, where $u_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of the standard Normal distribution and \hat{s}_{ii} denotes the i -th diagonal element of $\text{Var} \hat{\theta}$. By exponentiating both limits of this confidence interval, we immediately obtain the $1 - \alpha$ confidence interval for $\beta_i = \exp \theta_i$.

The construction of the estimator guarantees that the matrix \mathcal{X}_Δ has full rank—this implies that $\mathcal{X}_\Delta^\top \mathcal{X}_\Delta$ is invertible and the asymptotic variance matrix $\text{Var} \hat{\theta}$ always exists. If the number of observations is equal to the number of distinct strike prices (if there is only one option price for each strike price), it may happen that $\hat{\sigma}^2 = 0$ and the confidence intervals degenerate to a single point.

3.5. Put–Call parity

The prices of put options can be easily included in our estimation technique by applying the Put–Call parity of the option prices. Assuming that there are no dividends or costs connected with the ownership of the stock, each put option with price $P_t(K, T)$ corresponds to a call option with price

$$C_t(K, T) = P_t(K, T) + S_t - Ke^{-r(T-t)}.$$

In this way, the prices of the put options can be converted into the prices of call options and used in our model (Stoll, 1969). Statistically speaking, these additional observations will increase the precision of the SPD and will lead to more stable results.

In Germany, the Put–Call parity might be biased by an effect of the DAX index calculation which is based on the assumption that

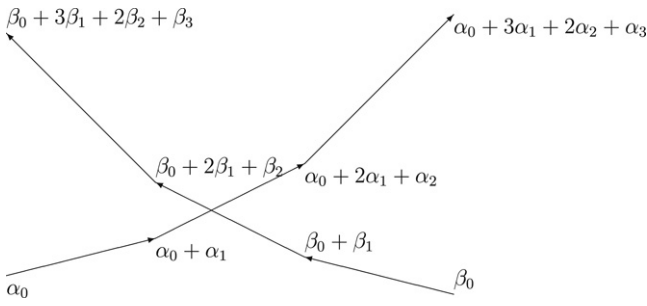


Fig. 2. Illustration of the dummy variables for both call (β) and put (α) options.

the dividends are reinvested after deduction of corporate income tax. As the income tax of some investors might be different, the value of the DAX has to be corrected before using Put–Call parity in subsequent analysis. For the exact description of this correction we refer to Hafner and Wallmeier (2000) who were analyzing the same data set.

The construction of our estimates allows us to include the put option prices in a more direct way by fitting the two curves separately using two sets of parameters. The situation is displayed in Fig. 2. Our assumption that the same SPD drives both the put and call option prices is naturally translated in terms of the coefficients α_i and β_i :

$$\alpha_i = \beta_{p-i+1}, \quad \text{for } i = 2, \dots, p - 1$$

$$\alpha_1 = 1 - \sum_{i=1}^{p-1} \beta_i.$$

The problem of estimating regression functions under such linear equality constraints is solved, for example, in Rao (1973). In Section 4.3, we will also investigate the covariance of the observed call and put option prices, and the suggested model will be presented in detail.

3.6. Expected value constraints

In Section 2.3, we have explained that the parameters $\beta_2, \dots, \beta_{p-1}$ can be interpreted as estimates of the state price density in points k_2, \dots, k_{p-1} . From the construction of the estimator, see also Fig. 1, it follows that parameter β_1 can be interpreted as the mass of the SPD lying to the right of k_{p-1} . Assuming that the observed strike prices entirely cover the support of the SPD, the mass β_1 could be attributed to the point k_p . Notice that the reparameterization introduced in Section 3 guarantees that $\sum_{i=1}^{p-1} \beta_i(\xi) < 1$, and it immediately follows that interpreting β_1 as the estimate of the SPD in point k_p does not violate any constraints described in Section 2.2.

Referring to Section 3.5, it is clear that the parameter $\beta_p \equiv \alpha_1 = 1 - \sum_{i=1}^{p-1} \beta_i$ can be interpreted as the estimator of the SPD in k_1 . The parameterization of the problem now guarantees that $\sum_{i=1}^p \beta_i = 1$.

The expected value of the underlying stock under the risk-neutral measure can now be estimated as $E^{\text{SPD}} = \sum_{i=1}^p k_i \beta_{p-i+1}$. From economic theory it follows that E^{SPD} has to be equal to the forward price of the stock. This constraint can be easily implemented by using the fact that β_1 and β_p estimate the mass of the SPD respectively to the right of k_{p-1} and to the left of k_2 .

If E^{SPD} is smaller than the forward price $\exp\{r(T-t)\}S_t$ of the stock, it suffices to move the mass β_1 further to the right. If E^{SPD} is too large, we move the mass β_p to the left. More precisely, setting $\tilde{k}_1 = k_1 - \mathbf{I}(E^{\text{SPD}} > \exp\{r(T-t)\}S_t)(E^{\text{SPD}} - \exp\{r(T-t)\}S_t)/\beta_p$, $\tilde{k}_p = k_p + \mathbf{I}(E^{\text{SPD}} < \exp\{r(T-t)\}S_t)(\exp\{r(T-t)\}S_t - E^{\text{SPD}})/\beta_1$,

we get

$$\exp\{r(T-t)\}S_t = \tilde{k}_1 \beta_p + \sum_{i=2}^{p-1} k_i \beta_{p-i+1} + \tilde{k}_p \beta_1.$$

This choice of \tilde{k}_1 and \tilde{k}_p guarantees that the expected value corresponding to the estimator β_1, \dots, β_p is equal to the forward price S_t of the stock; see the beginning of Section 6 for an application of this technique.

In Sections 4 and 5, we will concentrate on the properties of $\beta_2, \dots, \beta_{p-1}$ and further improvements in the estimation procedure.

4. Covariance structure

In this section, we use a model for the SPD development throughout the day to derive the covariance structure of the observed option prices depending on the strike prices and time of the trade. Considering the covariance structure in the estimation procedure solves the problems with heteroscedasticity and correlation of residuals that will be demonstrated in Section 5.3.

In this model, most recent option prices have the smallest variance and thus the largest weight in the estimation procedure. Similarly, the covariance of two option prices with the same strike price at approximately the same time is larger than the covariances of prices of some more dissimilar options.

We start by rewriting the model with iid error terms so that it can be more easily generalized. In Section 4.1, we present a model that accounts for heteroscedasticity and which is further developed in Sections 4.2 and 4.3, where an approximation of the covariance is calculated for any two options prices using only their strike prices and time of the trade. In Section 4.4, we suggest decomposing the error term into two parts, and we show how to estimate these additional parameters by the maximum likelihood method. The analysis of the resulting standardized residuals in Section 5.4 suggests that this covariance structure is applicable to our dataset.

Until now, we have assumed that the i -th option price (on a fixed day t) satisfies

$$C_i(k_j) = \Delta_j \tilde{\beta} + \varepsilon_i \tag{16}$$

or

$$C_i(k_j) = \Delta_j \tilde{\beta}_i + \varepsilon_i, \tag{17}$$

$$\tilde{\beta}_i = \tilde{\beta}_{i-1},$$

where ε_i are iid random errors with zero mean and constant variance σ^2 , $\tilde{\beta} = \tilde{\beta}_1 = \dots = \tilde{\beta}_i$ denotes the column vector of the unknown parameters, and Δ_j denotes the j -th row of the matrix Δ defined in (6), i.e.,

$$\Delta_j = (1, \Delta_p^j, \Delta_{p-1}^j, \dots, \Delta_{j+1}^j, \underbrace{0, \dots, 0}_{(j-1)}).$$

The residual analysis in Section 5.3 clearly demonstrates that the random errors ε_i are not independent and homoscedastic, and we have to consider some generalizations that lead to a better fit of the data set.

4.1. Heteroscedasticity

Assume that the i -th observation, corresponding to the j -th smallest exercise price k_j , can be written as

$$C_i(k_j) = \Delta_j \tilde{\beta}_i, \tag{18}$$

$$\tilde{\beta}_i = \tilde{\beta} + \varepsilon_i, \tag{19}$$

i.e., there are iid random vectors ε_i having iid components with zero mean and variances σ^2 in the state price density $\tilde{\beta}_i$. Clearly, the variance matrix of the vector of the observed option prices C is then

$$\text{Var } C = \sigma^2 \text{diag}(\mathcal{X}_\Delta \mathcal{X}_\Delta^T), \tag{20}$$

where \mathcal{X}_Δ is the design matrix in which each row of the matrix Δ is repeated n_j times, $j = 1, \dots, p$.

Remark 2. Assuming that the observed option prices have the covariance structure (20), the least squares estimates do not change, and

$$\text{Var } \hat{\beta} = \sigma^2 \{ \mathcal{X}_\Delta^\top \text{diag}(\mathcal{X}_\Delta \mathcal{X}_\Delta^\top)^{-1} \mathcal{X}_\Delta \}.$$

Another possible model for the heteroscedasticity would assume that the changes are multiplicative rather than additive.

$$C_i(k_j) = \Delta_j \tilde{\beta}_i \\ \log \tilde{\beta}_i = \log \tilde{\beta} + \varepsilon_i.$$

This model leads to a variance of $C_i(k_j)$ that depends on the value of the SPD:

$$\text{Var } C_i(k_j) = \sigma^2 \{ \beta_0^2 + (\Delta_p^j)^2 \beta_1^2 + (\Delta_{p-1}^j)^2 \beta_2^2 + (\Delta_{p-2}^j)^2 \beta_3^2 \\ + \dots + (\Delta_{j+1}^j)^2 \beta_j^2 \}.$$

It is straightforward that Remark 2 also applies in this situation.

4.2. Covariance

Let us now assume that there are random changes in the state price density coefficients β_i over time so that we have

$$C_i(k_j) = \Delta_j \tilde{\beta}_i, \\ \tilde{\beta}_i = \tilde{\beta}_{i-1} + \varepsilon_i, \tag{21}$$

where, for fixed i , $\tilde{\beta}_i$ is the parameter vector and ε_k , $k = i, i-1, \dots$, are iid random vectors having iid components with zero mean and variances σ^2 . For nonequidistant time points, let δ_i denote the time between the i -th and $(i-1)$ -th observation. The model is

$$C_i(k_j) = \Delta_j \tilde{\beta}_i, \\ \tilde{\beta}_i = \tilde{\beta}_{i-1} + \delta_i^{1/2} \varepsilon_i, \tag{22}$$

and it leads to the covariance matrix with elements

$$\text{Cov}\{C_{i-u}(k_j), C_{i-v}(k_i)\} = \text{Cov}(\Delta_j \tilde{\beta}_{i-u}, \Delta_i \tilde{\beta}_{i-v}) \\ = \sigma^2 \Delta_j \Delta_i^\top \sum_{l=1}^{\min(u,v)} \delta_{i+1-l}. \tag{23}$$

When we observe the i -th observation, we are usually interested in the estimation of the current value of the vector of parameters $\tilde{\beta}_i$.

4.3. Including put options

Similarly, we obtain the covariance for the price of the put options, $P_i(k_j)$. Using the relations between the α and β parameters, $\alpha_k = \beta_{p-k+1}$, for $k = 2, \dots, p-1$, and after some simplifications, we can write the model for the price of the put options, $P_i(k_j)$, as

$$P_i(k_j) = \Delta_j \tilde{\alpha}_i, \\ \tilde{\alpha}_i = \tilde{\alpha}_{i-1} + \delta_i^{1/2} \varepsilon_i, \tag{24}$$

where $\tilde{\alpha} = (\alpha_0, \alpha_1, \beta_{p-1}, \beta_{p-2}, \dots, \beta_2)^\top$ and Δ_j^p denotes the corresponding row of the design matrix, i.e.,

$$\Delta_j^p = (1, \Delta_j^1, \Delta_j^2, \dots, \Delta_j^{j-1}, \underbrace{0, \dots, 0}_{(p-j)}).$$

In this way, we obtain a joint estimation strategy for both the call and put option prices:

$$C_i(k_j) = \Delta_j \tilde{\beta}_i, \\ P_i(k_j) = \Delta_j^p \tilde{\alpha}_i, \\ \begin{pmatrix} \tilde{\beta}_i \\ \tilde{\alpha}_i \end{pmatrix} = \begin{pmatrix} \tilde{\beta}_{i-1} \\ \tilde{\alpha}_{i-1} \end{pmatrix} + \delta_i^{1/2} \varepsilon_i, \tag{25}$$

which directly leads to covariances

$$\text{Cov}\{P_{i-u}(k_j), P_{i-v}(k_i)\} = \text{Cov}(\Delta_j^p \tilde{\alpha}_{i-u}, \Delta_i^p \tilde{\alpha}_{i-v}) \\ = \sigma^2 \Delta_j^p (\Delta_i^p)^\top \sum_{l=1}^{\min(u,v)} \delta_{i+1-l} \tag{26}$$

and

$$\text{Cov}\{C_{i-u}(k_j), P_{i-v}(k_i)\} = \text{Cov}(\Delta_j \tilde{\beta}_{i-u}, \Delta_i^p \tilde{\alpha}_{i-v}) \\ = \sigma^2 \sum_{l=1}^{\min(u,v)} \delta_{i+1-l} \sum_{k=2}^{p-1} \Delta_{p+1-k}^j \Delta_i^{p+1-k}. \tag{27}$$

Together with (23), Eq. (26) and (27) allow us to calculate the covariance matrix of all observed option prices using only their strike prices and the times between the transactions.

4.4. Error term for option prices

Using the model (25) would mean that all changes observed in the option prices are due only to changes in the SPD. It seems natural to add another error term, η_i , as a description of the error in the option price:

$$C_i(k_j) = \Delta_j \tilde{\beta}_i + \eta_i, \\ P_i(k_j) = \Delta_j^p \tilde{\alpha}_i + \eta_i, \\ \begin{pmatrix} \tilde{\beta}_i \\ \tilde{\alpha}_i \end{pmatrix} = \begin{pmatrix} \tilde{\beta}_{i-1} \\ \tilde{\alpha}_{i-1} \end{pmatrix} + \delta_i^{1/2} \varepsilon_i, \tag{28}$$

where $\eta_i \sim N(0, v^2)$ are iid random variables independent of the random vectors ε_i . Here, normality assumptions are added both for η_i and ε_i so that the variance components parameters v^2 and σ^2 may be estimated by the maximum likelihood method.

Next, in order to simplify the notation, let us fix the index i , and let Y denote the vector of observed call and put option prices, \mathcal{X}_Δ the corresponding design matrix consisting of the corresponding rows Δ_j and Δ_j^p , and $\tilde{\gamma}$ the combined vector of unknown parameters. Denoting by Σ_i the matrix containing the covariances defined in (23), (26) and (27), we can rewrite model (25) as

$$Y = \mathcal{X}_\Delta \tilde{\gamma} + \xi, \tag{29}$$

where $\text{Var } \xi = \text{Var } Y = \sigma^2 \Sigma_i + v^2 \mathbf{I}_n = \sigma^2 (\Sigma_i + \psi^2 \mathbf{I}_n) = \sigma^2 V$, where $\psi^2 = v^2 / \sigma^2$. Differentiating the log-likelihood

$$l(\beta, \sigma^2, \psi^2) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\sigma^2 V| \\ - \frac{1}{2\sigma^2} (Y - \mathcal{X}_\Delta \tilde{\gamma})^\top V^{-1} (Y - \mathcal{X}_\Delta \tilde{\gamma}),$$

we obtain

$$\frac{\partial l(\beta, \sigma^2, \psi^2)}{\partial \psi^2} \\ = -\frac{1}{2} \text{tr}(V^{-1}) + \frac{1}{2\sigma^2} (Y - \mathcal{X}_\Delta \tilde{\gamma})^\top V^{-2} (Y - \mathcal{X}_\Delta \tilde{\gamma}). \tag{30}$$

For any fixed value of the parameter ψ^2 , it is straightforward to calculate the optimal σ^2 and $\tilde{\gamma}$. Hence, the numerical maximization of the log-likelihood can be based on a search for a root (zero) of the one-dimensional function (30).

Moreover, the variance components parameters σ^2 and $v^2 = \psi^2 \sigma^2$ have a very natural econometric interpretation: σ^2 describes the speed of change of the SPD and v^2 the error in observed option prices.

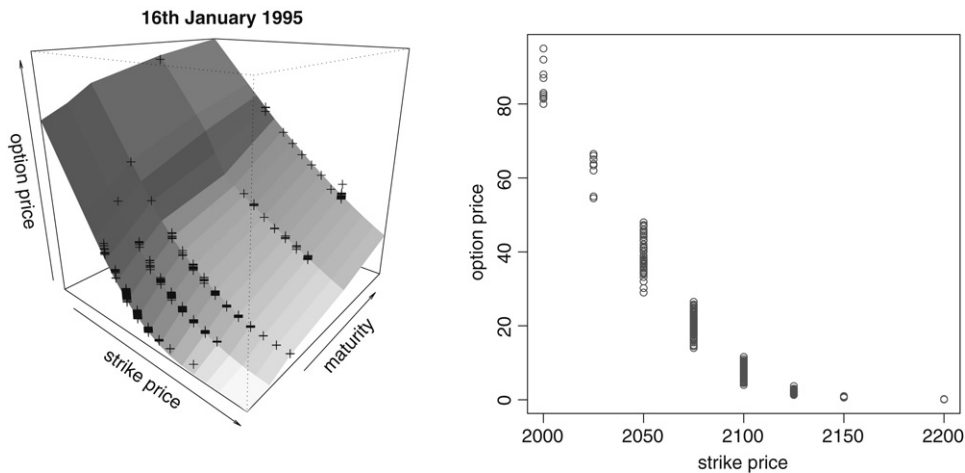


Fig. 3. Option prices plotted against strike price and time to maturity with a two-dimensional kernel regression surface (left) in January 1995 and the ensemble of the call option prices with shortest time to expiry against strike price (right) on 16 January 1995. SFB and CASE data base: sfb649.wiwi.hu-berlin.de.

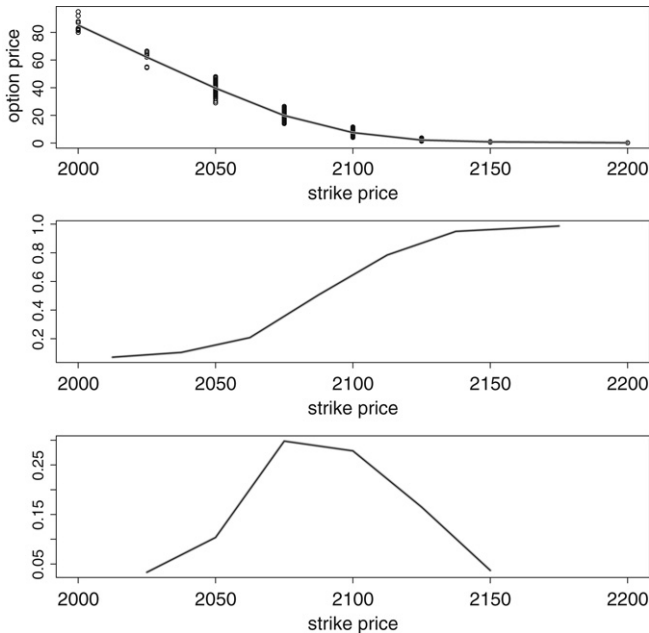


Fig. 4. On 16 January 1995, the unconstrained estimate satisfies the constraints. Hence, it is equal to the constrained estimate. The top panel shows the original data with the fitted call pricing functions. The second and the third panels show the estimates of the first and second derivatives, respectively.

5. Application to DAX data

We analyze a data set containing observed option prices for various strike prices and maturities. Other variables are the interest rate, date, and time. In 1995, one observed every day about 500 trades; in today's more liquid option markets this number has increased approximately 10 times. In our empirical study we will consider the time period from 1995 to 2003, thus also covering more recent liquid option market.

Fig. 3 displays the observed prices of European call options written on the DAX for the 16 January 1995. The left panel shows the ensemble of call option prices for different strikes and maturities as a free structure together with a smooth surface. The typical shape of dependency of the option price on the strike price can be observed in the right panel, containing the option prices only for the shortest time to expiry, $\tau = T - t = 4$ days.

In order to illustrate the method, we apply it to DAX option prices on two consecutive days. These days (16 and 17 January

1995) were selected since they provide a nice insight into the behavior of the presented methods.

5.1. Estimator with iid random errors

We start by a comparison of the unconstrained and constrained estimator described respectively in Sections 2.3 and 3.1.

For the European call option prices displayed in the right-hand plot in Fig. 3, we obtain the estimates plotted in Fig. 4. The top plot displays the original data, the second plot shows the estimate of the first derivative, and the third plot shows the estimate of the second derivative, i.e., the state price density. Actually, all plots contain two curves, both obtained using model (8). The thick line is calculated using the parameters β_i without constraints, whereas the thin line uses the reparameterization $\beta_i(\xi)$ given in Section 3.1. In Fig. 4, these two estimates coincide since the model maximizing the likelihood without constraints, by chance, fulfills the constraints ($\exists \xi : \beta_i = \beta_i(\xi), i = 0, \dots, p - 1$), and hence it is clear that the same parameters also maximize the constrained likelihood.

The situation, in which the call pricing functions fitted with and without constraints differ, is displayed in Fig. 5. Notice that the difference between the two regression curves is small, whereas the difference between the estimates of the state price density (i.e., the second derivative of the curve) is surprisingly large. The unconstrained estimate shows very unstable behavior on the left-hand side of the plot. The constrained version behaves more reasonably. Very small differences between the fitted call pricing functions in the top plot in Fig. 5 lead to huge differences in the estimates of the second derivative.

We therefore conclude that a small error in the estimate of the call pricing function may lead to large scale error in the estimates of the first and second derivatives. The scale of this type of error seems to be limited by imposing the shape constraints given in Section 2.2.

5.2. Confidence intervals

In Figs. 6 and 7, we plot both estimates together with the 95% confidence intervals. Notice that, in the unconstrained model, the estimates of the values of the SPD are just the parameters of the linear regression model. Hence, the confidence intervals for the parameters are, at the same time, also confidence intervals for the SPD. These confidence intervals for 16 and 17 January 1995 are displayed in the upper plots in Figs. 6 and 7. The drawbacks of this method are clearly visible. In Fig. 6, the lower bounds of the confidence intervals only asymptotically satisfy the condition of

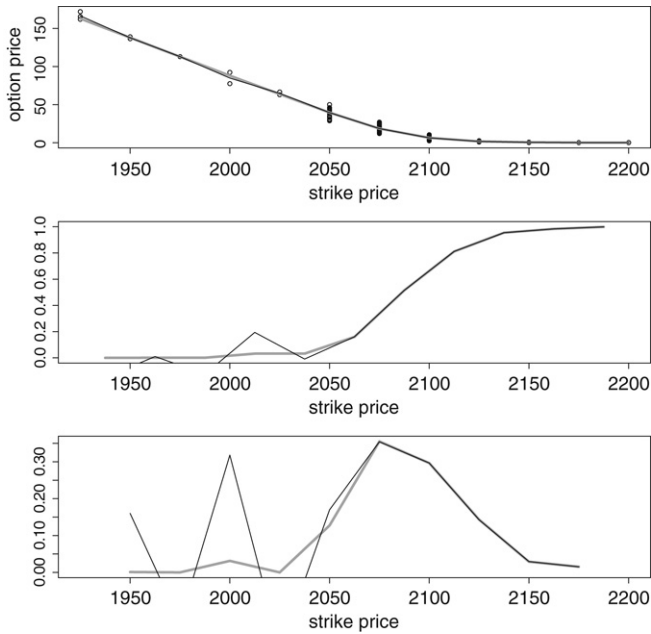


Fig. 5. On 17 January 1995, the unconstrained estimate, displayed using the thin line, does not satisfy the constraints. The top panel shows the original data with the two fitted call pricing functions. The estimates of the first derivative in the second panel look rather different. The constrained estimate of the second derivative in the bottom panel is clearly much more stable than the unconstrained estimate.

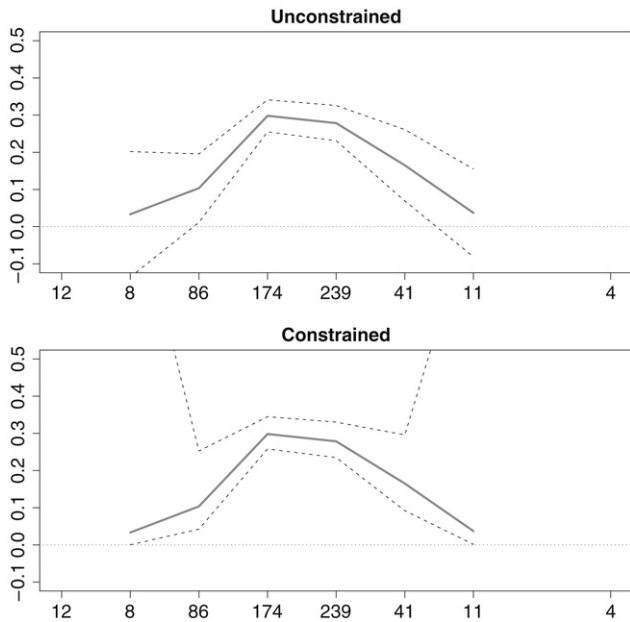


Fig. 6. The unconstrained and constrained confidence intervals for the SPD on 16 January 1995. The description on the x-axis shows the number of observations in each point.

positivity. In Fig. 7, we observe large variability on the left-hand side of the plot (the region with low number of observations). Again, some of the lower bounds are not positive. Clearly, the confidence intervals based on the unconstrained model make sense only if the constraints are, by chance, satisfied. Even if this is the case, there is no guarantee that the lower bounds will be positive. The lower panels in Figs. 6 and 7 display the nonnegative asymptotic confidence intervals calculated according to Section 3.4.

In Fig. 6, both types of confidence interval provide very similar results. The only difference is at the minimum and maximum value

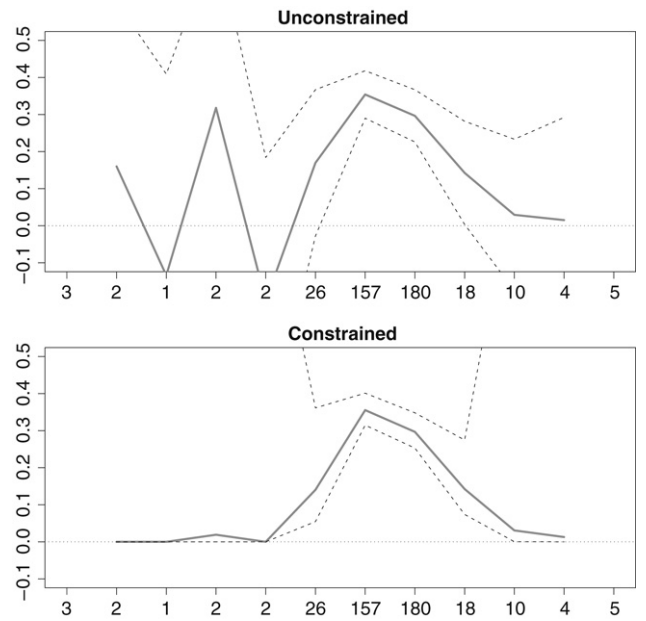


Fig. 7. Confidence intervals for SPD on 17 January 1995. The description on the x-axis shows the number of observations in each point.

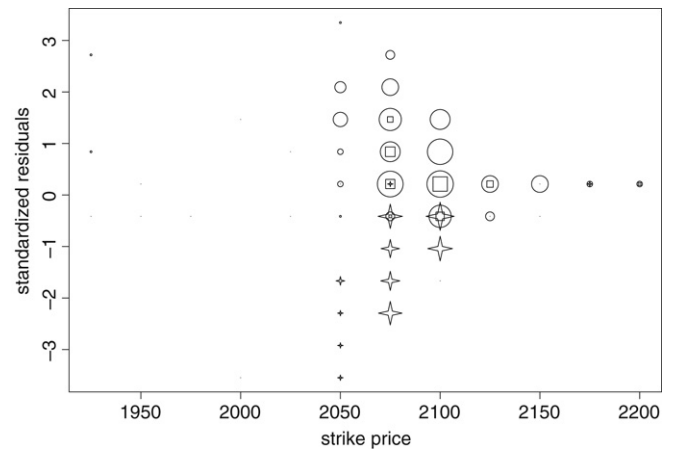


Fig. 8. The time dependency and the heteroscedasticity of the residuals during one day. The circle, square, and star denote the trades carried out in the morning, midday, and afternoon, respectively. The size of the symbols denotes the number of residuals.

of the independent variable (strike price), where the unconstrained method provides negative lower bounds and the conditional method leads to very large upper bounds of the confidence intervals.

In Fig. 7, we plot the confidence intervals for 17 January 1995. In the central region of the graphics, both types of confidence interval are quite similar. On the left-hand and right-hand sides, both methods tend to provide confidence intervals that seem to be overly wide. For the constrained method, we observe that the length of the confidence intervals explodes when the estimated value of the SPD is very close to zero and, at the same time, the number of observation in that region (see the description of the horizontal axis) is small.

5.3. Residual analysis

The residuals on 17 January 1995 are plotted in Fig. 8. The time of trade (in hours) is denoted by the plotting symbol. The circle, square, and star denote the trades carried out in the morning, midday, and afternoon, respectively. The size of the symbols

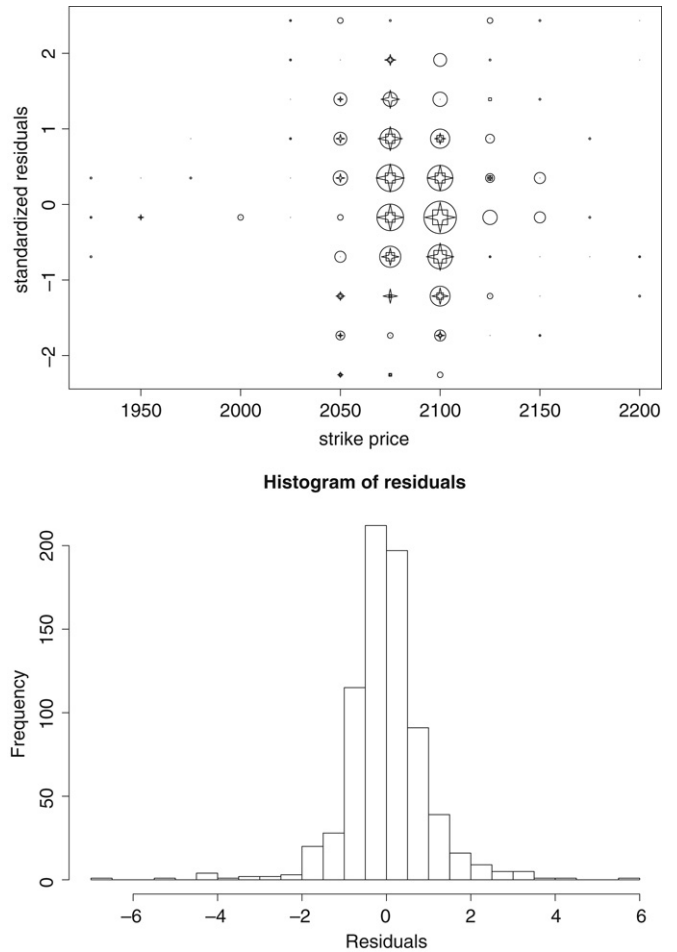
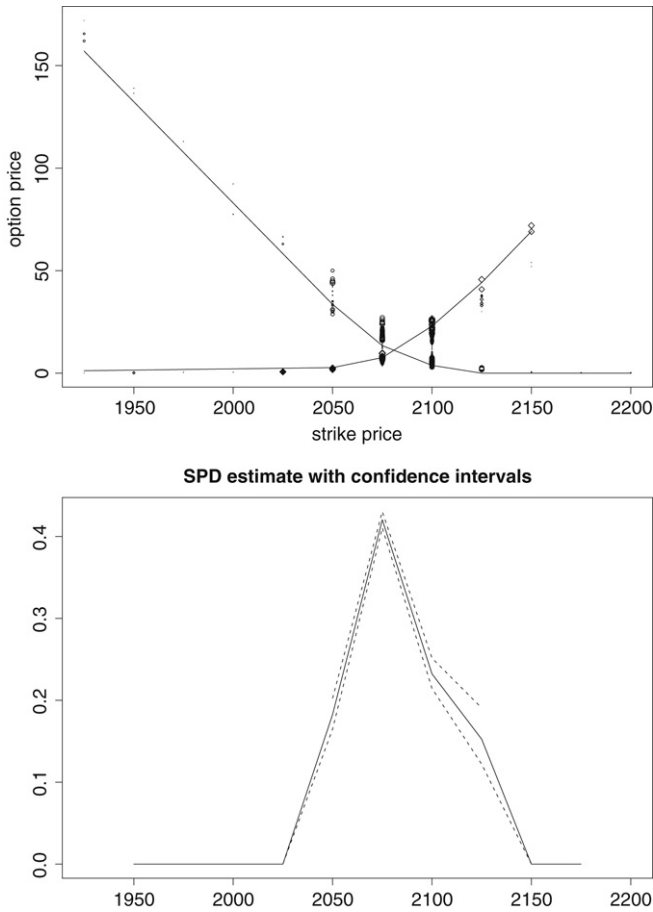


Fig. 9. Estimate using the covariance structure (28) on 17 January 1995. The upper plot shows the observed option prices and the constrained estimate. The size of the plotting symbols corresponds to the weight of the observations. The lower plot shows the estimated SPD with confidence intervals.

Fig. 10. The development of the standardized residuals resulting from the model with the covariance structure (28) on 17 January 1995 during the day, where circles, squares, and stars denote the residuals from morning, midday, and afternoon, and a histogram of the standardized residuals.

corresponds to the number of residuals lying in the respective areas.

The majority of the residuals correspond to the strike prices of 2075DEM and 2100DEM. The variance of the residuals is very low on the right-hand side of the plot and it rapidly increases when moving towards smaller strike prices. On the left-hand side of the plot, for strike prices smaller than 2000, we have only very few observations, and cannot judge the residual variability reliably.

Apart from the obvious heteroscedasticity we also observe a very strong systematic movement in the SPD throughout the day: the circles, corresponding to the first third of the day, are positive, and all stars, denoting the afternoon residuals, are negative. Similar patterns can be observed every day—residuals corresponding to the same time have the same sign.

We conclude that the assumption of iid random errors is obviously not fulfilled as the option prices tend to follow the changes of the market during the day.

5.4. Application of the covariance structure

In Fig. 9, we present the estimator combining both put and call option prices and using the covariance structure proposed in Section 4.4. In comparison with the results plotted in Fig. 7, we observe shorter length of the confidence intervals.

The estimates of the variance components parameters are $\hat{\psi}^2 = 17.77$, $\hat{\sigma}^2 = 0.0041$, and $\hat{\nu}^2 = 0.0722$. For interpretation, it is more natural to consider $\hat{\nu} = 0.2687$, suggesting that 95% of the option prices were on 17 January 1995 not further than 0.5DEM from the correct option price implied by the current (unobserved) SPD.

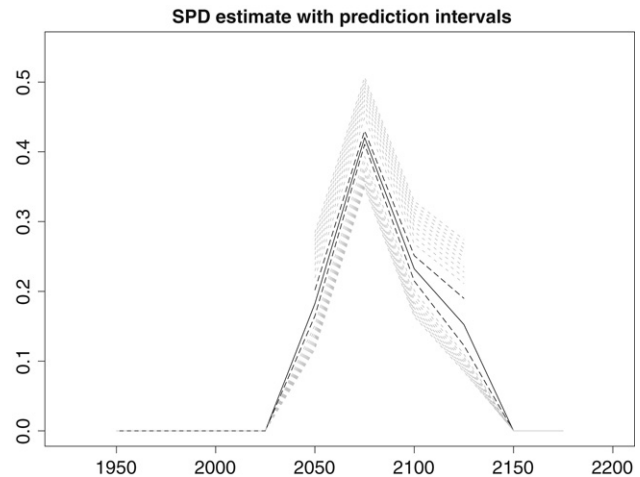


Fig. 11. SPD estimate on 17 January 1995 with prediction intervals for the next 5 h calculated for every 30 min.

The standardized residuals in the top panel of Fig. 10 were plotted using the same technique as the residuals in Fig. 8. Whereas the residuals for the iid model showed strong correlations and heteroscedasticity, the structure of the standardized residuals looks much better. It is natural that the residuals are larger in the central part since more than 90% of observations have strike price

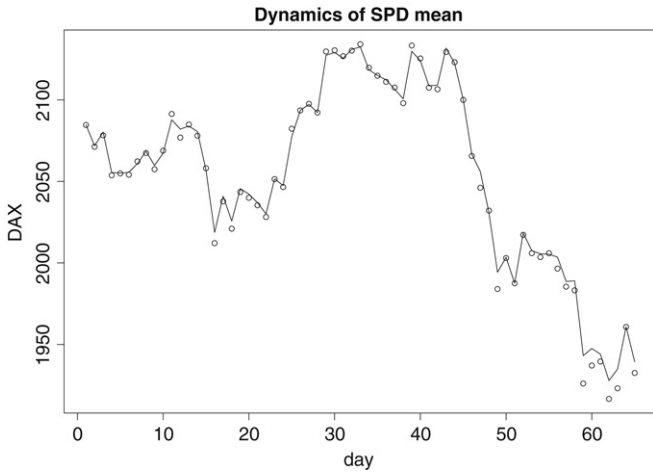


Fig. 12. Daily development of the expected value of the uncorrected SPD from January to March 1995. The circles denote the corresponding closing value of the DAX.

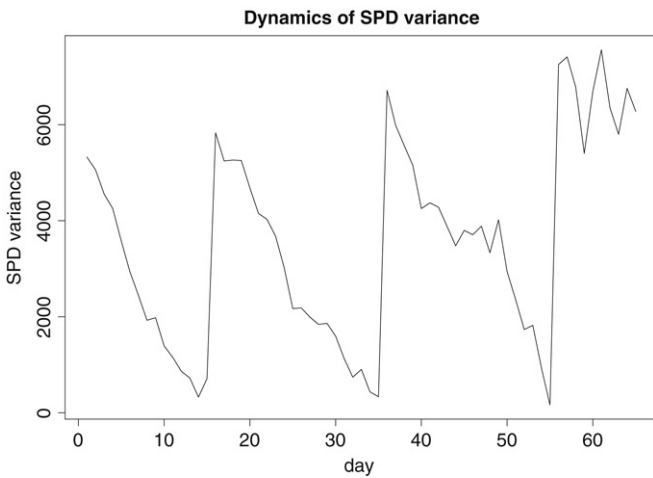


Fig. 13. Daily development of the SPD variance from January to March 1995.

between 2050 and 2100. The largest residuals were omitted in the residual plot so that the structure in the central part is more visible, but the lower panel of Fig. 10 displays the histogram of all residuals. The distribution of the residuals seems to be symmetric, and its shape is not too far from Normal distribution. However, the kurtosis of this distribution is too large, and formal tests reject normality.

In Fig. 11, we plot prediction intervals for the SPD obtained only by recalculating the covariance structure (28) with respect to some future time. More precisely, the prediction intervals are obtained from option prices observed until \hat{i} . Then, using the notation of Section 4.4, we have, for the future β_{i+1} and $\tilde{\alpha}_{i+1}$,

$$\begin{aligned} C_i(k_j) &= \Delta_j \tilde{\beta}_i + \eta_i, \\ P_i(k_j) &= \Delta_j^p \tilde{\alpha}_i + \eta_i, \\ \begin{pmatrix} \tilde{\beta}_{i+1} \\ \tilde{\alpha}_{i+1} \end{pmatrix} &= \begin{pmatrix} \tilde{\beta}_i \\ \tilde{\alpha}_i \end{pmatrix} + \delta_{i+1}^{1/2} \varepsilon_{i+1}. \end{aligned} \tag{31}$$

It is now easy to see that the only modification that has to be done for estimating $\tilde{\beta}_{i+1}$ is to add the length of the forecasting horizon δ_{i+1} to the sum in (23), (26) and (27), and to recalculate the confidence regions using this variance matrix with the same estimates of the variance parameters σ^2 and ν^2 . In Fig. 11, the 95% confidence intervals for the true SPD are denoted by the black

dashed line. The grey dashed lines denote the prediction intervals calculated for each 30 min for the next 5 h. In this way, we can obtain a simple approximation for future short-term fluctuations of the SPD. In the long run, the prediction intervals become too wide to be informative.

6. Dynamics of the SPD

In order to study the dynamics of SPDs, we calculated the basic moment characteristics of the estimated SPDs. Note that the estimator does not allow one to estimate the SPD in the tails of the distribution. We can only estimate the probability mass lying to the left ($1 - \sum_{i=1}^{p-1} \beta_i$) and to the right (β_1) of the available strike price range. Hence, the moments calculated in this section are only approximations which cannot be calculated more precisely without additional assumptions, for example, on the tail behavior or parametric shape of the SPD.

The estimated mean and variance in the first quarter of 1995 are plotted as lines in Figs. 12–13. Note that the SPDs in this period were always estimated using the options with shortest time to maturity. This means that the time to maturity is decreasing linearly in both plots, but it jumps up whenever the option with the shortest time to maturity expires. These jumps occurred at days 16, 36, and 56.

From no-arbitrage considerations, it follows that the mean of the SPD should correspond to the value of the DAX,

$$\widehat{E}^{SPD} = \int S_T f(S_T) dS_T = \exp\{r(T - t)\} S_t.$$

See also the discussion in Section 3.6. In Fig. 12, the observed values of the DAX multiplied by the factor $\exp\{r(T - t)\}$ are plotted as circles for the first 65 trading days in 1995, and we observe that the estimated means of the SPD estimates, displayed as the line, follow the theoretical value very closely. A small difference is mainly due to the fact that, in 1995, the observed strike prices do not entirely cover the support of the SPD. For example, on day 16, the difference between the SPD mean (2018.7) and the DAX multiplied by the discount factor (2012.1) is equal to 6.6. The fact that there are not any trades for strike prices smaller than 1925 means that we only know that the probability mass lying to the left from 1950 is equal to 0.25. In the calculation of the estimate of the SPD mean plotted in Fig. 12, this probability mass is assigned to the value 1925, as this is the leftmost observed strike price. Obviously, assigning this probability mass rather to the value $1925 - (6.6/0.25) = 1898.6$ leads a more realistic estimate of the SPD and to the equality of the SPD mean and the discounted DAX.

In Fig. 13, we see that the variance of the SPD decreases linearly as the option moves closer to its maturity. This observation suggests that SPD estimates calculated for neighboring maturities can be linearly interpolated in order to obtain an SPD estimate with arbitrary time to maturity. Such an estimate is important for making the SPD estimates comparable and for studying the development of the market expectations.

6.1. Estimate with the fixed time to expiry

The variances displayed in Fig. 13 suggest that the variance of the SPD estimates changes approximately linearly in time when moving closer to the date of expiry.

Hence, from the estimates $f_{\tau_1}(\cdot)$ and $f_{\tau_2}(\cdot)$ of centered SPDs corresponding to the times of expiry $\tau_1 < \tau_2$, we construct an estimate $f_\tau(\cdot)$ for any time of expiry $\tau \in (\tau_1, \tau_2)$ as

$$f_\tau(\cdot) = \frac{(\tau_2 - \tau)f_{\tau_1}(\cdot) + (\tau - \tau_1)f_{\tau_2}(\cdot)}{\tau_2 - \tau_1}. \tag{32}$$

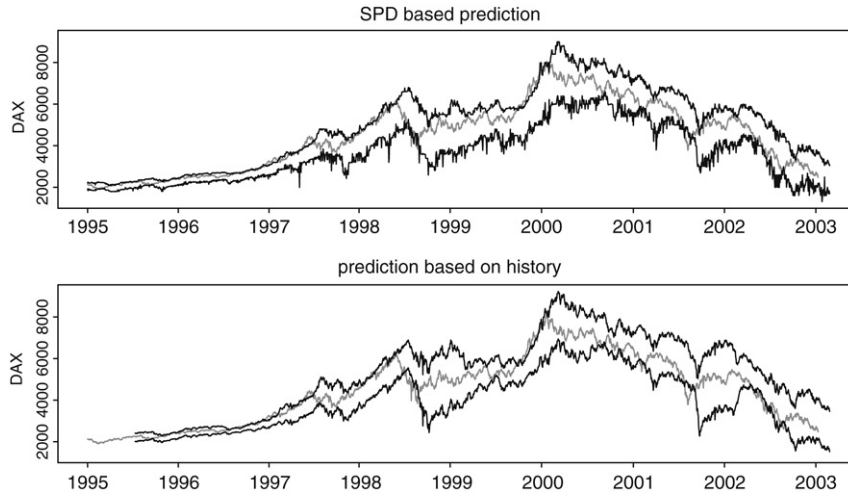


Fig. 14. Prediction intervals for the DAX based on SPDs and historical simulation from January 1995 to March 2003.

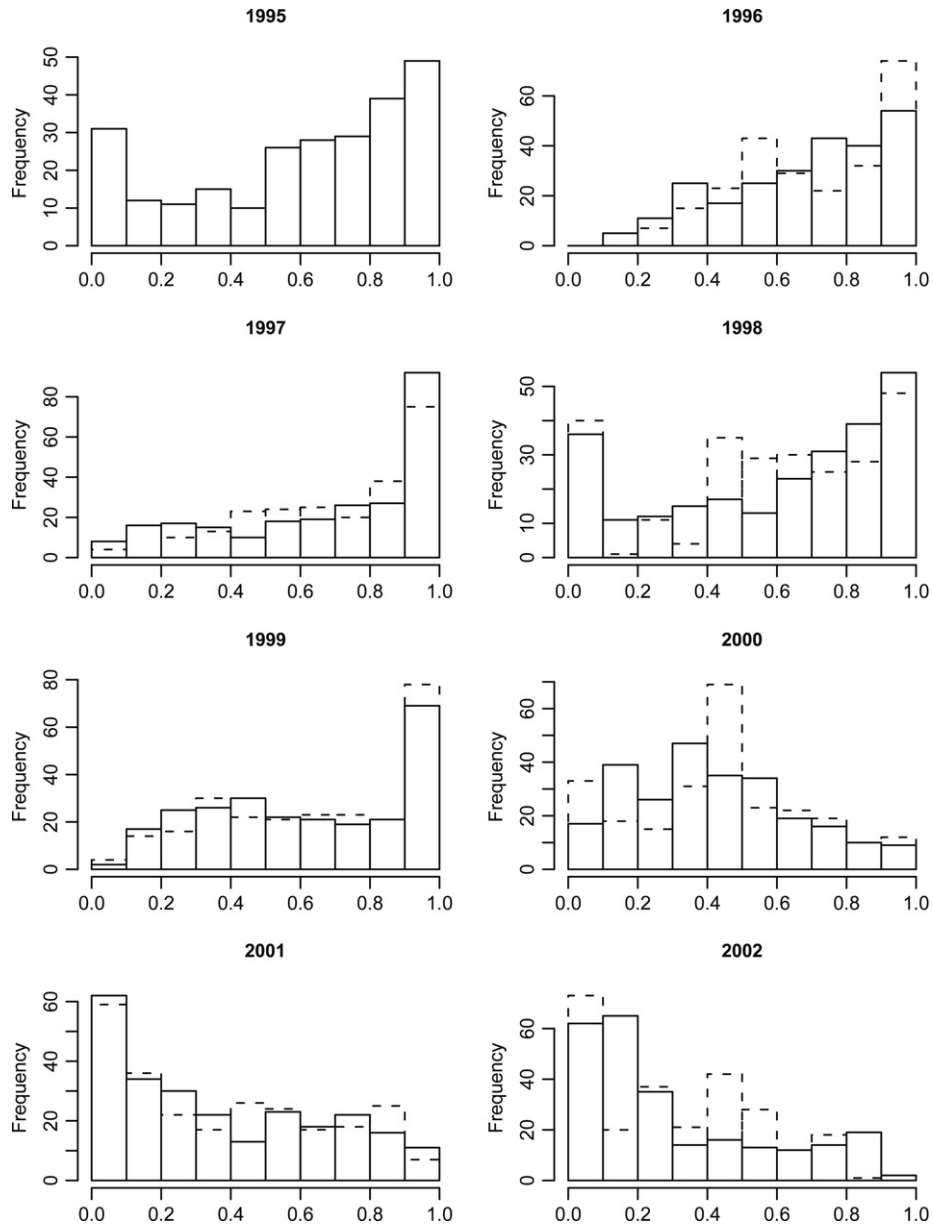


Fig. 15. Histograms for the SPDs (full line) and historical simulation (dashed line).

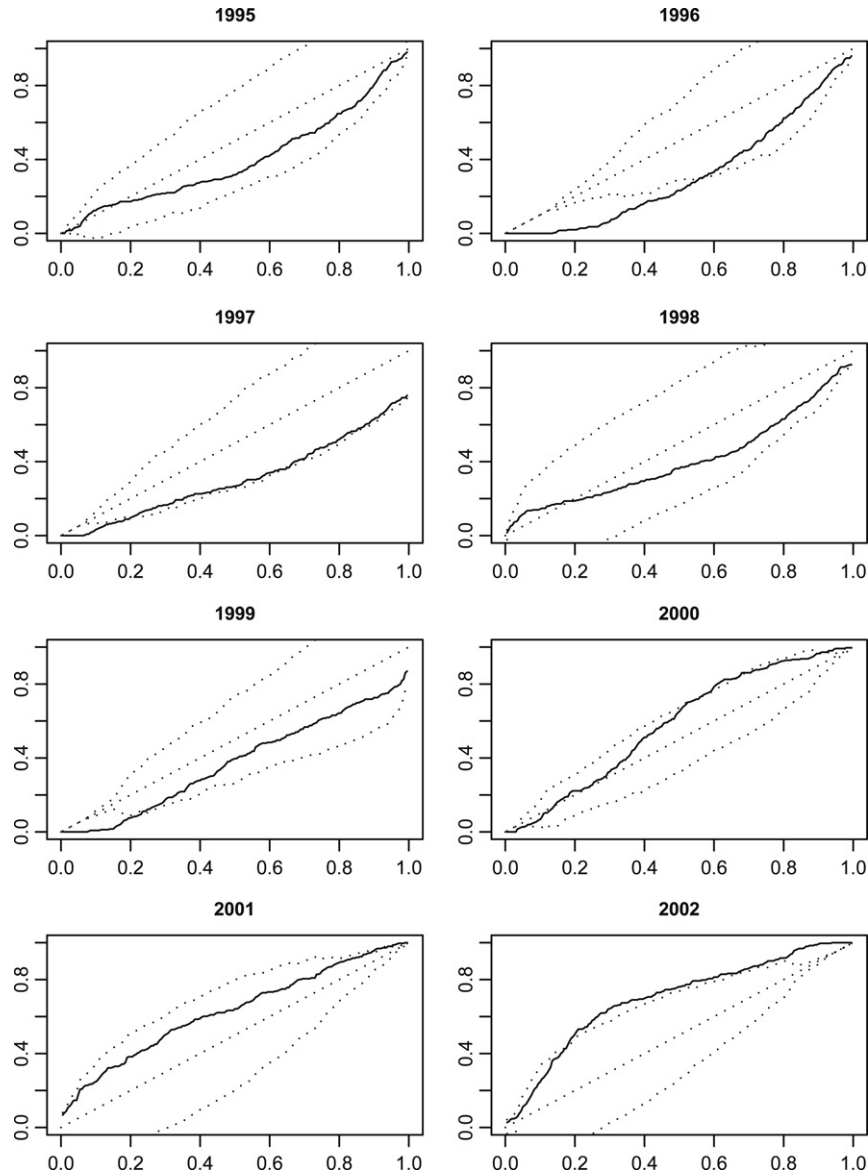


Fig. 16. Integral transformation for estimated SPDs.

In this way, the variance, V_τ , of the centered SPD with time to expiry equal to τ can be expressed as

$$\begin{aligned} V_\tau &= \int x^2 f_\tau(x) dx \\ &= \int x^2 \frac{(\tau_2 - \tau) f_{\tau_1}(x) + (\tau - \tau_1) f_{\tau_2}(x)}{\tau_2 - \tau_1} dx \\ &= \frac{(\tau_2 - \tau) V_{\tau_1} + (\tau - \tau_1) V_{\tau_2}}{\tau_2 - \tau_1}. \end{aligned}$$

We argue that such an estimate is reasonable since we observed in Fig. 13 that the SPD variances change linearly in time.

6.2. Verification of the market's expectations

Under the risk neutral (equivalent martingale) measure, the SPD reflects the market's expectation of the behavior of the value of the DAX in 45 days. Hence, it is interesting to use our data set to verify how these expectations compare with reality. In the left plot in Fig. 14, we plot intervals based on the SPD together with the true future value of the DAX: the black lines display the 2.5% and

97.5% quantiles of the estimated SPD; the future value of the DAX is displayed as a grey line. In the right plot, we show in the same way the 45-day ahead predictions based on the historical distribution of the 45-day absolute returns in the last 100 trading days; the 2.5% and 97.5% quantiles of this distribution are plotted as black lines.

Fig. 14 suggests that the method works well and that the DAX mostly stays well within the quantiles calculated from the estimated SPDs. The DAX was sometimes rising faster than the market expected from 1995 to mid-1998. After a fast decrease in the second half of 1998, the market increased again till the beginning of year 2000. Since then, the market has decreased. However, the changes stay mostly within or very close to the bounds predicted by our SPD estimates. The only exception is the large shock observed in September 2001, caused by the terrorist attack on the World Trade Center.

The upper quantiles, 97.5%, of the historical distribution of the 45-day absolute returns mostly agree with the upper quantiles of the SPD. The lower quantiles, 2.5%, of the SPDs seem to be much more variable than the same quantiles of the historical distribution. Both the lower and the upper quantiles of the historical distribution lie mostly above the corresponding quantiles of the estimated SPD, respectively in 69.44% and 81.75%.

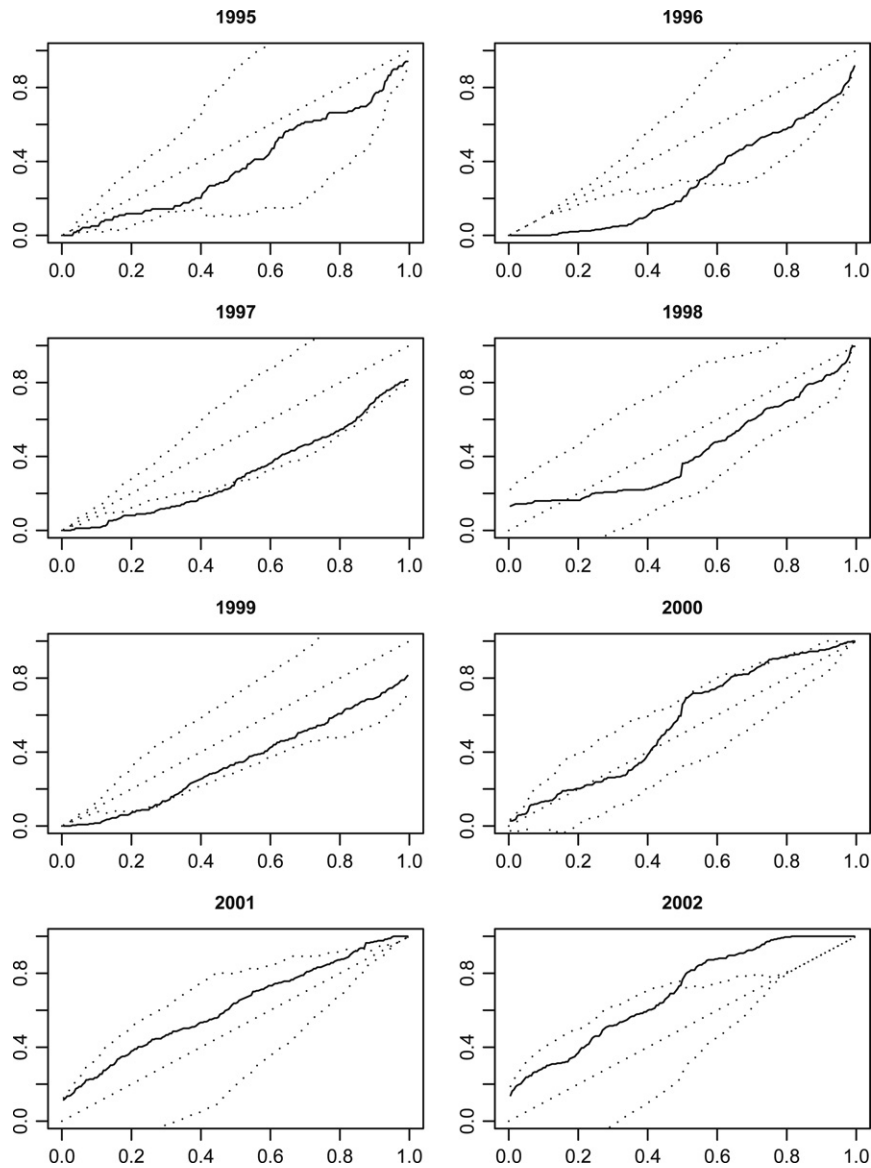


Fig. 17. Integral transformation for historical simulation.

Table 2
Fraction of the year that the DAX stays in the prediction corridor.

Year	1996	1997	1998	1999	2000	2001	2002
SPD (%)	84.40	66.13	75.30	74.60	97.22	85.66	94.84
Historical (%)	82.00	79.44	76.89	77.38	93.25	86.06	80.56

This observation just confirms the fact that the observed SPD includes effects of risk aversion.

In Table 2, we show the fraction of the year that the DAX stays in the prediction corridor. This suggests that the coverage is slightly better for the historical simulation if the DAX is increasing and better for the SPD based prediction if the DAX is decreasing (years 2000 and 2002).

6.3. Evaluation of the quality of the forecasts

The quality of the forecasts can be evaluated by comparing the true future observation with its predicted distribution (the SPD). Diebold et al. (1998) propose to evaluate density forecasts using the probability integral transformed observations $z_{h,t}$, where t denotes the time and h the forecasting horizon. More precisely,

we define

$$z_{h,t} = \int_{-\infty}^{X_{t+h}} \hat{f}_{h,t}(u) du,$$

where $\hat{f}_{h,t}(\cdot)$ denotes our estimate of the SPD h days ahead at time t and X_{t+h} is the future observation. In other words, $z_{h,t}$ is the probability value of X_{t+h} with respect to $\hat{f}_{h,t}(\cdot)$. Clearly, the $z_{h,t}$ should be uniformly $U(0, 1)$ distributed if the estimated SPD $\hat{f}_{h,t}(\cdot)$ is equal to the true density of X_{t+h} . In Fig. 15, we display the histograms of $z_{h,t}$'s for each year for the estimated SPDs and historical simulation using full and dashed histograms, respectively. Clearly, in the ideal case, the histograms should not be too far from a Uniform $U(0, 1)$ distribution. In our data, for the prediction horizon $h = 45$ days, we observe that the histograms look quite different from what we would expect. Especially in

years 1995–1999, the DAX was moving mainly in the upper quantiles of the predicted SPD. The forecasts based on the historical distribution of the 45-day returns behave similarly.

In order to account for the overlapping forecasting periods, we calculate the confidence limits for the empirical distribution function

$$\widehat{F}(u) = \frac{1}{T} \sum_{t=1}^T \mathbf{I}(z_{h,t} \leq u)$$

of $z_{h,t}$'s that take into account the autocorrelation structure.

$$\widehat{\text{Var}}\{\widehat{F}(u)\} = \frac{1}{T} \left\{ \widehat{\gamma}_u(0) + 2 \sum_{j=1}^h \left(1 - \frac{j}{T}\right) \widehat{\gamma}_u(j) \right\}, \quad (33)$$

where $\gamma_u(j)$ is the sample autocovariance of order j :

$$\gamma_u(j) = \frac{1}{T} \sum_{t=j+1}^T \{ \mathbf{I}(z_{h,t} \leq u) - \widehat{F}(u) \} \{ \mathbf{I}(z_{h,t-j} \leq u) - \widehat{F}(u) \}.$$

The empirical distribution functions $\widehat{F}(\cdot)$ are plotted separately for years 1995–2002 in Fig. 16. The distribution function of $U(0, 1)$ and the limits following from (33) are displayed as dotted lines. The year 2003 was not included since our dataset contains only two months of the year 2003, which did not leave enough observations to confirm the forecasts.

In 1996 and 1997, the market was growing much faster than the SPDs were indicating. In 1996, it never happened that the DAX fell below the 10% quantile of the SPD, and there were only a few days when this value was below 20%. The situation in 1998 and 1999 was less extreme even though the fast growth of the DAX continued. The distribution given by the SPD estimate $\widehat{f}_{t,h}(\cdot)$ for the horizon $h = 45$ days does not differ significantly from the true distribution of X_{t+h} in 2000–2001, but in 2002 we again observe significant differences. Thus, the DAX was growing faster than the option market expected in 1996, 1997, and 1999 and it was falling faster in 2002.

Fig. 17 shows the same graphics for the forecast based on the historical distribution of the returns. The deviations are more clearly visible but the overall picture is very similar; the only difference arises in 2001 when the predictions did not stay between the limits.

7. Conclusion

We have proposed a simple nonparametric model for arbitrage-free estimation of the SPD. Our procedure takes care of the daily changing covariance structure and involves both types of European option. Moreover, the covariance structure allows us to calculate prediction intervals capturing future behavior of the SPD. We analyze the moment dynamics of the SPD from 1995–2003. An application to DAX EUREX data for the years 1995–2003 produces a corridor that is compared to the future DAX index value. The proposed technique enables us not only to price exotic options but also to measure the risk and volatility ahead of us.

Acknowledgments

We thank Volker Krätschmer for useful comments concerning the existence and uniqueness of the constrained regression function and the anonymous referee for many insightful comments leading to substantial improvements in both the presentation and the content of the paper. The research was supported by Deutsche Forschungsgemeinschaft, SFB 649 “Ökonomisches Risiko”, by MSM0021620839, GAČR GA201/08/0486, and by MŠMT 1K04018.

References

- Ait-Sahalia, Y., Duarte, J., 2003. Nonparametric option pricing under shape restrictions. *Journal of Econometrics* 116, 9–47.
- Ait-Sahalia, Y., Lo, A.W., 1998. Nonparametric estimation of state-price densities implicit in financial asset prices. *Journal of Finance* 53, 499–547.
- Ait-Sahalia, Y., Lo, A.W., 2000. Nonparametric risk management and implied risk aversion. *Journal of Econometrics* 94, 9–51.
- Ait-Sahalia, Y., Wang, Y., Yared, F., 2000. Do option markets correctly price the probabilities of movement of the underlying asset? *Journal of Econometrics* 102, 67–110.
- Anděl, J., 1985. *Mathematical Statistics*. SNTL/Alfa, Prague (in Czech).
- Andersen, L.B.G., Brotherton-Ratcliffe, R., 1997. The equity option volatility smile: An implicit finite-difference approach. *Journal of Computational Finance* 1 (2), 5–37.
- Bondarenko, O., 2003. Estimation of risk-neutral densities using positive convolution approximation. *Journal of Econometrics* 116, 85–112.
- Breeden, D., Litzenberger, R., 1978. Prices of state-contingent claims implicit in option prices. *Journal of Business* 51, 621–651.
- Buehler, H., 2006. Expensive martingales. *Quantitative Finance* 6 (3), 207–218.
- Diebold, F.X., Gunther, T., Tay, A., 1998. Evaluating density forecasts, with applications to financial risk management. *International Economic Review* 39, 863–883.
- Dupire, B., 1994. Pricing with a smile. *RISK* 7 (1), 18–20.
- Fengler, M.R., 2005. *Semiparametric Modeling of Implied Volatility*. Springer, Heidelberg.
- Fengler, M.R., Härdle, W., Mammen, E., 2007. A dynamic semiparametric factor model for implied volatility string dynamics. *Journal of Financial Econometrics* 5 (2), 189–218.
- Hafner, R., Wallmeier, M., 2000. *The Dynamics of DAX Implied Volatilities*. University of Augsburg Working Paper. Available at SSRN: <http://ssrn.com/abstract=234829> or doi: 10.2139/ssrn.234829.
- Harrison, J., Pliska, S., 1981. Martingale and stochastic integral in the theory of continuous trading. *Stochastic Processes and their Applications* 11, 215–260.
- Hlávka, Z., Svojk, M., 2008. Application of extended Kalman filter to SPD estimation. In: Härdle, W., Hautsch, N., Overbeck, L. (Eds.), *Applied Quantitative Finance*. Springer, Berlin, pp. 233–247.
- Huynh, K., Kervella, P., Zheng, J., 2002. Estimating state-price densities with nonparametric regression. In: Härdle, W., Kleinow, T., Stahl, G. (Eds.), *Applied Quantitative Finance*. Springer, Heidelberg, pp. 171–196.
- Jackwerth, J.C., 1999. Option-implied risk-neutral distributions and implied binomial trees: A literature review. *Journal of Derivatives* 7, 66–82.
- Kahalé, N., 2004. An arbitrage-free interpolation of volatilities. *RISK* 17 (5), 102–106.
- Rao, C.R., 1973. *Linear Statistical Inference and Its Applications*. Wiley, New York.
- Renault, E., 1997. Econometric models of option pricing errors. In: Kreps, D.M., Wallis, K.F. (Eds.), *Advances in Economics and Econometrics: Theory and Applications*, Seventh World Congress, vol. III. Cambridge University Press, Cambridge, pp. 223–278.
- Robertson, T., Wright, F.T., Dykstra, R.L., 1988. *Order Restricted Statistical Inference*. Wiley, Chichester.
- Seber, G.A.F., Wild, C.J., 2003. *Nonlinear Regression*. Wiley, Hoboken, New Jersey.
- Serfling, R., 1980. *Approximation Theorems of Mathematical Statistics*. Wiley, New York.
- Stoll, H.R., 1969. The relationship between put and call option prices. *Journal of Finance* 24, 801–824.
- Yatchew, A., Härdle, W., 2006. Nonparametric state price density estimation using constrained least squares and the bootstrap. *Journal of Econometrics* 133 (2), 579–599.

Variable Selection and Oversampling in the Use of Smooth Support Vector Machines for Predicting the Default Risk of Companies

WOLFGANG HÄRDLE,¹ YUH-JYE LEE,²
DOROTHEA SCHÄFER^{3*} AND YI-REN YE²

¹ CASE, Humboldt University, Berlin, Germany

² Department of Computer Science Information Engineering,
National Taiwan University of Science and Technology, Taipei,
Taiwan

³ German Institute of Economic Research, Berlin, Germany

ABSTRACT

In the era of Basel II a powerful tool for bankruptcy prognosis is vital for banks. The tool must be precise but also easily adaptable to the bank's objectives regarding the relation of false acceptances (Type I error) and false rejections (Type II error). We explore the suitability of smooth support vector machines (SSVM), and investigate how important factors such as the selection of appropriate accounting ratios (predictors), length of training period and structure of the training sample influence the precision of prediction. Moreover, we show that oversampling can be employed to control the trade-off between error types, and we compare SSVM with both logistic and discriminant analysis. Finally, we illustrate graphically how different models can be used jointly to support the decision-making process of loan officers. Copyright © 2008 John Wiley & Sons, Ltd.

KEY WORDS insolvency prognosis; support vector machines; statistical learning theory; non-parametric classification

INTRODUCTION

Default prediction is at the core of credit risk management and has therefore always attracted special attention. It has become even more important since the Basel Committee on Banking Supervision (Basel II) established borrowers' rating as the crucial criterion for minimum capital requirements of banks. The methods for generating rating figures have developed significantly over the last 10 years (Krahn and Weber, 2001). The rationale behind the increased sophistication in predicting borrowers' default risk is the aim of banks to minimize their cost of capital and to mitigate their own bankruptcy risks.

*Correspondence to: Dorothea Schäfer, German Institute for Economic Research (DIW) Berlin, Mohrenstrasse 58, 10117 Berlin, Germany. E-mail: dschaefer@diw.de

In this paper we intend to contribute to the increasing sophistication by exploring the predicting power of smooth support vector machines (SSVM). SSVM are a variant of the conventional support vector machines (SVM). The working principle of SVM in general can be described very easily. Imagine a group of observations in distinct classes such as balance sheet data from solvent and insolvent companies. Assume that the observations are such that they cannot be separated by a linear function. Rather than fitting nonlinear curves to the data, SVM handle this problem by using a specific transformation function—the kernel function—that maps the data from the original space into a higher-dimensional space where a hyperplane can do the separation linearly. The constrained optimization calculus of SVM gives a unique optimal separating hyperplane and adjusts it in such a way that the elements of distinct classes possess the largest distance to the hyperplane. By re-transforming the separating hyperplane into the original space of variables, the typical nonlinear separating function emerges (Vapnik, 1995). The main difference between SSVM and SVM is the following: the SSVM technique formulates the problem as an unconstrained minimization problem. This formulation has mathematical properties such as strong convexity and desirable infinite differentiability.

Our aim is threefold when using SSVM. Firstly, we examine the power of the SSVM in predicting company defaults; secondly, we investigate how important factors that are exogenous to the model, such as selecting the appropriate set of accounting ratios, length of training period and structure of the training sample, influence the precision; and thirdly, we explore how oversampling and downsampling affect the trade-off between Type I and Type II errors. In addition, we illustrate graphically how loan officers can benefit from jointly considering the prediction results of different SSVM variants and different models.

There are basically three distinct approaches in predicting the risk of default: option theory-based approaches, parametric models and non-parametric methods. While the first class relies on the rule of no arbitrage, the latter two are based purely on statistic principles. The popular (Merton, 1974) model treats the company's equity as the underlying asset of a call option held by shareholders. In case of insolvency shareholders deny exercising. The probability of default is derived from an adapted Black–Scholes formula. Later, several authors (e.g., Longstaff and Schwartz, 1995; Mella-Barral and Perraudin, 1997; Leland and Toft, 1996; Zhou, 2001; to name only a few) proposed variations to ease the strict assumptions on the structure of the data imposed by the Merton model. These approaches are frequently denoted as structural models. However, the most challenging requirement is the knowledge of market values of debt and equity. This precondition is a severe obstacle to using the Merton model adequately as it is only satisfied in a minority of cases.

Parametric statistical models can be applied to any type of data, whether they are market based or book based. The first model introduced was discriminant analysis (DA) for univariate (Beaver, 1966) and multivariate models (Altman, 1968). After DA usage of the logit and probit approach for predicting default was proposed in Martin (1977) and Ohlson (1980). These approaches rely on the a priori assumed functional dependence between risk of default and predictor. DA requires a linear functional dependence, or a pre-shaped polynomial functional dependence in advanced versions. Logit and probit tools work with monotonic relationships between default event and predictors such as accounting ratios. However, such restrictions often fail to meet the reality of observed data. This fact makes it clear that there is a need for an approach that, in contrast to conventional methods, relaxes the requirements on data and/or lowers the dependence on heuristics. Semi-parametric models as in Hwang *et al.* (2007) are between conventional linear models and non-parametric approaches. Nonlinear classification methods such as support vector machines (SVM) or neural networks are even stronger candidates to meet these demands as they go beyond conventional

discrimination methods. Tam and Kiang (1992) and Altman *et al.* (1994) focus on neural networks. In contrast, we concentrate on SVM exclusively.

The SVM method is a relatively new technique and builds on the principles of statistical learning theory. It is easier to handle compared to neural networks. Furthermore, SVM have a wider scope of application as the class of SVM models includes neural networks (Schölkopf and Smola, 2002). The power of SVM technology becomes evident in a situation as depicted in Figure 1 where operating profit margin and equity ratio are used as explanatory variables. A separating function similar to a parabola (in black) appears in the two-dimensional space. The accompanying light-grey lines represent the margin boundaries whose shape and location determine the distance of elements from the separating function. In contrast, the logit approach and discriminant DA yield the (white) linear separating function (Härdle *et al.*, 2007a).

Selecting the best accounting ratios for executing the task of predicting is an important issue in practice but has not received appropriate attention in research. We address this issue of how important the chosen set of predictors is for the outcome. For this purpose we explore the prediction potential of SSVM within a two-step approach. First, we derive alternative sets of accounting ratios that are used as predictors. The benchmark set comes from Chen *et al.* (2006). A second set is defined by a 1-norm SVM, and the third set is based on the principle of adding only those variables that contain the most contrary information with respect to an initial set that is a priori chosen. We call the latter procedure the incremental forward selection of variables. As a result we are working with three variants of SSVM. In the second step, these variants are compared with respect to their prediction power. We also compare SSVM with two traditional methods: the logit model and linear discriminant analysis.

The analysis is built on 28 accounting ratios of 20,000 solvent and 1000 insolvent German companies. Our findings show that the different SSVM types have an overall good performance with the means of correct predictions ranging from 70% to 78%. The SSVM on the basis of incremental

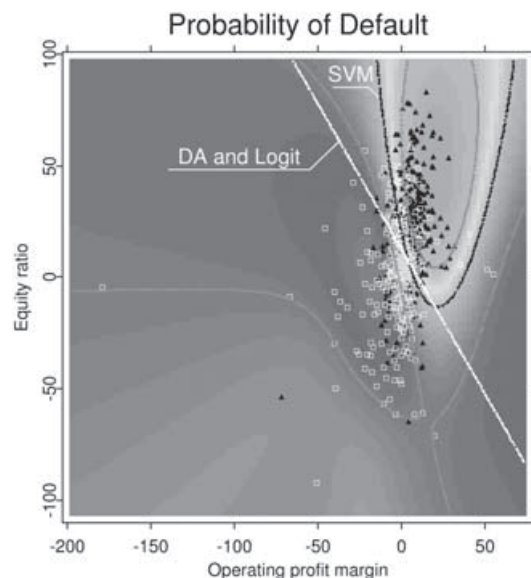


Figure 1. SVM-separating function (black) with margin in a two-dimensional space

forward selection clearly outperform the SSVM based on predictors selected by the 1-norm SVM. It is also found that oversampling influences the trade-off between Type I and Type II errors. Thus, oversampling can be used to make the relation of the two error types an issue of bank policy.

The rest of the paper is organized as follows. The following two sections describe the data, performance measures and SVM methodology. In the fourth section the variable selection technique and outcome are explained. The fifth section presents the experimental settings, estimation procedure and findings, and illustrates selected results. The sixth section concludes.

DATA AND MEASURES OF ACCURACY

In this study of the potential virtues of SVM in insolvency prognosis the CreditReform database is employed. The database consists of 20,000 financially and economically solvent and 1000 insolvent German companies observed once in the period from 1997 to 2002. Although the companies were randomly selected, accounting information dates most frequently in 2001 and 2002. Approximately 50% of the observations come from this period. The industry distribution of the insolvent companies is as follows: manufacturing 25.7%, wholesale and retail trade 20.1%, real estate 9.4%, construction 39.7% and others 5.1%. The latter includes businesses in agriculture, mining, electricity, gas and water supply, transport and communication, financial intermediation social service activities and hotels and restaurants. The 20,000 solvent companies belong to manufacturing (27.4%), wholesale and retail trade (24.8%), real estate (16.9%), construction (13.9%) and others (17.1%). There is only low coincidence between the industries represented in the insolvent and the solvent group of 'others'. The latter comprises many companies in industries such as publication administration and defense, education and health. Figure 2 shows the distribution of solvent and insolvent companies across industries. A set of balance sheet and income statement items describes each company. The ones we use for further analysis are described below:

- AD (amortization and depreciation)
- AP (accounts payable)
- AR (account receivable)

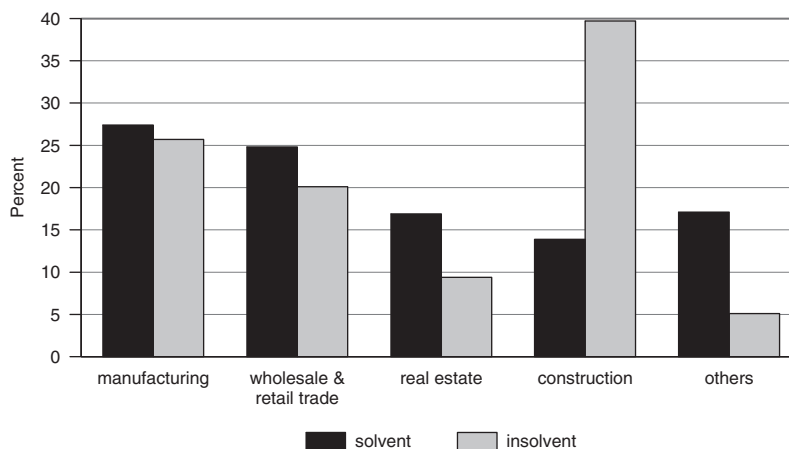


Figure 2. The distribution of solvent and insolvent companies across industries

- CA (current assets)
- CASH (cash and cash equivalents)
- CL (current liabilities)
- DEBT (debt)
- EBIT (earnings before interest and tax)
- EQUITY (equity)
- IDINV (growth of inventories)
- IDL (growth of liabilities)
- INTE (interest expense)
- INV (inventories)
- ITGA (intangible assets)
- LB (lands and buildings)
- NI (net income)
- OI (operating income)
- QA (quick assets)
- SALE (sales)
- TA (total assets)
- TL (total liabilities)
- WC (working capital (= CA – CL))

The companies appear in the database several times in different years; however, each year of balance sheet information is treated as a single observation. The data of the insolvent companies were collected 2 years prior to insolvency. The company sizes are measured by total assets. We construct 28 ratios to condense the balance sheet information (see Table I). However, before dealing with the CreditReform dataset, some companies whose behavior is very different from other ones are filtered out in order to make the dataset more compact. The data pre-processing procedure is described as follows:

1. We excluded companies whose total assets were not in the range of 10^5 – 10^7 EUR (remaining insolvent: 967; solvent: 15,834).
2. In order to compute the accounting ratios AP/SALE, OI/TA, TL/TA, CASH/TA, IDINV/INV, INV/SALE, EBIT/TA and NI/SALE, we have removed companies with zero denominators (remaining insolvent: 816; solvent 11,005).
3. We dropped outliers, that is, in the insolvent class companies with extreme values of financial indices have been removed (remaining insolvent: 811; solvent: 10,468).

After pre-processing, the dataset consists of 11,279 companies (811 insolvent and 10,468 solvent). In the following analysis, we focus on the revised dataset.

The performance of the SSVM is evaluated on the basis of three measures of accuracy: Type I error rate (%), Type II error rate (%) and total error rate (%). The Type I error is the ratio of the number of insolvent companies predicted as solvent ones to the number of insolvent companies. The Type II error is the ratio of the number of solvent companies predicted as insolvent ones to the number of solvent companies. Accordingly, the error-type rates (in percentage) are defined as follows

- Type I error rate = $FN/(FN + TP) \times 100$ (%);
- Type II error rate = $FP/(FP + TN) \times 100$ (%);
- Total error rate = $(FN + FP)/(TP + TN + FP + FN) \times 100$ (%);

Table I. Definitions of accounting ratios used in the analysis

Variable	Ratio	Indicator for
X1	NI/TA	Profitability
X2	NI/SALE	Profitability
X3	OI/TA	Profitability
X4	OI/SALE	Profitability
X5	EBIT/TA	Profitability
X6	(EBIT + AD)/TA	Profitability
X7	EBIT/SALE	Profitability
X8	EQUITY/TA	Leverage
X9	(EQUITY-ITGA)/ (TA-ITGA-CASH-LB)	Leverage
X10	CL/TA	Leverage
X11	(CL-CASH)/TA	Leverage
X12	TL/TA	Leverage
X13	DEBT/TA	Leverage
X14	EBIT/INTE	Leverage
X15	CASH/TA	Liquidity
X16	CASH/CL	Liquidity
X17	QA/CL	Liquidity
X18	CA/CL	Liquidity
X19	WC/TA	Liquidity
X20	CL/TL	Liquidity
X21	TA/SALE	Activity
X22	INV/SALE	Activity
X23	AR/SALE	Activity
X24	AP/SALE	Activity
X25	Log(TA)	Size
X26	IDINV/INV	Growth
X27	IDL/TL	Growth
X28	IDCASH/CASH	Growth

where

True positive (TP): Predict insolvent companies as insolvent ones

False positive (FP): Predict solvent companies as insolvent ones

True negative (TN): Predict solvent companies as solvent ones

False negative (FN): Predict insolvent companies as solvent ones

The following matrix explains the terms used in the definition of error rates:

		Predicted class	
		Positive	Negative
Actual Class	Positive	<i>True positive (TP)</i>	<i>False negative (FN)</i>
	Negative	<i>False positive (FP)</i>	<i>True negative (TN)</i>

SVM METHODOLOGY

In recent years, the so-called support vector machines (SVM), which have their roots in the theory of statistical learning (Burges, 1998; Christianini and Shawe-Taylor, 2000; Vapnik, 1995) have

become one of the most successful learning algorithms for classification as well as for regression (Drucker *et al.*, 1997; Mangasarian and Musicant, 2000; Smola and Schölkopf, 2004). Some features of SVM make them particularly attractive for predicting the default risk of companies. SVM are a non-parametric technique that learn the separating function from the data; they are based on a sound theoretical concept, do not require a particular distribution of the data, and deliver an optimal solution for the expected loss from misclassification. SVM estimate the separating hyperplane between defaulting and non-defaulting companies under the constraint of a maximal margin between the two classes (Vapnik, 1995; Schölkopf and Smola, 2002).

SVM can be formulated differently. However, in all variants either a constrained minimization problem or an unconstrained minimization problem is solved. The objective function in these optimization problems basically consists of two parts: a misclassification penalty part which stands for *model bias* and a regularization part which controls the *model variance*. We briefly introduce three different models: the smooth support vector machines (SSVM) (Lee and Mangasarian, 2001), the smooth support vector machines with reduced kernel technique (RSVM) and the 1-norm SVM. The SSVM will be used for classification and the 1-norm SVM will be employed for variable selection. The RSVM are applied for oversampling in order to mitigate the computational burden due to increasing the number of instances in the training sample.

Smooth support vector machines

The aim of the SVM technique is to find the separating hyperplane with the largest margin from the training data. This hyperplane is ‘optimal’ in the sense of statistical learning: it strikes a balance between overfitting and underfitting. Overfitting means that the classification boundary is too curved and therefore has less ability to classify unseen data correctly. Underfitting, on the other hand, gives a too simple classification boundary and leaves too many misclassified observations (Vapnik, 1995). We begin with linear support vector machines. Given a training dataset $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \subseteq \mathbb{R}^d \times \mathbb{R}$, where $\mathbf{x}_i \in \mathbb{R}^d$ is the input data and $y_i \in \{-1, 1\}$ is the corresponding class label, a conventional SVM separating hyperplane is generated by solving a convex optimization problem given as follows:

$$\begin{aligned} \min_{(w, b, \xi) \in \mathbb{R}^{d+1+n}} \quad & C \sum_{i=1}^n \xi_i + \frac{1}{2} \|w\|_2^2 \\ \text{s.t.} \quad & y_i(w^\top \mathbf{x}_i + b) + \xi_i \geq 1 \\ & \xi_i \geq 0, \quad \text{for } i = 1, 2, \dots, n \end{aligned} \quad (1)$$

where C is a positive parameter controlling the trade-off between the training error (model bias) and the part of maximizing the margin (model variance) that is achieved by minimizing $\|w\|_2^2$. In contrast to the conventional SVM of (1), smooth support vector machines minimize the square of the slack vector ξ with weight $\frac{C}{2}$. In addition, the SSVM methodology appends $\frac{b^2}{2}$ to the term that is to be minimized. This expansion results in the following minimization problem:

$$\begin{aligned} \min_{(w, b, \xi) \in \mathbb{R}^{d+1+n}} \quad & \frac{C}{2} \sum_{i=1}^n \xi_i^2 + \frac{1}{2} (\|w\|_2^2 + b^2) \\ \text{s.t.} \quad & y_i(w^\top \mathbf{x}_i + b) + \xi_i \geq 1 \\ & \xi_i \geq 0, \quad \text{for } i = 1, 2, \dots, n \end{aligned} \quad (2)$$

In a solution of (2), ξ is given by $\xi_i = \{1 - y_i(w^\top \mathbf{x}_i + b)\}_+$ for all i where the *plus* function x_+ is defined as $x_+ = \max\{0, x\}$. Thus, we can replace ξ_i in (2) by $\{1 - y_i(w^\top \mathbf{x}_i + b)\}_+$. This will convert the problem (2) into an unconstrained minimization problem as follows:

$$\min_{(w, b) \in \mathbb{R}^{d+1}} \frac{C}{2} \sum_{i=1}^n \{1 - y_i(w^\top \mathbf{x}_i + b)\}_+^2 + \frac{1}{2} (\|w\|_2^2 + b^2) \tag{3}$$

This formulation reduces the number of variables from $d + 1 + n$ to $d + 1$. However, the objective function to be minimized is not twice differentiable, which precludes the use of a fast Newton method. In the SSVM, the plus function x_+ is approximated by a smooth *p-function*, $p(x, \alpha) = x + \frac{1}{\alpha} \log(1 + e^{-\alpha x})$, $\alpha > 0$. Replacing the plus function with a very accurate smooth approximation *p-function* gives the smooth support vector machine formulation:

$$\min_{(w, b) \in \mathbb{R}^{d+1}} \frac{C}{2} \sum_{i=1}^n p(\{1 - y_i(w^\top \mathbf{x}_i + b)\}, \alpha)^2 + \frac{1}{2} (\|w\|_2^2 + b^2) \tag{4}$$

where $\alpha > 0$ is the smooth parameter. The objective function in problem (4) is strongly convex and infinitely differentiable. Hence, it has a unique solution and can be solved by using a fast Newton–Armijo algorithm. For the nonlinear case, this formulation can be extended to the nonlinear SVM by using the kernel trick as follows:

$$\min_{(u, b) \in \mathbb{R}^{n+1}} \frac{C}{2} \sum_{i=1}^n p\left(\left[1 - y_i \left\{ \sum_{j=1}^n u_j K(\mathbf{x}_i, \mathbf{x}_j) + b \right\}\right], \alpha\right)^2 + \frac{1}{2} (\|u\|_2^2 + b^2) \tag{5}$$

where $K(\mathbf{x}_i, \mathbf{x}_j)$ is a kernel function. This kernel function represents the inner product of $\phi(\mathbf{x}_i)$ and $\phi(\mathbf{x}_j)$, where ϕ is a certain mapping from input space \mathbb{R}^d to a feature space \mathcal{F} . We do not need to know the mapping of ϕ explicitly. This is the so-called kernel trick. The nonlinear SSVM classifier can be expressed in matrix form as follows:

$$\sum_{u_j \neq 0} u_j K(A_j^\top, \mathbf{x}) + b = K(\mathbf{x}, A^\top)u + b \tag{6}$$

where $A = [\mathbf{x}_1^\top; \dots; \mathbf{x}_n^\top]$ and $A_j = \mathbf{x}_j^\top$.

Reduced support vector machine

In large-scale problems, the full kernel matrix will be very large so it may not be appropriate to use the full kernel matrix when dealing with (5). In order to avoid facing such a big full kernel matrix, we brought in the reduced kernel technique (Lee and Huang, 2007). The key idea of the reduced kernel technique is to randomly select a portion of data and to generate a thin rectangular kernel matrix, then to use this much smaller rectangular kernel matrix to replace the full kernel matrix. In the process of replacing the full kernel matrix by a reduced kernel, we use the Nyström approximation (Smola and Schölkopf, 2000) for the full kernel matrix:

$$K(A, A^\top) \approx K(A, \tilde{A}^\top)K(\tilde{A}, \tilde{A}^\top)^{-1}K(\tilde{A}, A^\top) \tag{7}$$

where $K(A, A^\top) = K_{n \times n}$, $\tilde{A}_{\tilde{n} \times d}$ is a subset of A and $K(A, \tilde{A}) = \tilde{K}_{n \times \tilde{n}}$ is a reduced kernel. Thus, we have

$$K(A, A^\top)u \approx K(A, \tilde{A}^\top)K(\tilde{A}, \tilde{A}^\top)^{-1}K(\tilde{A}^\top, A)u = K(A, \tilde{A}^\top)\tilde{u} \quad (8)$$

where $\tilde{u} \in \mathbb{R}^{\tilde{n}}$ is an approximated solution of u via the reduced kernel technique. The reduced kernel method constructs a compressed model and cuts down the computational cost from $\mathcal{O}(n^3)$ to $\mathcal{O}(\tilde{n}^3)$. It has been shown that the solution of reduced kernel matrix approximates the solution of full kernel matrix well. The SSVM with the reduced kernel are called RSVM.

1-Norm support vector machine

The 1-norm support vector machine replaces the regularization term $\|w\|_2^2$ in (1) with the ℓ_1 -norm of w . The ℓ_1 -norm regularization term is also called the LASSO penalty (Tibshirani, 1996). It tends to shrink the coefficients w 's towards zeros in particular for those coefficients corresponding to redundant noise features (Zhu *et al.*, 2003; Williams and Seeger, 2001). This nice feature will lead to a way of selecting the important ratios in our prediction model. The formulation of 1-norm SVM is described as follows:

$$\begin{aligned} \min_{(w, b, \xi) \in \mathbb{R}^{d+1+n}} \quad & C \sum_{i=1}^n \xi_i + \|w\|_1 \\ \text{s.t.} \quad & y_i(w^\top \mathbf{x}_i + b) + \xi_i \geq 1 \\ & \xi_i \geq 0, \quad \text{for } i = 1, 2, \dots, n. \end{aligned} \quad (9)$$

The objective function of (9) is a piecewise linear convex function. We can reformulate it as the following linear programming problem:

$$\begin{aligned} \min_{(w, s, b, \xi) \in \mathbb{R}^{d+d+1+n}} \quad & C \sum_{i=1}^n \xi_i + \sum_{j=1}^d s_j \\ \text{s.t.} \quad & y_i(w^\top \mathbf{x}_i + b) + \xi_i \geq 1 \\ & -s_j \leq w_j \leq s_j, \quad \text{for } j = 1, 2, \dots, d, \\ & \xi_i \geq 0, \quad \text{for } i = 1, 2, \dots, n \end{aligned} \quad (10)$$

where s_j is the upper bound of the absolute value of w_j . In the optimal solution of (10) the sum of s_j is equal to $\|w\|_1$.

The 1-norm SVM can generate a very sparse solution w and lead to a parsimonious model. In a linear SVM classifier, solution sparsity means that the separating function $f(\mathbf{x}) = w^\top \mathbf{x} + b$ depends on very few input attributes. This characteristic can significantly suppress the number of nonzero coefficient w 's, especially when there are many redundant noise features (Fung and Mangasarian, 2004; Zhu *et al.*, 2003). Therefore the 1-norm SVM can be a very promising tool for the variable selection tasks. We will use it to choose the important financial indices for our bankruptcy prognosis model.

SELECTION OF ACCOUNTING RATIOS

In principle any possible combination of accounting ratios could be used as explanatory variables in a bankruptcy prognosis model. Therefore, appropriate performance measures are needed to gear the process of variable selection towards picking the ratios with the highest separating power. In

Chen *et al.* (2006) accuracy ratio (AR) and conditional information entropy ratio (CIER) determine the selection procedure's outcome. It turned out that the ratio 'accounts payable divided by sales', X24 (AP/SALE), has the best performance values for a univariate SVM model. The second selected variable was the one combined with X24 that had the best performance in a bivariate SVM model. This is the analogue of forward selection in linear regression modeling. Typically, improvement declines if new variables are added consecutively. In Chen *et al.* (2006) the performance indicators started to decrease after the model included eight variables. The described selection procedure is quite lengthy, since there are at least 216 accounting ratio combinations to be considered. We will not employ the procedure here but use the chosen set of eight variables as the benchmark set V1. Table II presents V1 in the first column.

We propose two different approaches for variable selection that will simplify the selection procedure. The first one is based on 1-norm SVM introduced above. The SVM were applied to the period from 1997 to 1999. We selected the variables according to the size of the absolute values of the coefficients w from the solution of the 1-norm SVM. Table II displays the eight selected variables as V2. We obtain eight variables out of 28. Note that five variables, X2, X3, X5, X15 and X24, are also in the benchmark set V1.

The second variable selection scheme is incremental forward variable selection. The intuition behind this scheme is that a new variable will be added into the already selected set, if it brings in the most extra information. We measure the extra information for an accounting ratio using the distance between this new ratio vector and the space spanned by the current selected ratio subset. This distance can be computed by solving a least-squares problem (Lee *et al.*, 2008). The ratio with the farthest distance will be added into the selected accounting ratio set. We repeat this procedure until a certain stopping criterion is satisfied. The accounting ratio X24 (AP/SALE) is used as the initial selected accounting ratio. Then we follow the procedure seven times to select seven more extra accounting ratios. The variable set generated is called V3. We will use these three variable sets, V1, V2 and V3, for further data analysis in the next section. The symbol $+$ denotes the variables that are common to all sets: X2, X3, X5 and X24.

Table II. Selected variables

Variable	Definition	V1	V2	V3
X2 ⁺	NI/SALE	x	x	x
X3 ⁺	OI/TA	x	x	x
X4	OI/SALE			x
X5 ⁺	EBIT/TA	x	x	x
X6	(EBIT + AD)/TA		x	
X7	EBIT/SALE			x
X8	EQUITY/TA		x	
X12	TL/TA	x		
X13	DEBT/TA			x
X15	CASH/TA	x	x	
X21	TA/SALE			x
X22	INV/SALE	x		
X23	AR/SALE		x	
X24 ⁺	AP/SALE	x	x	x
X26	IDINV/INV	x		

EXPERIMENTAL SETTING AND RESULTS

In this section we present our experimental setting and results. We compare the performance of three sets of accounting ratios, V1, V2 and V3, in our SSVM-based insolvency prognosis model. The performance is measured by Type I error rate, Type II error rate and total error rate. Fortunately, in reality, there is only a small number of insolvent companies compared to the number of solvent companies. Due to the small share in a sample that reflects reality, a simple classification such as naive Bayesian or a decision tree tends to classify every company as solvent. Such a classification would imply accepting all companies' loan applications and would thus lead to a very high Type I error rate while the total error rate and the Type II error rate are very small. Such models are useless in practice.

Our cleaned dataset consists of around 10% of insolvent companies. Thus, the sample is fairly unbalanced although the share of insolvent companies is higher than in reality. In order to deal with this problem, insolvency prognosis models usually start off with more balanced training and testing samples than reality can provide. For example, Härdle *et al.* (2007b) employ a downsampling strategy and work with balanced (50%/50%) samples. The chosen bootstrap procedure repeatedly randomly selects a fixed number of insolvent companies from the training set and adds the same number of randomly selected solvent companies. However, in this paper we adopt an oversampling strategy, to balance the size between the solvent and the insolvent companies, and refer to the downsampling procedure primarily for reasons of reference.

Oversampling duplicates the number of insolvent companies a certain number of times. In this experiment, we duplicate in each scenario the number of insolvent companies as many times as necessary to reach a balanced sample. Note that in our oversampling scheme every solvent and insolvent company's information is utilized. This increases the computational burden due to increasing the number of training instances. We employ the reduced kernel technique introduced above to mediate this problem.

All classifiers we need in these experiments are reduced SSVM with the Gaussian kernel, which is defined as

$$K(\mathbf{x}, \mathbf{z}) = e^{-\gamma \|\mathbf{x} - \mathbf{z}\|_2^2}$$

where γ is the width parameter. In nonlinear SSVM, we need to determine two parameters: the penalty term C and γ . The 2D grid search will consume a lot of time. In order to cut down the search time, we adopt the uniform design model selection method (Huang *et al.*, 2007) to search an appropriate pair of parameters.

Performance of SSVM

We conduct the experiments in a scenario in which we always train the SSVM bankruptcy prognosis model from the data at hand and then use the trained SSVM to predict the following year's cases. This strategy simulates the real task of prediction which binds the analyst to use past data for forecasting future outcomes. The experimental setting is described in Table III. The number of periods which enter the training set changes from 1 year (S1) to 5 years (S5).

In Tables IV and V we report the results for the oversampling and downsampling strategy respectively. Mean and standard deviation of Type I, Type II and total error rates (misclassification rates) are shown. We perform these experiments for the three variable sets, V1 to V3, and compare the oversampling and downsampling scheme in each experiment. All experiments are repeated 30 times

Table III. The scenario of our experiments

Scenario	Observation period of training set	Observation period of testing set
S1	1997	1998
S2	1997–1998	1999
S3	1997–1999	2000
S4	1997–2000	2001
S5	1997–2001	2002

Table IV. Results of oversampling for three variable sets (RSVM)

Set of accounting ratios	Scenario	Type I error rate		Type II error rate		Total error rate	
		Mean	SD	Mean	SD	Mean	SD
V1	S1	33.16	0.55	26.15	0.13	26.75	0.12
	S2	31.58	0.01	29.10	0.07	29.35	0.07
	S3	28.11	0.73	26.73	0.16	26.83	0.16
	S4	30.14	0.62	25.66	0.17	25.93	0.15
	S5	24.24	0.56	23.44	0.13	23.48	0.13
V2	S1	29.28	0.92	27.20	0.24	27.38	0.23
	S2	28.20	0.29	30.18	0.18	29.98	0.16
	S3	27.41	0.61	29.67	0.19	29.50	0.17
	S4	28.12	0.74	28.32	0.19	28.31	0.15
	S5	23.91	0.62	24.99	0.10	24.94	0.10
V3	S1	29.28	0.83	25.11	0.25	25.46	0.21
	S2	31.27	0.62	29.79	0.34	29.94	0.35
	S3	30.91	0.13	27.21	0.19	27.48	0.18
	S4	32.00	0.54	25.19	0.17	25.61	0.14
	S5	26.98	0.42	22.90	0.11	23.08	0.11

Table V. Results of downsampling for three variable sets (SSVM with Gaussian kernel)

Set of accounting ratios	Scenario	Type I error rate		Type II error rate		Total error rate	
		Mean	SD	Mean	SD	Mean	SD
V1	S1	32.20	3.12	28.98	1.70	29.26	1.46
	S2	29.74	2.29	28.77	1.97	28.87	1.57
	S3	30.46	1.88	26.23	1.33	26.54	1.17
	S4	31.55	1.52	23.89	0.97	24.37	0.87
	S5	28.81	1.53	23.09	0.73	23.34	0.69
V2	S1	29.94	2.91	28.07	2.15	28.23	1.79
	S2	28.77	2.58	29.80	1.89	29.70	1.52
	S3	29.88	1.88	27.19	1.32	27.39	1.19
	S4	29.06	1.68	26.26	1.00	26.43	0.86
	S5	26.92	1.94	25.30	1.17	25.37	1.06
V3	S1	30.87	3.25	26.61	2.45	26.98	2.11
	S2	33.31	2.16	28.60	2.01	29.08	1.65
	S3	31.82	1.52	26.41	1.45	26.80	1.31
	S4	35.0	2.13	24.29	0.77	24.96	0.68
	S5	30.66	1.60	21.92	0.96	22.30	0.92

because of the randomness in the experiments. The randomness is very obvious in the downsampling scheme (see Table V). Each time we only choose negative instances with the same size of the whole positive instances. The observed randomness in our oversampling scheme (Table IV) is due to applying the reduced kernel technique to solving the problem. We use the training set in the downsampling scheme as the reduced set. That is, we use all the insolvent instances and the equal number of solvent instances as our reduced set in generating the reduced kernel. Then we duplicate the insolvent part of the kernel matrix to balance the size of insolvent and solvent companies.

Both tables reveal that different variable selection schemes produce dissimilar results with respect to both precision and deviation of predicting. The oversampling scheme shows better results in the Type I error rate but has slightly bigger total error rates. It is also obvious that in almost all models a longer training period works in favor of accuracy of prediction. Clearly, the oversampling schemes have much smaller standard deviations in the Type I error rate, Type II error rate, and total error rate than the downsampling one. According to this observation, we conclude that the oversampling scheme will generate a more robust model than the downsampling scheme.

Figure 3 illustrates the development (learning curve) of the Type I error rate and total error rate with regard to variable set V3 for both oversampling and downsampling. The bullets on the lines

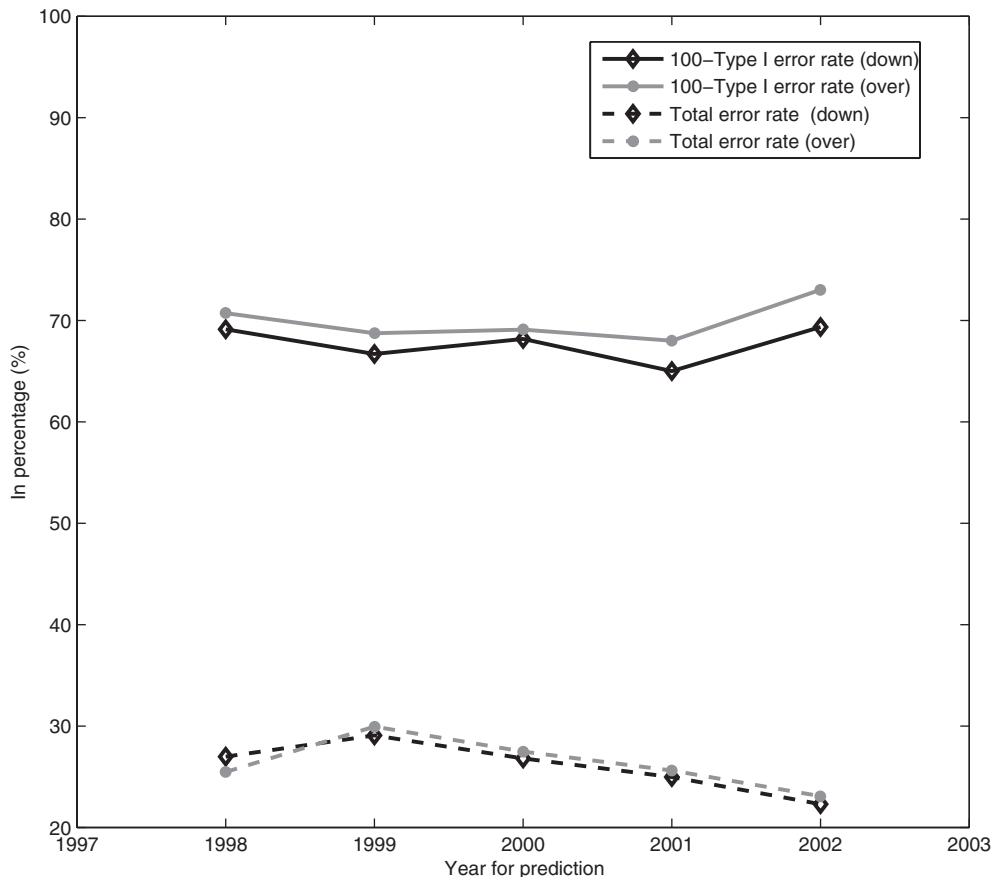


Figure 3. Learning curve for variables set V3

mark the different training scenarios. For example, the first bullets from the left represent S1 (training set from 1997, testing set from 1998), the second bullets illustrate S2 (training set from 1997 to 1998, testing set from 1999) etc. For the purpose of better visibility, the Type I error rate is only indirectly displayed as $100 - \text{Type I error rate}$. The upper solid line in gray represents the oversampling scheme and the black solid line the downsampling one. Note that the performance in terms of the Type I error rate is worse the higher the distance between the upper end of the diagram and the solid lines. The learning curve over the time frame the training sample covers shows an upward tendency between S1 and S5 for the number $100 - \text{Type I error rate}$. However, the curves are non-monotonic. There is a disturbance for the forecast of year 1999 that is based on training samples that cover 1997 to 1998, and also one for the forecast of year 2001 based on training samples covering 1997 to 2000. Both disturbances may have been caused by the reform of the German insolvency code that came into force in 1999. The most important objective of the reform was to allow for more company restructuring and less liquidation than before. This reform considerably changed the behavior of German companies towards declaring insolvency, and thus most likely the nature of balance sheets that are associated with insolvent companies.

The disturbances are less visible with respect to the overall performance. The dashed lines near the lower edge of the diagram box show total error rates, gray for the oversampling and black for the downsampling scheme. There is a clear tendency towards a lower total error rate from S2 to S5 for both schemes. The downsampling line is slightly below the oversampling one, representing a slightly better performance in terms of the mean of the total error rate. However, this result has to be seen in the light of the trade-off between magnitude and stability of results, as oversampling yields much more stable results. The standard deviations for V3 are only a small portion of the numbers generated by the downsampling procedure across all training scenarios (Tables IV and V).

Table VI presents the comparison between the sets by focusing on the total error rate. It indicates by an asterisk whether the differences in means are significant at the 10% level via t -test and, in addition, gives the set which is superior in the dual comparison. Variable set V2 is nearly absent in Table VI. Thus V2 is clearly outperformed by both sets V1 and V3. There is no clear distinction between V1 and V3 except for Scenario S5. Given the long training period V3 is superior in both the downsampling and oversampling scenarios and generates the lowest total error rate in absolute terms.

In order to investigate the effect of the oversampling versus the downsampling scheme we follow the setting as above, but we use the V3 variable set. For each training–test pair, we carry out oversampling for positive instances from 6 to 15 times. We show the trend and effect in Figure 4. It is

Table VI. Statistical significance in differences in means (10% level) between the three variable sets: total error

Sets	S1	S2	S3	S4	S5
<i>Oversampling</i>					
V1 vs. V2	V1*	V1*	V1*	V1*	V1*
V1 vs. V3	V3*	V1*	V1*	V3*	V3*
V2 vs. V3	V3*		V3*	V3*	V3*
<i>Downsampling</i>					
V1 vs. V2	V2*	V1*	V1*	V1*	V1*
V1 vs. V3	V3*			V1*	V3*
V2 vs. V3	V3*		V3*	V3*	V3*

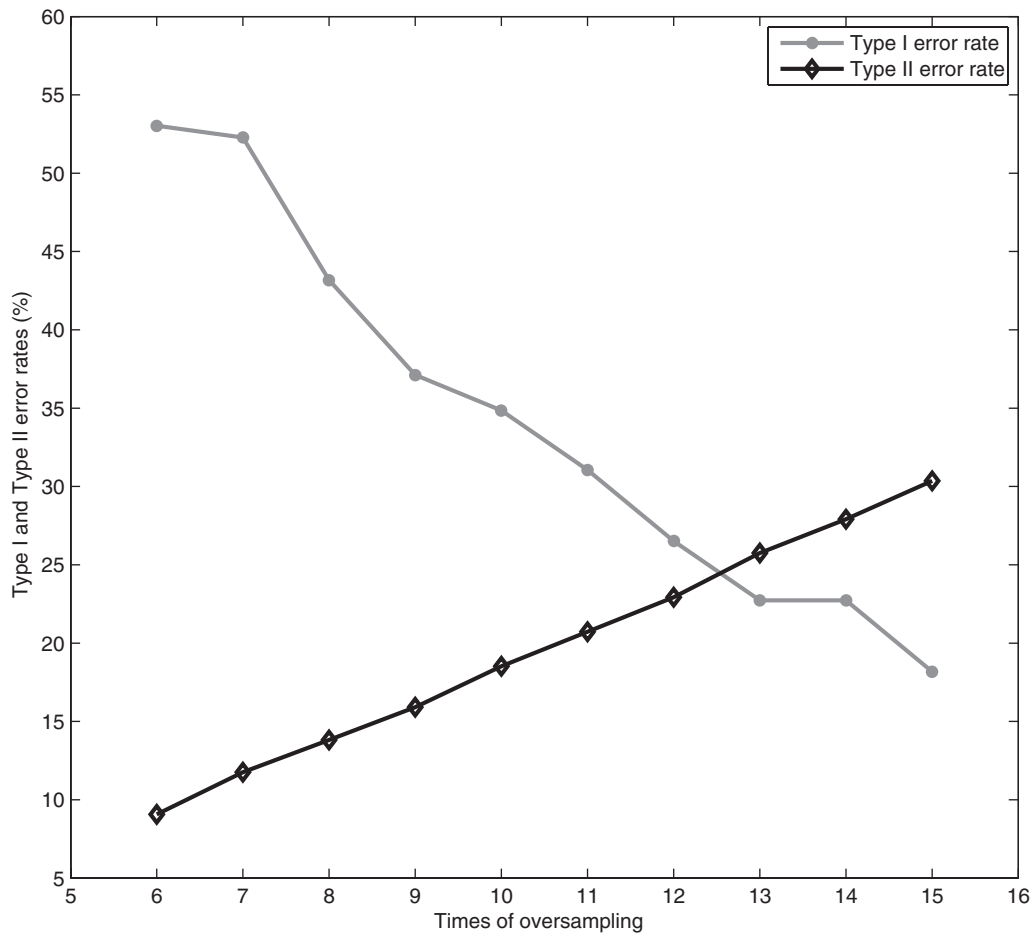


Figure 4. The effect of oversampling on Type I and Type II error rates for scenario S5 and variables set V3

easy to see that the Type I (II) error rate decreases (increases) as the oversampling times increase. This feature implies that the machine would have a tendency of classifying all companies as solvent if the training sample had realistic shares of insolvent and solvent companies. Such behavior would produce a Type I error rate of 100%. The more balanced the sample is, the higher the penalty for classifying insolvent companies as solvent. This fact is illustrated in Figure 4 by the decreasing curve with respect to the number of duplications of insolvent companies.

Often banks favor a strategy that allows them to minimize the Type II errors for a given number of Type I errors. The impact of oversampling on the trade-off between the two types of errors—shown in Figure 4—implies that the number of oversampling times is a strategic variable in training the machine. This number can be determined by the bank's aim regarding the relation of Type I and Type II errors.

Comparison with logit and linear discriminant analysis

The examination of SSVM is incomplete without comparing it to highly used traditional methods such as the logistic model (LM) and linear discriminant analysis (DA). Therefore, we replicate the research design of the previous section with both traditional models. In addition, we test whether the difference in means in the total error rate is statistically significant. The comparison of means with regard to the total error rate is presented in Tables VII and VIII for the oversampling and downsampling strategy respectively. Table IX summarizes the comparison of the approaches and displays the statistical significance of their mean differences. Asterisks indicate the out-performance

Table VII. Comparison of the total error rate (%) as generated by SSVM with LM and DA: oversampling for three variable sets

Set of accounting ratios	Scenario	SSVM	LM	DA
		Mean	Mean	Mean
V1	S1	26.75	26.50	25.60
	S2	29.35	28.96	27.22
	S3	26.83	28.94	27.42
	S4	25.93	26.20	25.55
	S5	23.48	26.95	28.23
V2	S1	27.38	26.80	26.20
	S2	29.98	28.63	28.70
	S3	29.50	29.52	29.46
	S4	28.31	28.43	28.08
	S5	24.94	29.22	31.42
V3	S1	25.46	25.07	23.65
	S2	29.94	28.29	27.02
	S3	27.48	27.89	25.84
	S4	25.61	26.60	24.85
	S5	23.08	25.32	26.15

Table VIII. Comparison of the total error rate (%) as generated by SSVM with LM and DA: downsampling for three variable sets

Set of accounting ratios	Scenario	SSVM	LM	DA
		Mean	Mean	Mean
V1	S1	29.26	26.86	27.34
	S2	28.87	28.62	28.26
	S3	26.54	27.54	28.22
	S4	24.37	24.80	25.47
	S5	23.34	24.81	25.86
V2	S1	28.23	27.28	28.62
	S2	29.70	29.29	29.65
	S3	27.39	28.56	29.58
	S4	26.43	26.41	27.96
	S5	25.37	26.52	29.69
V3	S1	26.98	26.03	25.47
	S2	29.08	28.04	27.22
	S3	26.80	26.60	26.51
	S4	24.96	25.25	25.44
	S5	22.30	23.96	24.31

Table IX. Statistical significance in differences of means (10% level) between SSVM and LM and SSVM and DA, respectively, for the sets V1 to V3: total error rate

V1	S1	S2	S3	S4	S5
<i>Oversampling</i>					
SSVM vs. LM			*	*	*
SSVM vs. DA			*		*
<i>Downsampling</i>					
SSVM vs. LM			*	*	*
SSVM vs. DA			*	*	*
V2	S1	S2	S3	S4	S5
<i>Oversampling</i>					
SSVM vs. LM				*	*
SSVM vs. DA					*
<i>Downsampling</i>					
SSVM vs. LM			*		*
SSVM vs. DA			*	*	*
V3	S1	S2	S3	S4	S5
<i>Oversampling</i>					
SSVM vs. LM			*	*	*
SSVM vs. DA					*
<i>Downsampling</i>					
SSVM vs. LM					*
SSVM vs. DA				*	*

of the logistic model or discriminant analysis by SSVMs at the 10% level via t -test. It is obvious that the SSVM technique yields the better results, the longer the period is from which the training observations are taken. In fact, the results show that the SSVM works significantly better than LM and DA in most cases in S3 to S5, with the clearest advantage for testing sets S4 and S5, where the accounting information of the predicted companies dates most frequently in 2001 and 2002.

We also investigate the effect of oversampling on LM and DA. We follow the same setting in the previous section, doing oversampling for positive instances from 6 to 15 times. Unlike the SSVM-based insolvency prognosis model, the DA approach is insensitive in both Type I and Type II error rates to the replication of positive instances. The result for DA is illustrated in Figure 5. The LM approach has very similar results to the SSVM model. We will not show the result here.

More data visualization

Each SSVM model has its own output value. We use this output to construct 2D coordinate systems. Figure 6 shows an example for scenario S5 where the scores of the SSVM_{V3} model (SSVM_{V1} model) are represented by the horizontal (vertical) line. A positive (negative) value indicates predicted insolvency (solvency). We then map all insolvent companies in the testing set onto the coordinate systems. There are 132 insolvent companies and 2866 solvent companies in this testing set. We also randomly choose the same amount of solvent companies from the testing set.

The plus points in the lower left quadrant and the circle points in the upper right quadrant show the number of Type I errors and Type II errors, respectively, in both models. Plus points in the upper right quadrant and circle points in the lower left quadrant reflect those companies that are predicted

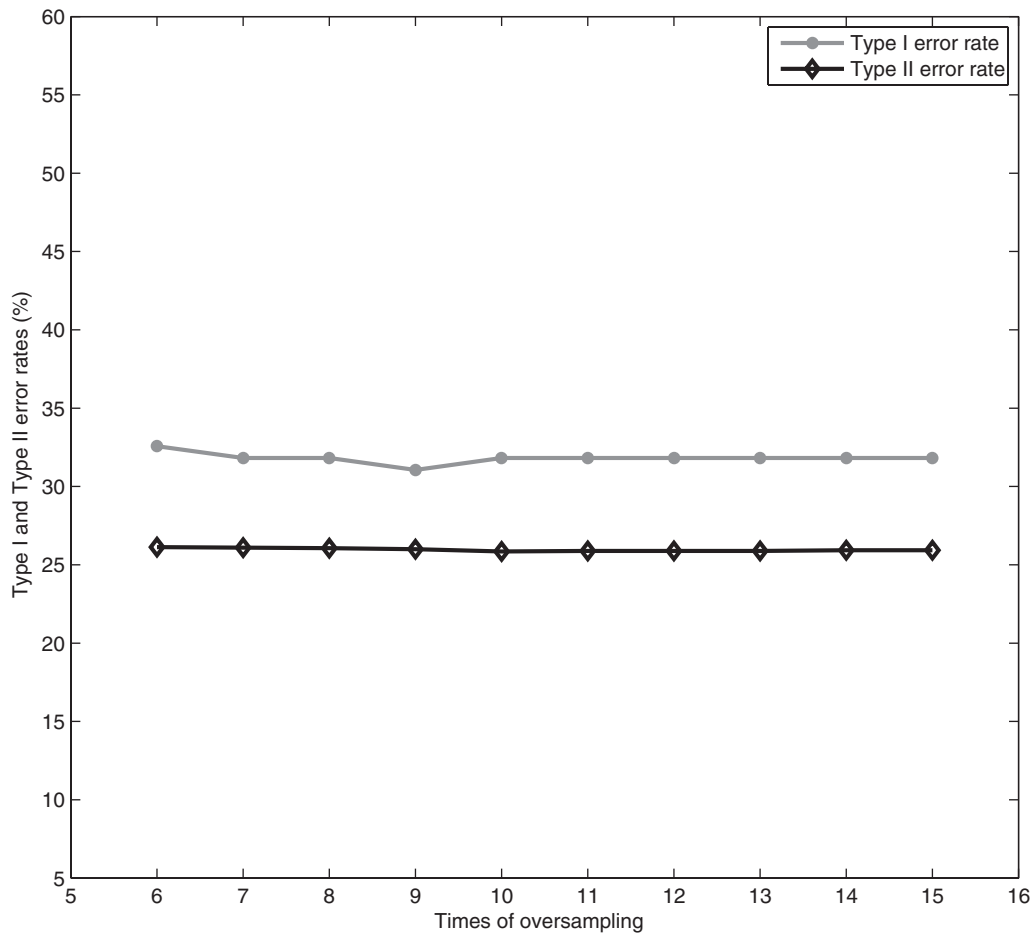


Figure 5. The effect of oversampling on Type I and Type II error rates for scenario S5 and variables set V3 in DA

correctly by both models. Circles and plus points in the lower right quadrant (upper left quadrant) represent conflicting prognoses. We also report the number of insolvent companies and the number of solvent companies in each quadrant of Figure 6. The two different insolvency prognosis models based on V1 and V3, respectively, can be considered as alternative experts. The two forecasts for each instance in the testing set is plotted in the diagram. The proposed visualization scheme could be used to support loan officers in their final decision on accepting or rejecting a client's application. Furthermore, this data visualization scheme can also be applied to two different learning algorithms, such as $SSVM_{V3}$ vs. LM_{V3} and $SSVM_{V3}$ vs. DA_{V3} . We show these data visualization plots in Figures 7 and 8. If the loan application has been classified as solvent or insolvent by alternative machines, it is most likely that the prognosis meets reality (the plus points in the upper right quadrant and the circle points in the lower left quadrant). Opposing forecasts, however, should be taken as a hint to evaluate the particular company more thoroughly, for example by employing an expert team, or even by using a third model.

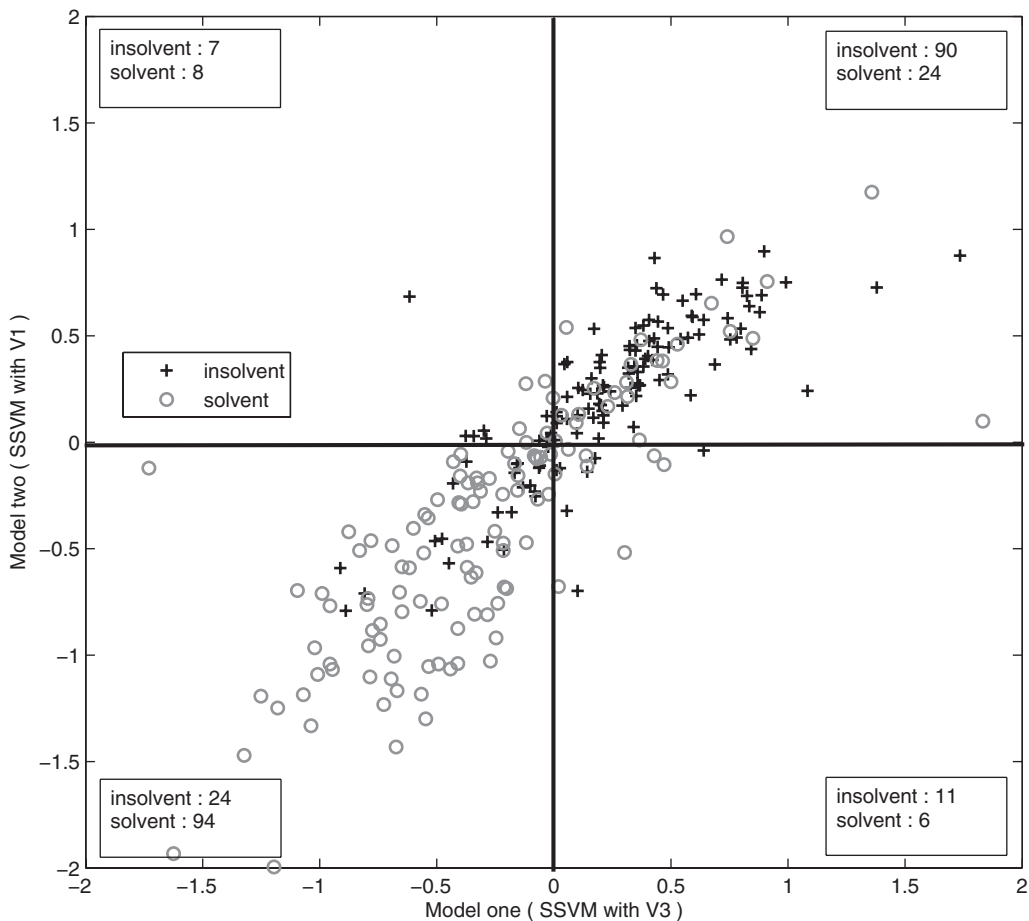


Figure 6. Data visualization via model one (generated by SSVM with V3) and model two (generated by SSVM with V1) in scenario S5

CONCLUSION

In this paper we apply different variants of support vector machines to a unique dataset of German solvent and insolvent companies. We use a priori a given set of predictors as a benchmark, and suggest two further variable selection procedures; the first procedure uses the 1-norm SVM and the second, incremental way consecutively selects the variable that is the farthest one from the column space of the current variable set. Given the three SSVM based on distinct variable sets, the relative performance of the types of smooth support vector machines is tested. The performance is measured by error rates. The two sets of variables newly created here lead to a dissimilar performance of SSVM. The selection of variables by the 1-norm SVM clearly underperforms compared to the incremental selection scheme. This difference in accuracy hints at the need for further research with respect to the variable selection. The training period makes a clear difference, though. Results improve considerably if more years of observation are used in training the machine. The SSVM

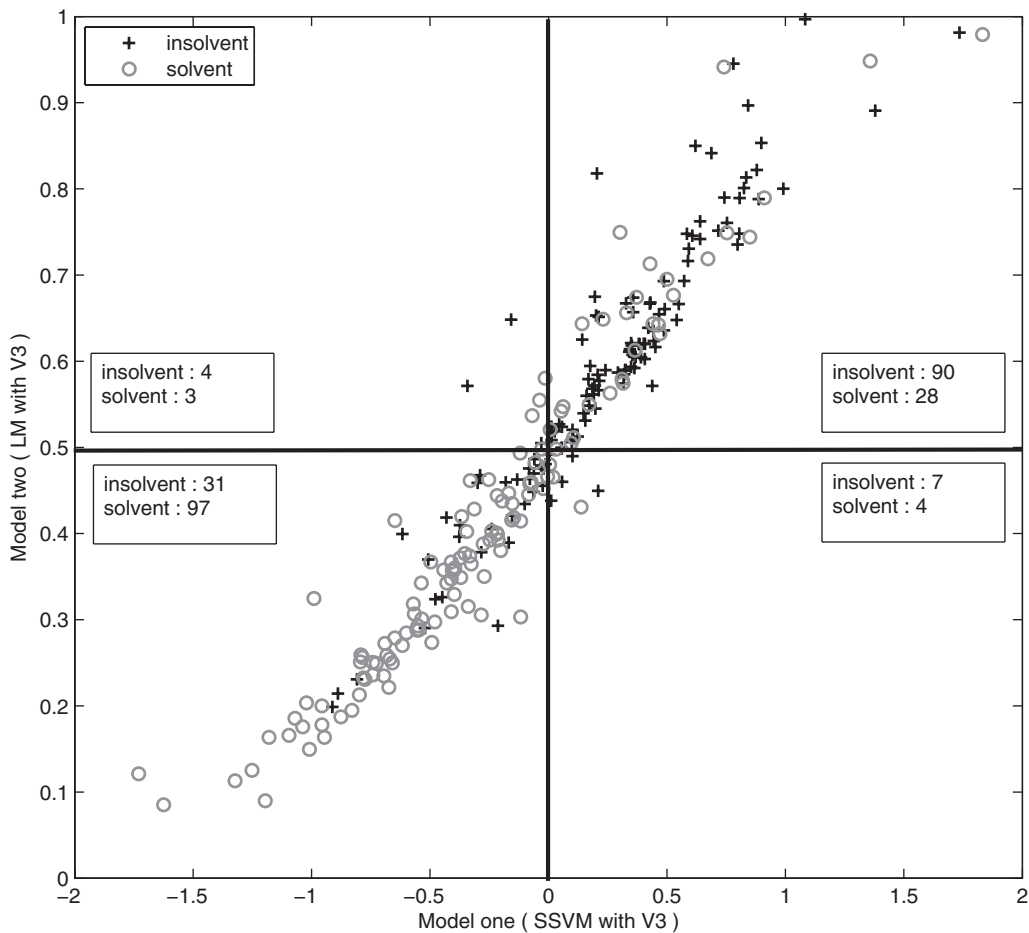


Figure 7. Data visualization via model one (generated by SSVM with V3) and model two (generated by LM with V3) in scenario S5

model benefits more from longer training periods than traditional methods do. As a consequence the logit model and discriminant analysis are both outperformed by the SSVM in long-term training scenarios. Moreover, the oversampling scheme works very well in dealing with unbalanced datasets. It provides flexibility to control the trade-off between Type I and Type II errors, and is therefore a strategic instrument in a bank's hand. The results generated are very stable in terms of small deviations of Type I, Type II and total error rates.

Finally, we want to stress that SSVM should be considered not as a substitute for traditional methods but rather as a complement which, when employed side by side with either the logit model or discriminant analysis, can generate new information that helps practitioners select those companies that are difficult to predict and, therefore, need more attention and further treatment.

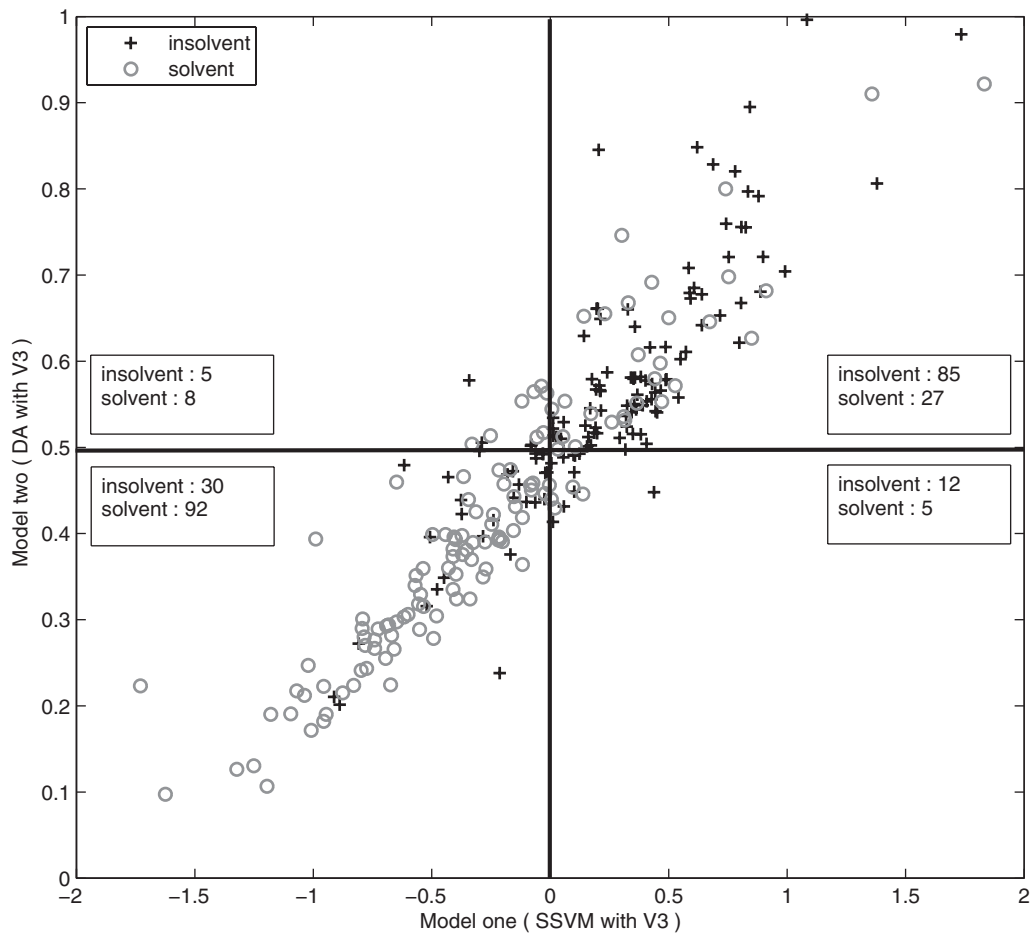


Figure 8, Data visualization via model one (generated by SSVM with V3) and model two (generated by DA with V3) in scenario S5

ACKNOWLEDGEMENTS

This research was supported by the ‘Stiftung Geld und Wahrung’ and by the Deutsche Forschungsgemeinschaft through the SFB 649 ‘Economic Risk’.

REFERENCES

- Altman E. 1968. Financial ratios, discriminant analysis and the prediction of corporate bankruptcy. *Journal of Finance* **23**(4): 589–609.
- Altman E, Marco G, Varetto F. 1994. Corporate distress diagnosis: comparisons using linear discriminant analysis and neural networks (the italian experience). *Journal of Banking and Finance* **18**: 505–529.

- Beaver W. 1966. Financial ratios as predictors of failures: empirical research in accounting: selected studies. *Journal of Accounting Research* **4**: 71–111.
- Burges CJC. 1998. A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery* **2**(2): 121–167.
- Chen S, Härdle W, Moro RA. 2006. Estimation of default probabilities with support vector machines. *SFB 649 Discussion Paper 2006-077*.
- Cristianini N, Shawe-Taylor J. 2000. *An Introduction to Support Vector Machines*. Cambridge University Press: Cambridge, UK.
- Drucker H, Burges CJC, Kaufman L, Smola A, Vapnik V. 1997. Support vector regression machines. In *Advances in Neural Information Processing Systems 9*, Mozer MC, Jordan MI, Petsche T (eds). MIT Press: Cambridge, MA; 155–161.
- Fung G, Mangasarian OL. 2004. A feature selection Newton method for support vector machine classification. *Computational Optimization and Applications* **28**(2): 185–202.
- Härdle W, Moro R, Schäfer D. 2007a. Graphical data representation in bankruptcy analysis based on support vector machines. In *Handbook of Data Visualization*, Chen C, Härdle W, Unwin A (eds). Springer: Heidelberg; 853–872.
- Härdle W, Moro RA, Schäfer D. 2007b. Estimating probabilities of default with support vector machines. *SFB 649 Discussion Paper 2007-035*.
- Huang CM, Lee YJ, Lin DKJ, Huang SY. 2007. Model selection for support vector machines via uniform design. *Computational Statistics and Data Analysis* **52**: 335–346. Special Issue on Machine Learning and Robust Data Mining (to appear).
- Hwang RC, Cheng KF, Lee JC. 2007. A semiparametric method for predicting bankruptcy. *Journal of Forecasting* **26**(5): 317–342.
- Krahen JP, Weber M. 2001. Generally accepted rating principles: a primer. *Journal of Banking and Finance* **25**(1): 3–23.
- Lee YJ, Huang SY. 2007. Reduced support vector machines: a statistical theory. *IEEE Transactions on Neural Networks* **18**: 1–13.
- Lee YJ, Mangasarian OL. 2001. SSVM: a smooth support vector machine. *Computational Optimization and Applications* **20**: 5–22.
- Lee YJ, Chang CC, Chao CH. 2008. Incremental forward feature selection with application to microarray gene expression. *Journal of Biopharmaceutical Statistics* **18**(5): 824–840.
- Leland H, Toft K. 1996. Optimal capital structure, endogenous bankruptcy, and the term structure of credit spreads. *Journal of Finance* **51**: 987–1019.
- Longstaff FA, Schwartz ES. 1995. A simple approach to valuating risky fixed and floating rate debt. *Journal of Finance* **50**: 789–819.
- Mangasarian OL, Musicant DR. 2000. Robust linear and support vector regression. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **22**(9): 950–955.
- Martin D. 1977. Early warning of bank failure: a logit regression approach. *Journal of Banking and Finance* **1**: 249–276.
- Mella-Barral P, Perraudin W. 1997. Strategic debt service. *Journal of Finance* **52**: 531–556.
- Merton R. 1974. On the pricing of corporate debt: the risk structure of interest rates. *Journal of Finance* **29**(2): 449–470.
- Ohlson J. 1980. Financial ratios and the probabilistic prediction of bankruptcy. *Journal of Accounting Research* **18**(1): 109–131.
- Schölkopf B, Smola AJ. 2002. *Learning with Kernels*. MIT Press: Cambridge, MA.
- Smola A, Schölkopf B. 2000. Sparse greedy matrix approximation for machine learning. In *Proceedings of the 17th International Conference on Machine Learning*, San Francisco, CA.
- Smola A, Schölkopf B. 2004. A tutorial on support vector regression. *Statistics and Computing* **14**: 199–222.
- Tam K, Kiang M. 1992. Managerial application of neural networks: the case of bank failure prediction. *Management Science* **38**(7): 926–947.
- Tibshirani R. 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society* **58**(1): 267–288.
- Vapnik VN. 1995. *The Nature of Statistical Learning Theory*. Springer: New York.
- Williams CKI, Seeger M. 2001. Using the Nyström method to speed up kernel machines. *Advances in Neural Information Processing Systems* **13**: 682–688.

- Zhou C. 2001. The term structure of credit spreads with jump risk. *Journal of Banking and Finance* **25**: 2015–2040.
- Zhu J, Rosset S, Hastie T, Tibshirani R. 2003. 1-Norm support vector machines. In *Advances in Neural Information Processing Systems* **16**: 49–56.

Authors' biographies:

Wolfgang Härdle did in 1982 his Dr. rer. nat. in Mathematics at Universität Heidelberg and in 1988 his Habilitation at Universität Bonn. He is currently chair professor of statistics at the Dept. of Economics and Business Administration, Humboldt-Universität zu Berlin. He is also director of CASE—Center for Applied Statistics & Economics and of the Collaborative Research Center 'Economic Risk'. His research focuses on dimension reduction techniques, computational statistics and quantitative finance. He has published 34 books and more than 200 papers in top statistical, econometrics and finance journals. He is one of the 'Highly cited Scientist' according to the Institute of Scientific Information.

Yuh-Jye Lee received his Master degree in Applied Mathematics from the National Tsing Hua University, Taiwan in 1992 and PhD degree in computer sciences from the University of Wisconsin-Madison in 2001. In 2002, Dr. Lee joined the Computer Science and Information Engineering Department, National Taiwan University of Science and Technology. He is an associate professor now. His research interests are in machine learning, data mining, optimization, information security and operations research. He developed new algorithms for large data mining problems such as classification problem, clustering, feature selection and dimension reduction. These algorithms have been used in intrusion detection systems (IDS), face detection, micro array gene expression analysis and breast cancer diagnosis and prognosis.

Dorothea Schäfer did in 1992 her Dr. rer. pol. in Economics and in the year 2000 her Habilitation at Freie Universität Berlin. She is currently coordinator of the research group Financial Markets and Financial Institutions and senior researcher at the German Institute for Economic Research (DIW) Berlin which she joined in 2002. She is managing editor of the Quarterly Journal of Economic Research (Vierteljahreshefte zur Wirtschaftsforschung) and adjunct lecturer at Freie Universität Berlin. Her research focuses on insolvency risk, financial management of firms and banks, and on behavioural finance.

Yi-Ren Yeh received the M.S. degree from the Department of Computer Science and Information Engineering, National Taiwan University of Science and Technology, Taipei, Taiwan, R.O.C., in 2006. He is currently working toward the PhD degree in the Department of Computer Science and Information Engineering, National Taiwan University of Science and Technology. His research interests include machine learning, data mining, and information security.

Authors' addresses:

Wolfgang Härdle, Center for Applied Statistics and Economics, Humboldt-Universität zu Berlin, Spandauer Straße 1, 10178 Berlin, Germany.

Yuh-Jye Lee and **Yi-Ren Yeh**, Department of Computer Science Information Engineering, National Taiwan University of Science and Technology, Taipei 106, Taiwan.

Dorothea Schäfer, German Institute for Economic Research (DIW) Berlin, Mohrenstrasse 58, 10117 Berlin, Germany.

Time Series Modelling With Semiparametric Factor Dynamics

Byeong U. PARK, Enno MAMMEN, Wolfgang HÄRDLE, and Szymon BORAK

High-dimensional regression problems, which reveal dynamic behavior, are typically analyzed by time propagation of a few number of factors. The inference on the whole system is then based on the low-dimensional time series analysis. Such high-dimensional problems occur frequently in many different fields of science. In this article we address the problem of inference when the factors and factor loadings are estimated by semiparametric methods. This more flexible modeling approach poses an important question: Is it justified, from an inferential point of view, to base statistical inference on the estimated time series factors? We show that the difference of the inference based on the estimated time series and “true” unobserved time series is asymptotically negligible. Our results justify fitting vector autoregressive processes to the estimated factors, which allows one to study the dynamics of the whole high-dimensional system with a low-dimensional representation. We illustrate the theory with a simulation study. Also, we apply the method to a study of the dynamic behavior of implied volatilities and to the analysis of functional magnetic resonance imaging (fMRI) data.

KEY WORDS: Asymptotic inference; Factor models; Implied volatility surface; Semiparametric models; Vector autoregressive process.

1 INTRODUCTION

Modeling for high-dimensional data is a challenging task in statistics especially when the data comes in a dynamic context and is observed at changing locations with different sample sizes. Such modeling challenges appear in many different fields. Examples are Stock and Watson (2005) in empirical macroeconomics, Lee and Carter (1992) in mortality analysis, Nelson and Siegel (1987) and Diebold and Li (2006) in bond portfolio risk management or derivative pricing, Martinussen and Scheike (2000) in biomedical research. Other examples include the studies of radiation treatment of prostate cancer by Kauermann (2000) and evoked potentials in Electroencephalogram (EEG) analysis by Gasser, Möcks, and Verleger (1983). In financial engineering, it is common to analyze the dynamics of implied volatility surface for risk management. For functional magnetic resonance imaging data (fMRI), one may be interested in analyzing the brain’s response over time as well as identifying its activation area, see Worsley et al. (2002).

A successful modeling approach utilizes factor type models, which allow low-dimensional representation of the data. In an orthogonal L -factor model an observable J -dimensional random vector $Y_t = (Y_{t,1}, \dots, Y_{t,J})^T$ can be represented as

$$Y_{t,j} = m_{0,j} + Z_{t,1}m_{1,j} + \dots + Z_{t,L}m_{L,j} + \varepsilon_{t,j}, \quad (1)$$

where $Z_{t,l}$ are common factors, $\varepsilon_{t,j}$ are errors or specific factors, and the coefficients $m_{l,j}$ are factor loadings. In most applications, the index $t = 1, \dots, T$ reflects the time evolution of the whole system, and Y_t can be considered as a multidimensional time series. For a method to identify common factors in this model we refer to Peña and Box (1987). The study of high-dimensional Y_t is then simplified to the modeling of $Z_t = (Z_{t,1},$

$\dots, Z_{t,L})^T$, which is a more feasible task when $L \ll J$. The model (1) reduces to a special case of the generalized dynamic factor model considered by Forni, Hallin, Lippi, and Reichlin (2000), Forni and Lippi (2001) and Hallin and Liska (2007), when $Z_{t,l} = a_{l,1}(B)U_{t,1} + \dots + a_{l,q}(B)U_{t,q}$ where the q -dimensional vector process $U_t = (U_{t,1}, \dots, U_{t,q})^T$ is an orthonormal white noise and B stands for the lag operator. In this case, the model (1) is expressed as $Y_{t,j} = m_{0,j} + \sum_{k=1}^q b_{k,j}(B)U_{t,k} + \varepsilon_{t,j}$, where $b_{k,j}(B) = \sum_{l=1}^L a_{l,k}(B)m_{l,j}$.

In a variety of applications, one has explanatory variables $X_{t,j} \in \mathbb{R}^d$ at hand that may influence the factor loadings m_l . An important refinement of the model (1) is to incorporate the existence of observable covariates $X_{t,j}$. The factor loadings are now generalized to functions of $X_{t,j}$, so that the model (1) is generalized to

$$Y_{t,j} = m_0(X_{t,j}) + \sum_{l=1}^L Z_{t,l}m_l(X_{t,j}) + \varepsilon_{t,j}, \quad 1 \leq j \leq J_t. \quad (2)$$

In this model, $Z_{t,l}$ for each $l: 1 \leq l \leq L$ enters into all $Y_{t,j}$ for j such that $m_l(X_{t,j}) \neq 0$. Note that the probability of the event that $m_l(X_{t,j}) = 0$ for some $1 \leq j \leq J$ equals zero if $m_l(x) = 0$ at countably many points of x and the density f_l of $X_{t,j}$ is supported on an interval with nonempty interior, as we assume at (A2) in Section 5.

The model (2) can be interpreted as a discrete version of the following functional extension of the model (1):

$$Y_t(x) = m_0(x) + \sum_{l=1}^L Z_{t,l}m_l(x) + \varepsilon_t(x), \quad (3)$$

where $\varepsilon_t(\cdot)$ is a mean zero stochastic process, and also regarded as a regression model with embedded time evolution. It is different from varying-coefficient models, such as in Fan, Yao, and Cai (2003) and Yang, Park, Xue, and Härdle (2006), because Z_t is unobservable. Our model also has some similarities to the one considered in Connor and Linton (2007) and Connor, Hagmann, and Linton (2007), which generalized the study of Fama and French (1992) on the common movements of stock price returns. There, the covariates, denoted by $X_{l,j}$, are

Byeong U. Park is Professor, Department of Statistics, Seoul National University Seoul 151-747, Korea (E-mail: bupark@stats.snu.ac.kr). Enno Mammen is Professor, Department of Economics, University of Mannheim, 68131 Mannheim, Germany (E-mail: emammen@rumms.uni-mannheim.de). Wolfgang Härdle is Professor, Institute for Statistics and Econometrics, Humboldt Universität zu Berlin, D-10178 Berlin, Germany (E-mail: haerdle@wiwi.hu-berlin.de). Szymon Borak is Ph.D. Student, Institute for Statistics and Econometrics, Humboldt Universität zu Berlin, D-10178 Berlin, Germany (E-mail: szymon.borak@gmail.de). The authors gratefully acknowledge financial support Deutsche Forschungsgemeinschaft and the Sonderforschungsbereich 649 “Ökonomisches Risiko.” Byeong U. Park’s research was supported by the Korea Research Foundation Grant funded by the Korean Government (MOEHRD) (KRF-2005-070-C00021). The authors thank the associate editor and referees for their helpful comments and suggestions.

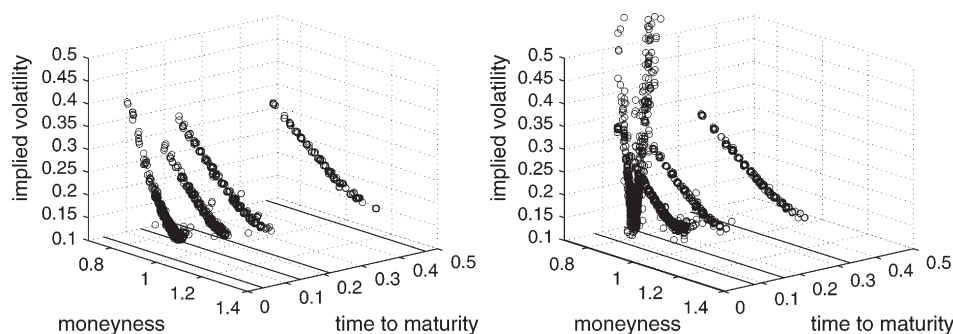


Figure 1. The typical IV data design on two different days. In the maturity direction observations appear in the discrete points for each particular day. Bottom solid lines indicate the observed maturities. Left panel: observations on 2004.07.08, $J_t = 5,606$. Right panel: observations on 2004.08.19, $J_t = 8,152$.

time-invariant and are different for different m_t , which allows a direct application of backfitting procedures and makes the problem quite different from our setting. Some linear models, which allow time-varying coefficients, as considered in Hansen, Nielsen, and Nielsen (2004) and Bruback and Rice (1998), may be recognized as a special case of (2).

In this article we consider the model (2) with unknown nonparametric functions m_t . We call this model a dynamic semiparametric factor model (DSFM). The evolution of complex high-dimensional objects may be described by (2), so that their analysis can be reduced to the study of a low-dimensional vector of factors Z_t . In the present article, we consider an efficient nonparametric method of fitting the model. We provide relevant theory for the method as well as illustrate its empirical aspects through a simulation and a real data application. Fengler, Härdle, and Mammen (2007) used a kernel smoothing approach for the same model, but it was focused on a particular data application without offering any discussion of numerical issues, statistical theory, and simulation analysis.

One of the main motivations for the model (2) comes from a special structure of the implied volatility (IV) data, as is observed in Figure 1. The IV is a volatility parameter that matches the observed plain vanilla option prices with the theoretical ones given by the formula of Black and Scholes (1973). Figure 1 shows the special “string” structure of the IV data obtained from the European option prices on the German stock index DAX (ODAX) for two different days. The volatility strings shift toward expiry, which is indicated by the bottom line in the figure. Moreover the shape of the IV strings is subject to stochastic deformation. Fengler et al. (2007) proposed to use the model (2) to describe the dynamics of the IV data, where $Y_{t,j}$ are the values of IV or those of its transformation on the day t , and $X_{t,j}$ are the two-dimensional vectors of the moneyness and time-to-maturity. For more details on the data design and econometric motivation, we refer to Fengler et al. (2007).

One may find another application of the model (2) in the analysis of functional magnetic resonance imaging (fMRI) data. The fMRI is a noninvasive technique of recording brain’s signals on spatial area in every particular time period (usually 1–4 sec). One obtains a series of three-dimensional images of the blood-oxygen-level-dependent (BOLD) fMRI signals, whereas an exercised person is subject to certain stimuli. An example of the images in 15 different slices at one particular time point is presented in Figure 2. For the more detailed

description on the fMRI methodology we refer to Logothetis and Wandell (2004). The main aims of the statistical methods in this field are identification of the brain’s activation areas and analysis of its response over time. For this purpose the model (2) can be applied. DSFM may be applied to many other problems, such as modeling of yield curve evolution where the standard approach is to use the parametric factor model proposed by Nelson and Siegel (1987).

Our methods produce estimates of the true unobservable Z_t , say \hat{Z}_t , as well as estimates of the unknown functions m_t . In practice, one operates on these estimated values of Z_t for further statistical analysis of the data. In particular, for the IV application, one needs to fit an econometric model to the estimated factors \hat{Z}_t . For example, Hafner (2004) and Cont and da Fonseca (2002) fitted an AR(1) process to each factor, and Fengler et al. (2007) considered a multivariate VAR(2) model. The main question that arises from these applications is whether the inference based on \hat{Z}_t is equivalent to the one based on Z_t . Attempting to give an answer to this question forms the core of this article.

It is worthwhile to note here that Z_t is not identifiable in the model (2). There are many versions of (Z_t, m) , where $m = (m_0, \dots, m_L)^T$, that give the same distribution of Y_t . This means that estimates of Z_t and m_t are not uniquely defined. We show that for any version of $\{Z_t\}$ there exists a version of $\{\hat{Z}_t\}$ whose lagged covariances are asymptotically the same as those of $\{Z_t\}$. This justifies the inference based on $\{\hat{Z}_t\}$ when $\{Z_t\}$ is a VAR process, in particular. We confirm this theoretical result by a Monte Carlo simulation study. We also discuss fitting the model to the real ODAX IV and fMRI data.

The article is organized as follows. In the next section we propose a new method of fitting DSFM and an iterative algorithm that converges at a geometric rate. In Section 3 we present the results of a simulation study that illustrate the theoretical findings given in Section 5. In Section 4 we apply the model to the ODAX IV and fMRI data. Section 5 is devoted to the asymptotic analysis of the method. Technical details are provided in the Appendix.

2. METHODOLOGY

We observe $(X_{t,j}, Y_{t,j})$ for $j = 1, \dots, J_t$ and $t = 1, \dots, T$ such that

$$Y_{t,j} = Z_t^T m(X_{t,j}) + \varepsilon_{t,j}. \quad (4)$$

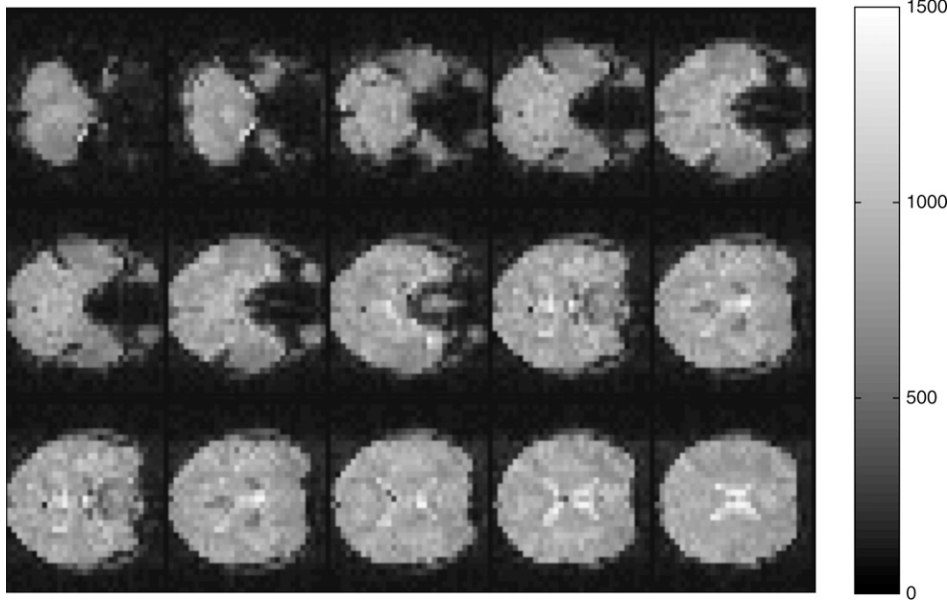


Figure 2. Typical fMRI data in one particular time point. The figure presents 15 parallel horizontal images. The brightness corresponds to the strength of the observed signals.

Here $Z_t^\top = (1, Z_t^\top)$ and $Z_t = (Z_{t,1}, \dots, Z_{t,L})^\top$ is an unobservable L -dimensional process. The function m is an $(L + 1)$ -tuple (m_0, \dots, m_L) of unknown real-valued functions m_l defined on a subset of \mathbb{R}^d . The variables $X_{1,1}, \dots, X_{T,J_T}, \varepsilon_{1,1}, \dots, \varepsilon_{T,J_T}$ are independent. The errors $\varepsilon_{t,j}$ have zero means and finite second moments. For simplicity of notation, we will assume that the covariates $X_{t,j}$ have support $[0, 1]^d$, and also that $J_t \equiv J$ do not depend on t .

For the estimation of m , we use a series estimator. For an integer $K \geq 1$, we choose functions $\psi_1, \dots, \psi_K: [0, 1]^d \rightarrow \mathbb{R}$, which are normalized so that $\int_{[0,1]^d} \psi_k^2(x) dx = 1$. For example, one may take $\{\psi_k: 1 \leq k \leq K\}$ to be a tensor B-spline basis (e.g., see de Boor 2001). Then, an $(L + 1)$ -tuple of functions $m = (m_0, \dots, m_L)^\top$ may be approximated by $\mathcal{A}\psi$, where $\mathcal{A} = (\alpha_{l,k})$ is an $(L + 1) \times K$ matrix and $\psi = (\psi_1, \dots, \psi_K)^\top$. We define the least squares estimators $\widehat{Z}_t = (\widehat{Z}_{t,1}, \dots, \widehat{Z}_{t,L})^\top$ and $\widehat{\mathcal{A}} = (\widehat{\alpha}_{l,k})$:

$$S(\mathcal{A}, z) \equiv \sum_{t=1}^T \sum_{j=1}^J \{Y_{t,j} - (1, z_t^\top) \mathcal{A} \psi(X_{t,j})\}^2 = \min_{\mathcal{A}, z}! \quad (5)$$

where $z = (z_1^\top, \dots, z_T^\top)^\top$ for L -dimensional vectors z_t . With $\widehat{\mathcal{A}}$ at hand, we estimate m by $\widehat{m} = \widehat{\mathcal{A}}\psi$.

We note that, given z or \mathcal{A} , the function S in (5) is quadratic with respect to the other variables, and thus has an explicit unique minimizer. However, minimization of S with respect to \mathcal{A} and z simultaneously is a fourth-order problem. The solution is neither unique nor explicit. It is unique only up to the values of $\widehat{z}_1^\top \widehat{\mathcal{A}}, \dots, \widehat{z}_T^\top \widehat{\mathcal{A}}$, where $\widehat{z}_t^\top = (1, \widehat{Z}_t^\top)$. We will come back to this identifiability issue later in this section.

To find a solution $(\widehat{\mathcal{A}}, \widehat{z})$ of the minimization problem (5), one might adopt the following iterative algorithm: (i) Given an initial choice $Z^{(0)}$, minimize $S(\mathcal{A}, Z^{(0)})$ with respect to \mathcal{A} , which is an ordinary least squares problem and thus has an explicit unique solution. Call it $\mathcal{A}^{(1)}$. (ii) Minimize $S(\mathcal{A}^{(1)}, z)$ with respect to z now, which is also an ordinary least squares

problem. (iii) Iterate (i) and (ii) until convergence. This is the approach taken by Fengler et al. (2007). However, the procedure is not guaranteed to converge to a solution of the original problem.

We propose to use a Newton-Raphson algorithm. Let $\alpha \equiv \alpha(\mathcal{A})$ denote the stack form of $\mathcal{A} = (\alpha_{l,k})$ [i.e., $\alpha = (\alpha_{0,1}, \dots, \alpha_{L,1}, \alpha_{0,2}, \dots, \alpha_{L,2}, \dots, \alpha_{0,K}, \dots, \alpha_{L,K})^\top$]. In a slight abuse of notation we write $S(\alpha, z)$ for $S(\mathcal{A}, z)$. Define

$$\begin{aligned} F_{10}(\alpha, z) &= \frac{\partial}{\partial \alpha} S(\alpha, z), & F_{01}(\alpha, z) &= \frac{\partial}{\partial z} S(\alpha, z), \\ F_{20}(\alpha, z) &= \frac{\partial^2}{\partial \alpha^2} S(\alpha, z), & F_{11}(\alpha, z) &= \frac{\partial^2}{\partial \alpha \partial z} S(\alpha, z), \\ F_{02}(\alpha, z) &= \frac{\partial^2}{\partial z^2} S(\alpha, z). \end{aligned}$$

Let $\Psi_t = [\psi(X_{t,1}), \dots, \psi(X_{t,J})]$ be a $K \times J$ matrix. Define A to be the $L \times K$ matrix obtained by deleting the first row of \mathcal{A} . Writing $\zeta_t^\top = (1, z_t^\top)$, it can be shown that

$$\begin{aligned} F_{10}(\alpha, z) &= 2 \sum_{t=1}^T [(\Psi_t \Psi_t^\top) \otimes (\zeta_t \zeta_t^\top)] \alpha - 2 \sum_{t=1}^T (\Psi_t Y_t) \otimes \zeta_t, \\ F_{20}(\alpha, z) &= 2 \sum_{t=1}^T [(\Psi_t \Psi_t^\top) \otimes (\zeta_t \zeta_t^\top)], \end{aligned}$$

$F_{01}(\alpha, z)^\top = 2(\zeta_1^\top \mathcal{A} \Psi_1 \Psi_1^\top A^\top - Y_1^\top \Psi_1^\top A^\top, \dots, \zeta_T^\top \mathcal{A} \Psi_T \Psi_T^\top A^\top - Y_T^\top \Psi_T^\top A^\top)$, and $F_{02}(\alpha, z)$ equals a $(TL) \times (TL)$ matrix that consists of T diagonal blocks $A \Psi_t \Psi_t^\top A^\top$ for $t = 1, \dots, T$. Here and later, \otimes denotes the Kronecker product operator. Also, by some algebraic manipulations it can be shown that

$$[(\Psi_t \Psi_t^\top) \otimes (\zeta_t \zeta_t^\top)] \alpha = (\Psi_t \Psi_t^\top A^\top \zeta_t) \otimes \zeta_t. \quad (6)$$

Let \mathcal{I} be an $(L + 1) \times L$ matrix such that $\mathcal{I}^\top = (0, I_L)$ and I_L denote the identity matrix of dimension L . Define $F_{11,t}(\alpha, z) = (\Psi_t \Psi_t^\top A^\top) \otimes \zeta_t + (\Psi_t \Psi_t^\top A^\top \zeta_t) \otimes \mathcal{I} - (\Psi_t Y_t) \otimes$

\mathcal{I} . Then, we get $F_{11}(\alpha, z) = 2(F_{11,1}(\alpha, z), F_{11,2}(\alpha, z), \dots, F_{11,T}(\alpha, z))$. Let

$$F(\alpha, z) = \begin{pmatrix} F_{10}(\alpha, z) \\ F_{01}(\alpha, z) \end{pmatrix}, F'(\alpha, z) = \begin{pmatrix} F_{20}(\alpha, z) & F_{11}(\alpha, z) \\ F_{11}(\alpha, z)^\top & F_{02}(\alpha, z) \end{pmatrix}.$$

We need to solve the equation $F(\alpha, z) = 0$ simultaneously for α and z . We note that the matrices $(\Psi_t \Psi_t^\top) \otimes (\zeta_t \zeta_t^\top) = (\Psi_t \otimes \zeta_t)(\Psi_t \otimes \zeta_t)^\top$ and $A \Psi_t \Psi_t^\top A^\top$ are nonnegative definite. Thus, by Miranda's existence theorem (for example, see Vrahatis 1989) the nonlinear system of equations $F(\alpha, z) = 0$ has a solution.

Given $(\alpha^{\text{OLD}}, Z^{\text{OLD}})$, the Newton-Raphson algorithm gives the updating equation for $(\alpha^{\text{NEW}}, Z^{\text{NEW}})$:

$$\begin{pmatrix} \alpha^{\text{NEW}} \\ Z^{\text{NEW}} \end{pmatrix} = \begin{pmatrix} \alpha^{\text{OLD}} \\ Z^{\text{OLD}} \end{pmatrix} - F'_*(\alpha^{\text{OLD}}, Z^{\text{OLD}})^{-1} F(\alpha^{\text{OLD}}, Z^{\text{OLD}}), \quad (7)$$

where $F'_*(\alpha, z)$ for each given (α, z) is the restriction to \mathcal{F}_* of the linear map defined by the matrix $F'(\alpha, z)$ and \mathcal{F}_* is the linear space of values of (α, z) with $\sum_{t=1}^T z_t = 0$ and $\sum_{t=1}^T Z_t^{(0)}(z_t - Z_t^{(0)})^\top = 0$. We denote the initial value of the algorithm by $(\alpha^{(0)}, Z^{(0)})$. We will argue later that under mild conditions, $(\hat{\alpha}, \hat{Z})$ can be chosen as an element of \mathcal{F}_* .

The algorithm (7) is shown to converge to a solution of (5) at a geometric rate under some weak conditions on the initial choice $(\alpha^{(0)}, Z^{(0)})$, as is demonstrated by Theorem 1 later. We collect the conditions for the theorem.

(C1) It holds that $\sum_{t=1}^T Z_t^{(0)} = 0$. The matrix $\sum_{t=1}^T Z_t^{(0)} Z_t^{(0)\top}$ and the map $F'_*(\alpha^{(0)}, Z^{(0)})$ are invertible.

(C2) There exists a version $(\hat{\alpha}, \hat{Z})$ with $\sum_{t=1}^T \hat{Z}_t = 0$ such that $\sum_{t=1}^T \hat{Z}_t Z_t^{(0)\top}$ is invertible. Also, $\hat{\alpha}_l = (\hat{\alpha}_{l1}, \dots, \hat{\alpha}_{lK})^\top$ for $l = 0, \dots, L$ are linearly independent.

Let $\alpha^{(k)}$ and $Z^{(k)}$ denote the k th updated vectors in the iteration with the algorithm (7). Also, we write $\mathcal{A}^{(k)}$ for the matrix that corresponds to $\alpha^{(k)}$, and $\mathcal{Z}_t^{(k)\top} = (1, Z_t^{(k)\top})$.

Theorem 1. Let T, J and K be held fixed. Suppose that the initial choice $(\alpha^{(0)}, Z^{(0)})$ satisfies (C1) and (C2). Then, for any constant $0 < \gamma < 1$ there exist $r > 0$ and $C > 0$, which are random variables depending on $\{(X_{t,j}, Y_{t,j})\}$, such that, if $\sum_{t=1}^T \|\mathcal{Z}_t^{(0)\top} \mathcal{A}^{(0)} - \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}}\|^2 \leq r$, then

$$\sum_{t=1}^T \|\mathcal{Z}_t^{(k)\top} \mathcal{A}^{(k)} - \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}}\|^2 \leq C 2^{-2(k-1)} \gamma^{2(2^k-1)}.$$

We now argue that under (C1) and (C2), $(\hat{\alpha}, \hat{Z})$ can be chosen as an element of \mathcal{F}_* . Note first that one can always take $Z_t^{(0)}$ and \hat{Z}_t so that $\sum_{t=1}^T Z_t^{(0)} = 0$ and $\sum_{t=1}^T \hat{Z}_t = 0$. This is because, for any version $(\hat{\alpha}, \hat{Z})$, one has

$$\begin{aligned} \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}} &= \hat{\alpha}_0^\top + \sum_{l=1}^L \hat{Z}_{t,l} \hat{\alpha}_l^\top = \left(\hat{\alpha}_0^\top + \sum_{l=1}^L \bar{\bar{Z}}_{t,l} \bar{\bar{\alpha}}_l^\top \right) \\ &+ \sum_{l=1}^L (\hat{Z}_{t,l} - \bar{\bar{Z}}_{t,l}) \hat{\alpha}_l^\top \stackrel{\text{let}}{=} \hat{\alpha}_0^{*\top} + \sum_{l=1}^L \hat{Z}_{t,l}^* \hat{\alpha}_l^\top = \hat{\mathcal{Z}}_t^{*\top} \hat{\mathcal{A}}^*, \end{aligned}$$

where $\bar{\bar{Z}}_t = T^{-1} \sum_{t=1}^T \bar{\bar{Z}}_{t,l}$, $\hat{\mathcal{Z}}_t^{*\top} = (1, \hat{Z}_{t,l}^{*\top})$ and $\hat{\mathcal{A}}^*$ is the matrix obtained from $\hat{\mathcal{A}}$ by replacing its first row by $\hat{\alpha}_0^{*\top}$. Furthermore, the minimization problem (5) has no unique solution. If $(\hat{Z}_t, \hat{\mathcal{A}})$ or $(\hat{Z}_t, \hat{m} = \hat{\mathcal{A}}\psi)$ is a minimizer, then also $(B^\top \hat{Z}_t, \tilde{B}^{-1} \hat{m})$ is a minimizer. Here

$$\tilde{B} = \begin{pmatrix} 1 & 0 \\ 0 & B \end{pmatrix} \quad (8)$$

and B is an arbitrary invertible matrix. The special structure of \tilde{B} assures that the first component of $\tilde{B}^\top \hat{Z}_t$ equals 1. In particular, with the choice $B = (\sum_{t=1}^T Z_t^{(0)} \hat{Z}_t^\top)^{-1} \sum_{t=1}^T Z_t^{(0)} Z_t^{(0)\top}$ we get for $\hat{Z}_t^* = B^\top \hat{Z}_t$ that $\sum_{t=1}^T Z_t^{(0)} (\hat{Z}_t^* - Z_t^{(0)})^\top = 0$.

In Section 5, we will show that, for any solution \hat{Z}_t and for any version of true Z_t , there exists a random matrix B such that $\tilde{Z}_t = B^\top \hat{Z}_t$ has asymptotically the same covariance structure as Z_t . This means that the difference of the inferences based on \tilde{Z}_t and Z_t is asymptotically negligible.

We also note that one can always choose $\hat{m} = \hat{\mathcal{A}}\psi$ such that the components $\hat{m}_1, \dots, \hat{m}_L$ are orthonormal in $L_2([0, 1]^d)$ or in other L_2 [e.g., in $L_2(T^{-1} \sum_{t=1}^T \hat{f}_t)$ where \hat{f}_t is a kernel estimate of the density of $X_{t,j}$]. If one selects \hat{m} in this way, then the matrix B should be an orthogonal matrix and the underlying time series Z_t is estimated up to such transformations.

In practice one needs to choose an initial estimate $(\alpha^{(0)}, Z^{(0)})$ to run the algorithm. One may generate normal random variates for $Z_{t,l}^{(0)}$, and then find the initial $\alpha^{(0)}$ by solving the equation $F_{10}(\alpha, Z^{(0)})$. This initial choice was found to work well in our numerical study presented in Sections 3 and 4.

As an alternative way of fitting the model (2), one may extend the idea of the principal component method that is used to fit the orthogonal factor model (1). In this way, the data $\{Y_{t,j}; 1 \leq j \leq J\}$ are viewed as the values of a functional datum $Y_t(\cdot)$ observed at $x = X_{t,j}$, $1 \leq j \leq J$, and the functional factor model given at (3) may be fitted with smooth approximations of Y_t obtained from the original dataset. If one assumes $E Z_t = 0$, $\text{var}(Z_t) = I_L$, as is typically the case with the orthogonal factor model (1), then one can estimate m_t and Z_t by performing functional principal component analysis with the sample covariance function

$$\hat{K}(x, x') = T^{-1} \sum_{t=1}^T \{Y_t(x) - \bar{Y}(x)\} \{Y_t(x') - \bar{Y}(x')\},$$

where $\bar{Y}(x) = T^{-1} \sum_{t=1}^T Y_t(x)$. There are some limitations for this approach. First, it requires initial fits to get smooth approximations of $Y_t(\cdot)$, which may be difficult when the design points $X_{t,j}$ are sparse as is the case with the IV application. Our method avoids the preliminary estimation and shifts the discrete representation directly to the functions m_t . Second, for the method to work one needs at least stationarity of Z_t and ε_t , whereas our theory does not rely on these assumptions.

3. SIMULATION STUDY

In Theorem 3 we will argue that the inference based on the covariances of the unobserved factors Z_t is asymptotically equivalent to the one based on $B^\top \hat{Z}_t$ for some invertible B . In this section we illustrate the equivalence by a simulation study. We compare the covariances of Z_t and $\tilde{Z}_t \equiv B^\top \hat{Z}_t$, where

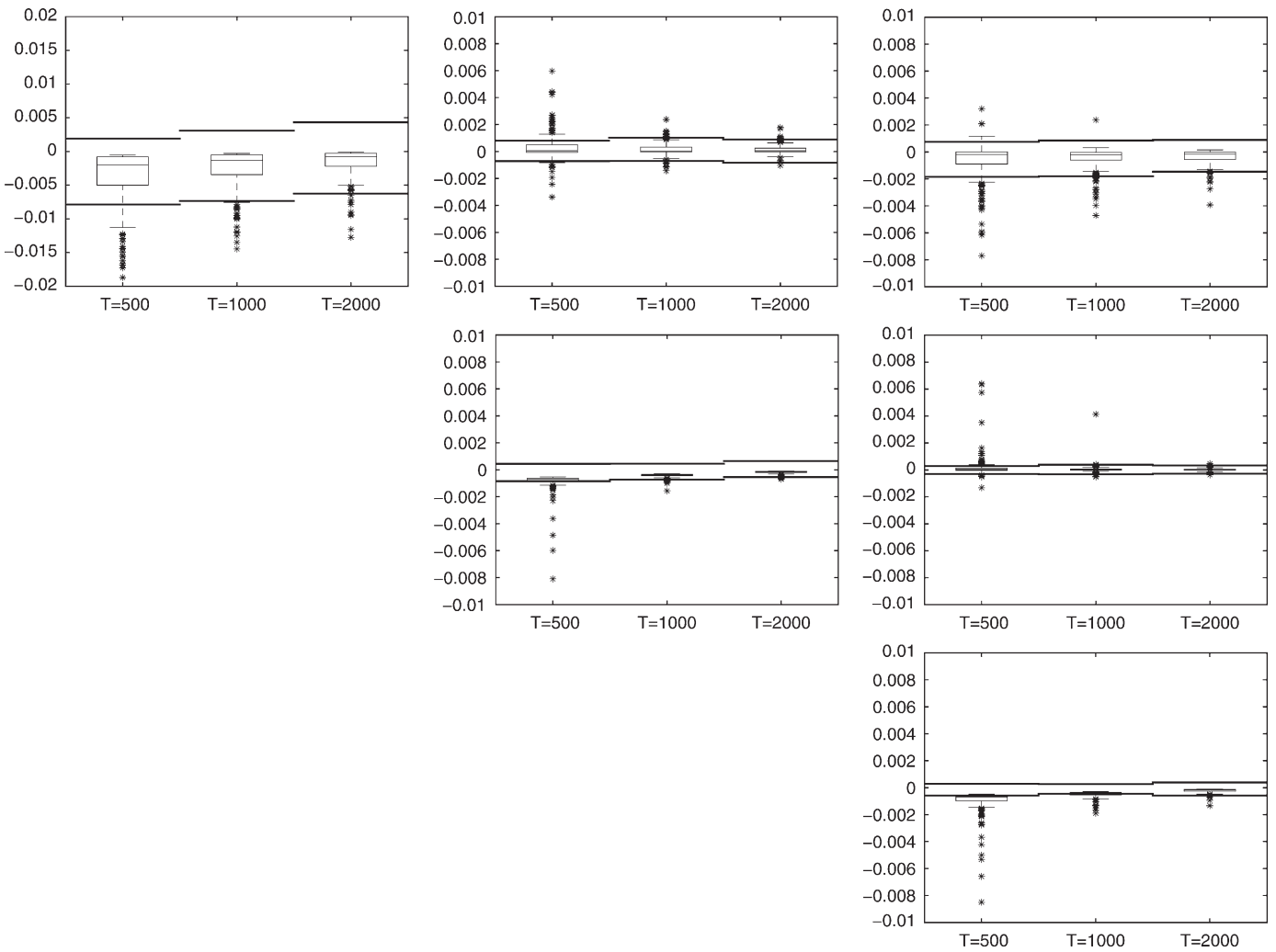


Figure 3. The boxplots based on 250 values of the entries of the scaled difference of the covariance matrices given at (10). The lengths of the series Z_t and \tilde{Z}_t were 500, 1,000, 2,000. The thick lines represent the upper and lower quartiles of (11).

$$B = \left(T^{-1} \sum_{t=1}^T Z_{c,t} \hat{Z}_{c,t}^\top \right)^{-1} T^{-1} \sum_{t=1}^T Z_{c,t} Z_{c,t}^\top, \quad (9)$$

$Z_{c,t} = Z_t - T^{-1} \sum_{s=1}^T Z_s$ and $\hat{Z}_{c,t} = \hat{Z}_t - T^{-1} \sum_{s=1}^T \hat{Z}_s$. Note that B at (9) minimizes $\sum_{t=1}^T \| \hat{Z}_{c,t} - (B^\top)^{-1} Z_{c,t} \|^2$. In the Appendix we will prove that Theorem 3 holds with the choice at (9).

We took $T = 500, 1,000, 2,000, J = 100, 250, 1,000$ and $K = 36, 49, 64$. We considered $d = 2, L = 3$ and the following tuple of 2-dimensional functions:

$$\begin{aligned} m_0(x_1, x_2) &= 1, & m_1(x_1, x_2) &= 3.46(x_1 - .5), \\ m_2(x_1, x_2) &= 9.45 \left\{ (x_1 - .5)^2 + (x_2 - .5)^2 \right\} - 1.6, \\ m_3(x_1, x_2) &= 1.41 \sin(2\pi x_2). \end{aligned}$$

The coefficients in these functions were chosen so that m_1, m_2, m_3 are close to orthogonal. We generated Z_t from a centered VAR(1) process $Z_t = \mathcal{R}Z_{t-1} + U_t$, where U_t is $N_3(0, \Sigma_U)$ random vector, the rows of \mathcal{R} from the top equal $(0.95, -0.2, 0), (0, 0.8, 0.1), (0.1, 0, 0.6)$, and $\Sigma_U = 10^{-4}I_3$. The design points $X_{t,j}$ were independently generated from a uniform distribution on the unit square, $\varepsilon_{t,j}$ were iid $N(0, \sigma^2)$ with $\sigma = 0.05$, and $Y_{t,j}$

were obtained according to the model (4). The simulation experiment was repeated 250 times for each combination of (T, J, K) . For the estimation we employed, for ψ_j , the tensor products of linear B-splines. The one-dimensional linear B-splines $\tilde{\psi}_k$ are defined on a consecutive equidistant knots x^k, x^{k+1}, x^{k+2} by $\tilde{\psi}_k(x) = (x - x^k)/(x^{k+1} - x^k)$ for $x \in (x^k, x^{k+1}]$, $\tilde{\psi}_k(x) = (x^{k+2} - x)/(x^{k+2} - x^{k+1})$ for $x \in (x^{k+1}, x^{k+2}]$, and $\tilde{\psi}_k(x) = 0$ otherwise. We chose $K = 8 \times 8 = 64$.

We plotted in Figure 3 the entries of the scaled difference of the covariance matrices

$$\tilde{D} = \frac{1}{\sqrt{T}} \left\{ \sum_{t=1}^T (\tilde{Z}_t - \bar{\tilde{Z}}) (\tilde{Z}_t - \bar{\tilde{Z}})^\top - \sum_{t=1}^T (Z_t - \bar{Z}) (Z_t - \bar{Z})^\top \right\}. \quad (10)$$

Each panel of Figure 3 corresponds to one entry of the matrix \tilde{D} , and the three boxplots in each panel represent the distributions of the 250 values of the corresponding entry for $T = 500, 1,000, 2,000$. In the figure we also depicted, by thick lines, the upper and lower quartiles of

$$D = \frac{1}{\sqrt{T}} \left\{ \sum_{t=1}^T (Z_t - \bar{Z}) (Z_t - \bar{Z})^\top - T\Gamma \right\}, \quad (11)$$

where Γ is the true covariance matrix of the simulated VAR process. We refer to Lütkepohl (1993) for a representation of Γ .

Our theory in Section 5 tells that the size of \tilde{D} is of smaller order than the normalized error D of the covariance estimator based on Z_t . It is known that the latter converges to a non-degenerate law as $T \rightarrow \infty$. This is well supported by the plots in Figure 3 showing that the distance between the two thick lines in each panel is almost invariant as T increases. The fact that the additional error incurred by using \tilde{Z}_t instead of Z_t is negligible for large T is also confirmed. In particular, the long stretches at tails of the distributions of \tilde{D} get shorter as T increases. Also, the upper and lower quartiles of each entry of \tilde{D} , represented by the boxes, lie within those of the corresponding entry of D , represented by the thick lines, when $T = 1,000$ and $2,000$.

4. APPLICATIONS

This section presents an application of DSFM. We fit the model to the intraday IV based on ODAX prices and to fMRI data.

For our analysis we chose the data observed from July 1, 2004 to June 29, 2005. The one year period corresponds to the financial regulatory requirements. The data were taken from Financial and Economic Data Center of Humboldt-Universität zu Berlin. The IV data were regressed on the two-dimensional space of future moneyness and time-to-maturity, denoted by $(\kappa_t, \tau_t)^\top$. The future moneyness κ_t is a monotone function of the strike price K : $\kappa_t = K/(S_t e^{-r_t \tau_t})$, where S_t is the spot price at time t and r_t is the interest rate. We chose r_t as a daily Euro Interbank Offered Rate (EURIBOR) taken from the Ecwin Reuters database. The time-to-maturity of the options were measured in years. We took all trades with $10/365 < \tau < 0.5$. We limit also the moneyness range to $\kappa \in [0.7, 1.2]$.

The structure of the IV data, described already in Section 1, requires a careful treatment. Apart from the dynamic degeneration, one may also observe nonuniform frequency of the trades with significantly greater market activities for the options closer to expiry or at-the-money. Here, ‘‘at-the-money’’ means a condition in which the strike price of an option equals the spot price of the underlying security (i.e., $K = S_t$). To avoid the computational problems with the highly skewed empirical distribution of $X_t = (\kappa_t, \tau_t)$, we transformed the initial space $[0.7, 1.2] \times [0.03, 0.5]$ to $[0, 1]^2$ by using the marginal empirical distribution functions. We applied the estimation algorithm to the transformed space, and then transformed back the results to the original space.

Because the model is not nested, the number of the dynamic functions needs to be determined in advance. For this, we used

$$RV(L) = \frac{\sum_t^T \sum_j^{J_t} \left\{ Y_{t,j} - \hat{m}_0(X_{t,j}) - \sum_{l=1}^L \hat{Z}_{t,l} \hat{m}_l(X_{t,j}) \right\}^2}{\sum_t^T \sum_j^{J_t} (Y_{t,j} - \bar{Y})^2}, \quad (12)$$

although one may construct an Akaike information (AIC) or Bayesian information (BIC) type of criterion, where one penalizes the number of the dynamic functions in the model, or performs some type of cross-validation. The quantity $1 - RV(L)$ can be interpreted as a proportion of the variation explained by

Table 1. Proportion of the explained variation by the models with $L = 1, \dots, 5$ dynamic factors

No. factors	$L = 1$	$L = 2$	$L = 3$	$L = 4$	$L = 5$
$1 - RV(L)$	0.848	0.969	0.976	0.978	0.980

the model among the total variation. The computed values of $RV(L)$ are given in Table 1 for various L . Because the third, fourth, and fifth factor made only a small improvement in the fit, we chose $L = 2$.

For the series estimators of \hat{m}_l we used tensor B-splines that are cubic in the moneyness and quadratic in the maturity direction. In the transformed space we placed 10×5 knots, 10 in the moneyness and 5 in the maturity direction. We found that the results were not sensitive to the choice of the number of knots and the orders of splines. For several choices of knots in the range $5 \times 5 - 15 \times 10$ and for the spline orders (2, 1), (2, 2), (3, 2), the values of $1 - RV(2)$ were between 0.949 and 0.974. Because the model is identifiable only up to the transformation (8), one has a freedom for the choice of factors. Here, we chose the approach taken by Fengler et al. (2007) with $L_2[0,1]^2$ norm. Specifically, we orthonormalized \hat{m}_l and transformed \hat{Z}_t according to their Equation (19) with $\Gamma = \int \hat{m}(x) \hat{m}(x)^\top dx$, where $\hat{m} = (\hat{m}_1, \dots, \hat{m}_L)^\top$. Call them \hat{m}_l^* and \hat{Z}_t^* , respectively. Then, we transformed them further by $\hat{m}_l^{**} = p_l^\top \hat{m}_l^*$ and $\hat{Z}_{t,l}^{**} = p_l^\top \hat{Z}_t^*$, where p_l were the orthonormal eigenvectors of the matrix $\sum_{t=1}^T \hat{Z}_t^* \hat{Z}_t^{*\top}$ that correspond to the eigenvalues $\lambda_1 > \lambda_2$. Note that $\hat{Z}_t^{*\top} \hat{m}^* = \hat{Z}_t^{**\top} \hat{m}^{**}$. In this way, $\{\hat{Z}_{t,1}^{**} \hat{m}_1^{**}\}$ makes a larger contribution than $\{\hat{Z}_{t,2}^{**} \hat{m}_2^{**}\}$ to the total variation $\sum_{t=1}^T \int (\hat{Z}_t^{**\top} \hat{m}^{**})^2$ because $\sum_{t=1}^T \int (\hat{Z}_{t,1}^{**} \hat{m}_1^{**})^2 = \lambda_1$ and $\sum_{t=1}^T \int (\hat{Z}_{t,2}^{**} \hat{m}_2^{**})^2 = \lambda_2$. Later, we continue to write \hat{Z}_t and \hat{m} for such \hat{Z}_t^{**} and \hat{m}^{**} , respectively.

The estimated functions \hat{m}_1 and \hat{m}_2 are plotted in Figure 4 in the transformed estimation space. The intercept function \hat{m}_0 was almost flat around zero, thus is not given. By construction, $\hat{m}_0 + \hat{Z}_{t,1} \hat{m}_1$ explain the principal movements of the surface. It was observed by Cont and da Fonseca (2002) and Fengler et al. (2007) that most dominant innovations of the entire surface are parallel level shifts. Note that VDAX is an estimated at-the-money IV for an option with 45 days to maturity, and thus indicates up-and-down shifts. The left panel of Figure 5 shows the values of VDAX together with $\hat{m}_0(X_{t,0}) + \hat{Z}_{t,1} \hat{m}_1(X_{t,0})$, where $X_{t,0}$ is the moneyness and maturity corresponding to an option at-the-money with 45 days to maturity. The right panel of Figure 5 depicts the factor \hat{Z}_t , where one can find that \hat{Z}_t shows almost the same dynamic behavior as the index VDAX. This similarity supports that DSFM catches leading dynamic effects successfully. Obviously the model in its full setting explains other effects, such as skew or term structure changes, which are not explicitly stated here.

Statistical analysis on the evolution of a high-dimensional system ruling the option prices can be simplified to a low-dimensional analysis of the \hat{Z}_t . In particular, as our theory in Section 5 and the simulation results in Section 3 assert, the inference based on the \hat{Z}_t is well justified in the VAR context. To select a VAR model we computed the Schwarz (SC), the Hannan-Quinn (HQ), and the Akaike criterion, as given in

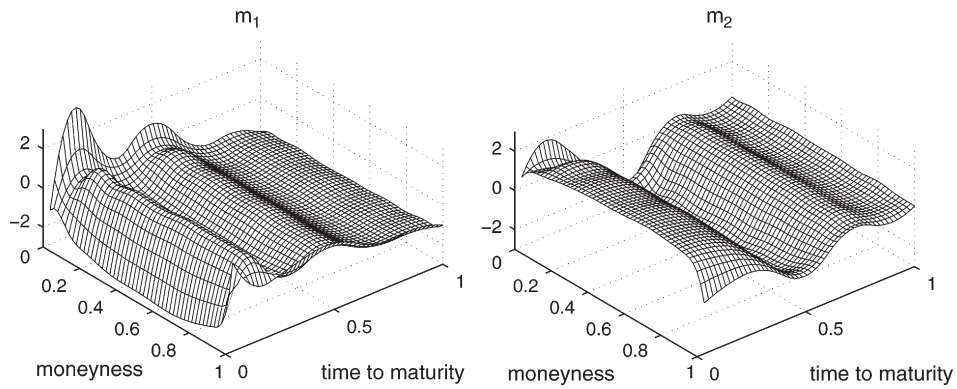


Figure 4. The estimated factor functions for the ODAX IV data in the period 20040701–20050629.

Table 2. One can find that SC and HQ suggest a VAR(1) process, whereas AIC selects VAR(2). The parameter estimates for each selected model are given in Table 3. The roots of the characteristic polynomial lie inside the unit circle, so the specified models satisfy the stationarity condition. For each of VAR(1) and VAR(2) models, we conducted a portmanteau test for the hypothesis that the autocorrelations of the error term at lags up to 12 are all zero, and also a series of LM tests, each of which tests whether the autocorrelation at a particular lag up to 5 equals zero. Some details on selection of lags for these tests can be found in Hosking (1980, 1981) and Brüggemann, Lütkepohl, and Saikkonen (2006). We found that in any test the null hypothesis was not rejected at 5% level. A closer inspection on the autocorrelations of the residuals, however, revealed that the autocorrelation of $\widehat{Z}_{t,2}$ residuals at lag one is slightly significant in the VAR(1) model, see Figure 6. But, this effect disappears in the VAR(2) case, see Figure 7. Similar analyses of characteristic polynomials, portmanteau and Lagrange multiplier (LM) tests supported VAR(2) as a successful model for \widehat{Z}_t .

As a second application of the model, we considered fitting an fMRI dataset. The data were obtained at Max-Planck Institut für Kognitions-und-Neurowissenschaften Leipzig by scanning a subject’s brain using a standard head coil. The scanning was done every two seconds on the resolution of $3 \times 3 \times 2 \text{ mm}^3$ with 1 mm gap between the slices. During the experiment, the

subject was exposed to three types of objects (bench, phone and motorbike) and rotated around randomly changing axes for four seconds, followed by relaxation phase of six to ten seconds. Each stimulus was shown 16 times in pseudo-randomized order. As a result, a series of 290 images with $64 \times 64 \times 30$ voxels was obtained.

To apply the model (2) to the fMRI data, we took the voxel’s index (i_1, i_2, i_3) as covariate $X_{t,j}$, and the BOLD signal as $Y_{t,j}$. For numerical tractability we reduced the original data to a series of $32 \times 32 \times 15$ voxels by taking every second slice in each direction. Thus, $J_t \equiv 32 \times 32 \times 15$ and $T = 290$. The voxels’ indices (i_1, i_2, i_3) for $1 \leq i_1, i_2 \leq 32 ; 1 \leq i_3 \leq 15$ are associated with $32 \times 32 \times 15$ equidistant points in \mathbb{R}^3 . The function m_0 represents the “average” signal as a function of the three-dimensional location, and m_l for each $l \geq 1$ determines the effect of the l th common factor $Z_{t,l}$ on the brain’s signal. In Figure 8, each estimated function \widehat{m}_l is represented by its sections on the 15 slices in the direction of i_3 [i.e., by those $\widehat{m}_l(\cdot, \cdot, x_3)$ for which x_3 are fixed at the equidistant points corresponding to $i_3 = 1, \dots, 15$]. We used quadratic tensor B-splines on equidistant knots. The number of knots in each direction was 8, 8, 4, respectively, so that $K = 9 \times 9 \times 5 = 405$. For the model identification we used the same method as in the IV application, but normalized \widehat{Z} to have mean zero.

In contrast to the IV application, there was no significant difference between the values of $1 - RV(L)$ for different $L \geq 1$.

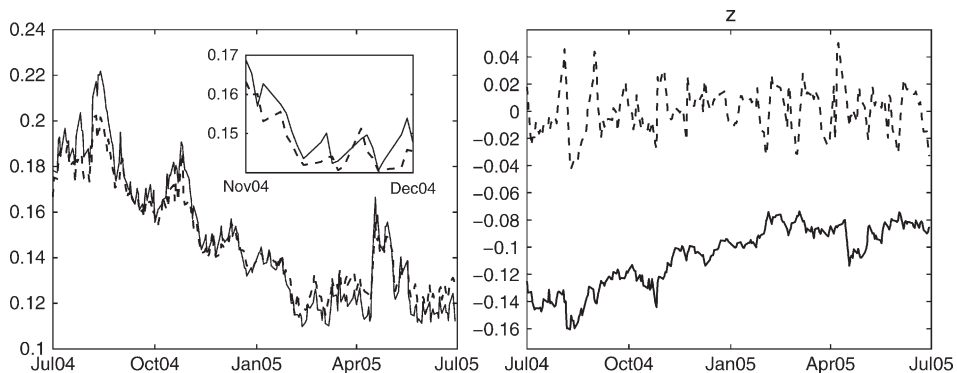


Figure 5. Left panel: VDX in the period 20040701–20050629 (solid) and the dynamics of the corresponding IV given by the submodel $\widehat{m}_0 + \widehat{Z}_{t,1}\widehat{m}_1$ (dashed). Right panel: The obtained time series \widehat{Z}_t on the ODAX IV data in the period 20040701–20050629. The solid line represents $\widehat{Z}_{t,1}$, the dashed line $\widehat{Z}_{t,2}$.

Table 2. The VAR model selection criteria. The smallest value for each criterion is marked by an asterisk

Order	AIC	SC	HQ
1	-14.06	-13.98*	-14.03*
2	-14.07*	-13.93	-14.02
3	-14.06	-13.86	-13.98
4	-14.06	-13.81	-13.96
5	-14.07	-13.76	-13.95

All the values for $L \geq 1$ were around 0.871. The fMRI signals $Y_{t,j}$ were explained mostly by $\hat{m}_0(X_{t,j}) + Z_{t,1}\hat{m}_1(X_{t,j})$, and the effects of the common factors $Z_{t,l}$ for $l \geq 2$ were relatively small. The slow increase in the value of $1 - RV(L)$ as $L \geq 1$ grows in the fMRI application, contrary to the case of the IV application, can be explained partly by the high complexity of human brain. Because the values of $1 - RV(L)$ were similar for $L \geq 1$, one might choose $L = 1$. However, we chose $L = 4$, which we think still allows relatively low complexity, to demonstrate some further analysis that might be possible with similar datasets. The estimated functions \hat{m}_l for $0 \leq l \leq 4$ and the time series $\hat{Z}_{t,l}$ for $1 \leq l \leq 4$ are plotted in Figures 8 and 9, respectively. The function \hat{m}_0 can be recognized as a smoothed version of the original signal. By construction the first factor and loadings incorporate the largest variation. One may see the strong positive trend in $\hat{Z}_{t,1}$ and relatively flat patterns of $\hat{Z}_{t,2}, \hat{Z}_{t,3}, \hat{Z}_{t,4}$. These effects could be typically explained by the mixture of several components, such as physiological pulsation, subtle head movement, machine noise, and so on. For a description of different artifacts, which significantly influence the fMRI signals, we refer to Biswal, Yetkin, Haughton, and Hyde (1995). The function estimates \hat{m}_l for $1 \leq l \leq 4$ appear to have a clear peak, and $\hat{Z}_{t,l}$ for $2 \leq l \leq 4$ show rather mild mean reverting behavior.

To see how the recovered signals interact with the given stimuli, we plotted $\hat{Z}_{t+s,l} - \hat{Z}_{s,l}$ against t in Figure 10, where s is the time when a stimulus appears. The mean changes of $\hat{Z}_{t,1}$ and $\hat{Z}_{t,3}$ show mild similarity, up to sign change, to the hemodynamic response (see Worsley et al. 2002). The case of $\hat{Z}_{t,4}$ has a similar pattern as those of $\hat{Z}_{t,1}$ and $\hat{Z}_{t,3}$ but with larger amplitude, whereas the changes in $\hat{Z}_{t,2}$ seem to be independent of the stimuli. In fitting the fMRI data, we did not use any external information on the signal. From the biological perspective it could be hardly expected that a pure statistical procedure gives full insight into understanding of the complex dynamics of MR images. For the latter one needs to incorporate into the procedure the shape of hemodynamic response, for example, or consider physiologically motivated identification of the fac-

tors. It goes however beyond the scope of this illustrative example.

5. ASYMPTOTIC ANALYSIS

In the simulation study and the real data application in Sections 3 and 4, we considered the case where Z_t is a VAR-process. Here, we only make some weak assumptions on the average behavior of the process. In our first theorem we allow that it is a deterministic sequence. In our second result we assume that it is a mixing sequence. For the asymptotic analysis, we let $K, J, T \rightarrow \infty$. This is a very natural assumption often also made in cross-sectional or panel data analysis. It is appropriate for data with many observations per data point that are available for many dates. It allows us to study how J and T have to grow with respect to each other for a good performance of a procedure. The distance between m and its best approximation $\mathcal{A}\psi$ does not tend to zero unless $K \rightarrow \infty$, see Assumption (A5) later. One needs to let $J \rightarrow \infty$ to get consistency of $\hat{Z}_t^\top \hat{A}$ and $\hat{m} = \hat{A}\psi$ as estimates of $Z_t^\top \mathcal{A}^*$ and m , respectively, where \mathcal{A}^* is defined at (A5). One should let $T \rightarrow \infty$ to describe the asymptotic equivalence between the lagged covariances of Z_t and those of \hat{Z}_t , see Theorem 3 below. In our analysis the dimension L is fixed. Clearly, one could also study our model with L growing to infinity. We treat the case where X_{it} are random. However, a theory for deterministic designs can be developed along the lines of our theory.

Our first result relies on the following assumptions.

(A1) The variables $X_{1,1}, \dots, X_{T,J}, \varepsilon_{1,1}, \dots, \varepsilon_{T,J}$, and Z_1, \dots, Z_T are independent. The process Z_t is allowed to be nonrandom.
 (A2) For $t = 1, \dots, T$ the variables $X_{t,1}, \dots, X_{t,J}$ are identically distributed, have support $[0, 1]^d$ and a density f_t that is bounded from below and above on $[0, 1]^d$, uniformly over $t = 1, \dots, T$.

(A3) We assume that $E\varepsilon_{t,j} = 0$ for $1 \leq t \leq T, 1 \leq j \leq J$, and for $c > 0$ small enough $\sup_{1 \leq t \leq T, 1 \leq j \leq J} E \exp(c\varepsilon_{t,j}^2) < \infty$.

(A4) The functions ψ_k may depend on the increasing indices T and J , but are normed so that $\int_{[0,1]^d} \psi_k^2(x) dx = 1$ for $k = 1, \dots, K$. Furthermore, it holds that $\sup_{x \in [0,1]^d} \|\psi(x)\| = \mathcal{O}(K^{1/2})$.

(A5) The vector of functions $m = (m_0, \dots, m_L)^\top$ can be approximated by ψ_k , i.e.,

$$\delta_K \equiv \sup_{x \in [0,1]^d} \inf_{\mathcal{A} \in \mathbb{R}^{(L+1) \times K}} \|m(x) - \mathcal{A}\psi(x)\| \rightarrow 0$$

as $K \rightarrow \infty$. We denote \mathcal{A} that fulfills $\sup_{x \in [0,1]^d} \|m(x) - \mathcal{A}\psi(x)\| \leq 2\delta_K$ by \mathcal{A}^* .

(A6) There exist constants $0 < C_L < C_U < \infty$ such that all eigenvalues of the matrix $T^{-1} \sum_{t=1}^T Z_t Z_t^\top$ lie in the interval $[C_L, C_U]$ with probability tending to one.

Table 3. The estimated parameters for VAR(1) and VAR(2) models. Those that are not significant at 5% level are marked by asterisk

	VAR(1)			VAR(2)				
	$\hat{Z}_{t-1,1}$	$\hat{Z}_{t-1,2}$	Const.	$\hat{Z}_{t-1,1}$	$\hat{Z}_{t-1,2}$	$\hat{Z}_{t-2,1}$	$\hat{Z}_{t-2,2}$	Const.
$\hat{Z}_{t,1}$	0.984	-0.029*	-0.001	0.913	-0.025	0.071	-0.004	-0.001
$\hat{Z}_{t,2}$	0.055	0.739	0.005	0.124	0.880	-0.065	-0.187*	0.006

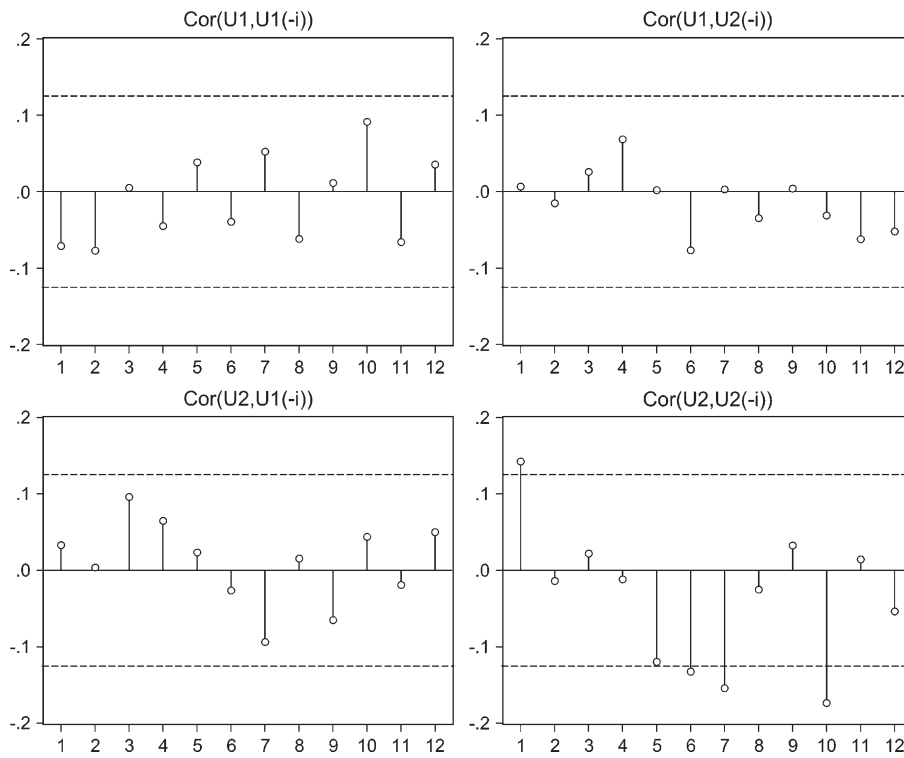


Figure 6. Cross-autocorrelogram for the VAR(1) residuals. The dashed line-bounds indicate $\pm 2 \times$ (standard deviations), which correspond to an approximate 95% confidence bound.

(A7) The minimization (5) runs over all values of (\mathcal{A}, z) with

$$\sup_{x \in [0,1]} \max_{1 \leq t \leq T} \| (1, z_t^T) \mathcal{A} \psi(x) \| \leq M_T,$$

where the constant M_T fulfils $\max_{1 \leq t \leq T} \|Z_t\| \leq M_T/C_m$ (with probability tending to one) for a constant C_m such that $\sup_{x \in [0,1]} \|m(x)\| < C_m$.

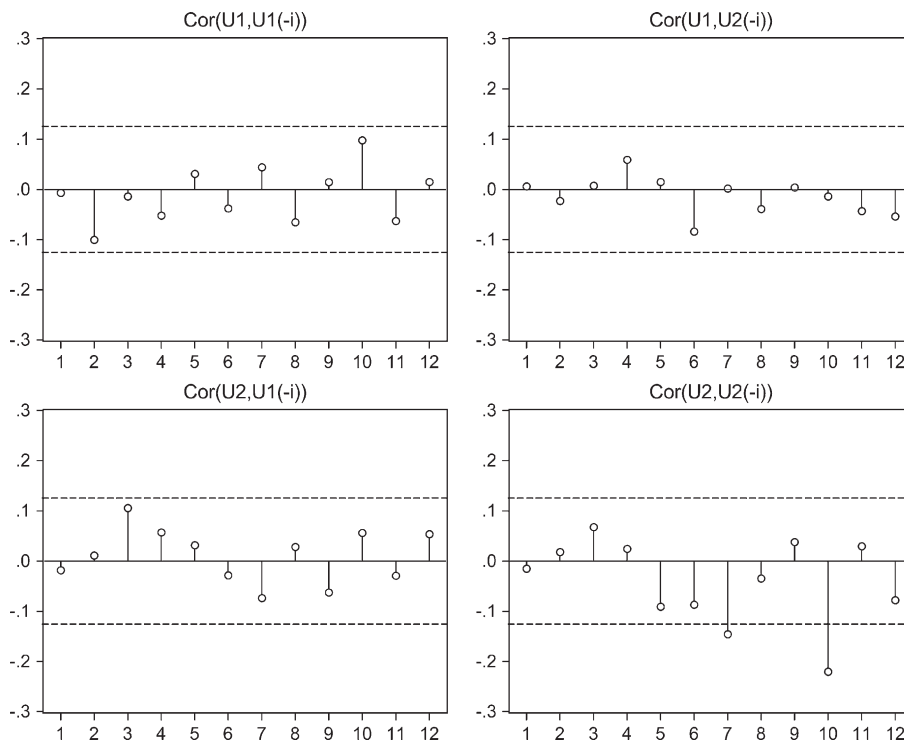


Figure 7. Cross-autocorrelogram for the VAR(2) residuals. The dashed line-bounds indicate $\pm 2 \times$ (standard deviations), which correspond to an approximate 95% confidence bound.

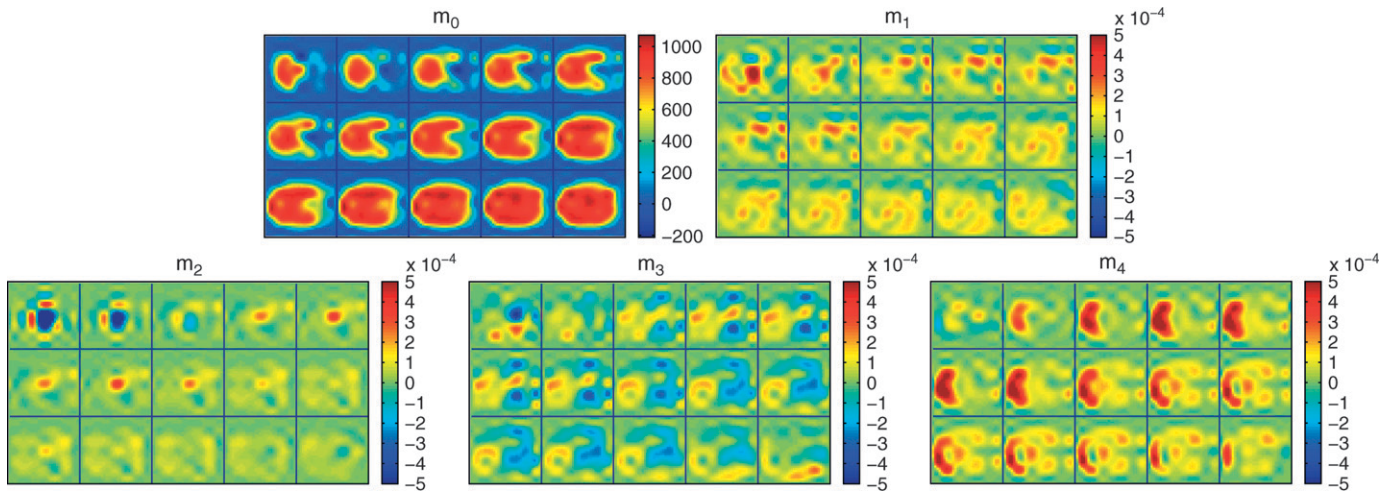


Figure 8. The estimated functions \hat{m}_i for the fMRI signals.

(A8) It holds that $\rho^2 = (K + T)M_T^2 \log(JTM_T)/(JT) \rightarrow 0$. The dimension L is fixed.

Assumption (A7) and the additional bound M_T in the minimization is introduced for purely technical reasons. We conjecture that to some extent the asymptotic theory of this article could be developed under weaker conditions. The independence assumptions in (A1) and Assumption (A3) could be relaxed to assuming that the errors $\epsilon_{t,j}$ have a conditional mean zero and have a conditional distribution with subgaussian tails, given the past values $X_{s,i}, Z_s$ ($1 \leq i \leq J, 1 \leq s \leq t$). Such a theory would

require an empirical process theory that is more explicitly designed for our model and it would also require a lot of more technical assumptions. We also expect that one could proceed with the assumption of subexponential instead of subgaussian tails, again at the cost of some additional conditions. Recall that the number of parameters to be estimated equals $TL + K(L + 1)$. Because L is fixed, Assumption (A8) requires basically that, neglecting the factor $M_T^2 \log(JTM_T)$, the number of parameters grows slower than the number of observations, JT .

Our first result gives rates of convergence for the least squares estimators \hat{Z}_t and \hat{A} .

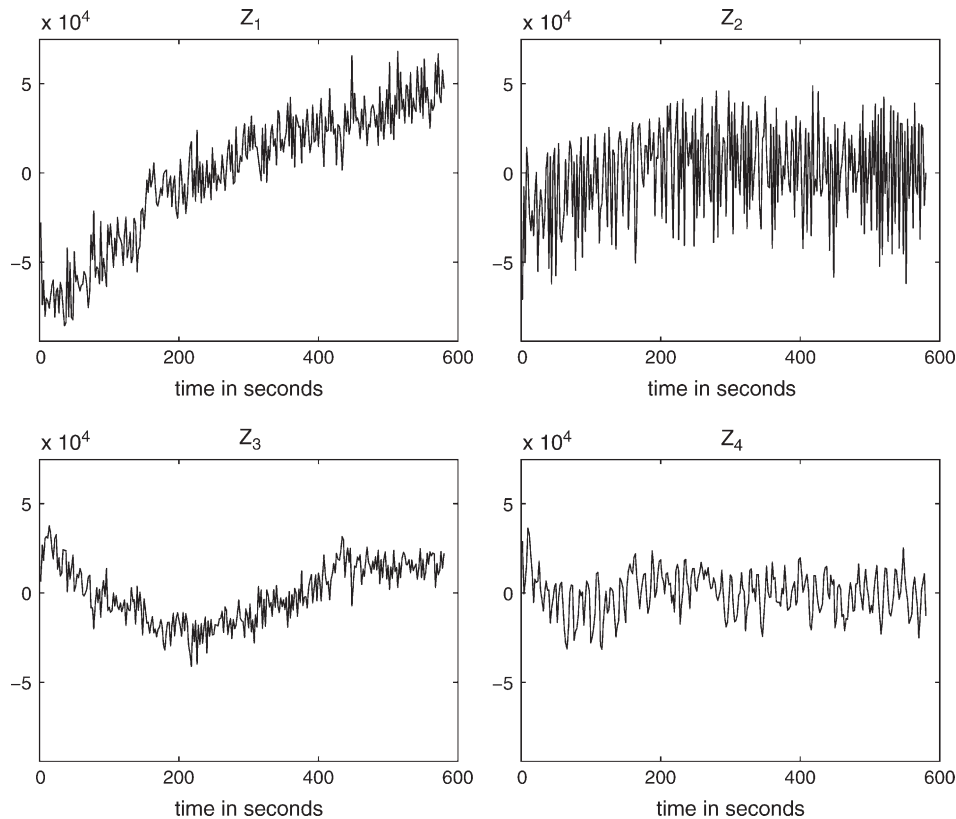


Figure 9. The estimated time series $\hat{Z}_{t,l}$ for the fMRI signals.

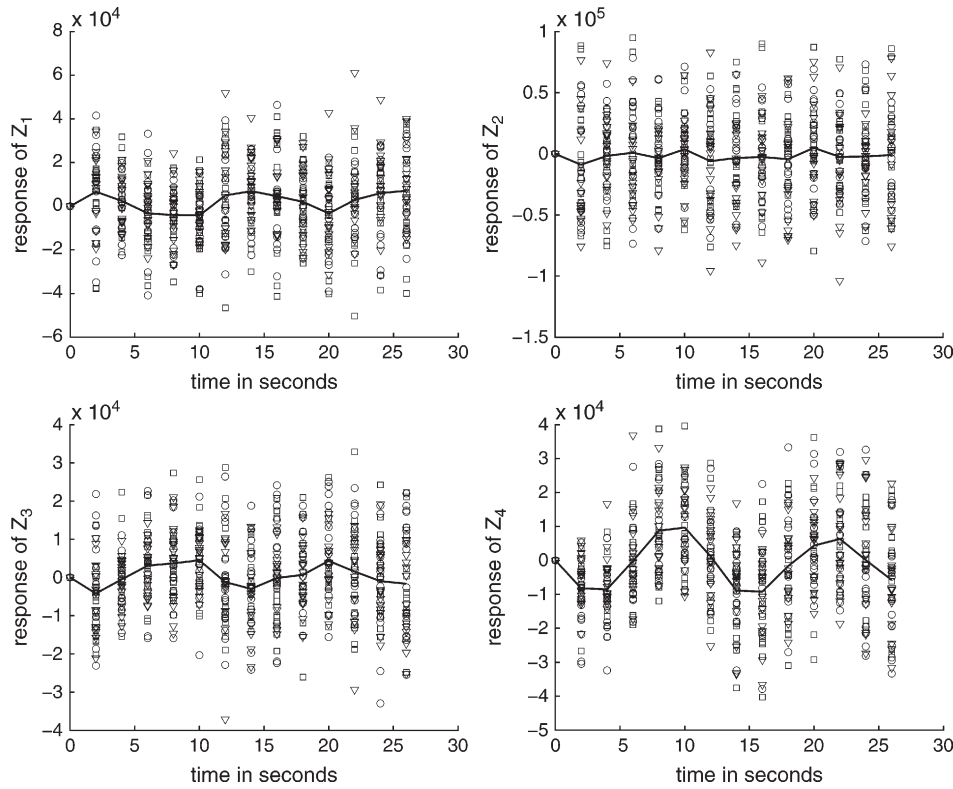


Figure 10. The responses of $\hat{Z}_{t,t}$ to the stimuli.

Theorem 2. Suppose that model (4) holds and that (\hat{Z}_t, \hat{A}) is defined by the minimization problem (5). Make the Assumptions (A1)–(A8). Then it holds that

$$\frac{1}{T} \sum_{1 \leq t \leq T} \left\| \hat{Z}_t^\top \hat{A} - Z_t^\top A^* \right\|^2 = \mathcal{O}_P(\rho^2 + \delta_K^2). \quad (13)$$

At this point we have made no assumptions on the sequence Z_t , $1 \leq t \leq T$, besides the bound in (A7). Up to now it is allowed to be a deterministic or a random sequence. We now assume that it is a random process. We discuss how a statistical analysis differs if inference on Z_t is based on \hat{Z}_t instead of using (the unobserved) process Z_t . We will show that the differences are asymptotically negligible (except an orthogonal transformation). This is the content of the following theorem, where we consider estimators of autocovariances and show that these estimators differ only by second order terms. This asymptotic equivalence carries over to classical estimation and testing procedures in the framework of fitting a vector autoregressive model. For the statement of the theorem we need the following assumptions:

- (A9) Z_t is a strictly stationary sequence with $E(Z_t) = 0$, $E(\|Z_t\|^\gamma) < \infty$ for some $\gamma > 2$. It is strongly mixing with $\sum_{i=1}^\infty \alpha(i)^{(\gamma-2)/\gamma} < \infty$. The matrix $E Z_t Z_t^\top$ has full rank. The process Z_t is independent of $X_{11}, \dots, X_{TJ}, \epsilon_{11}, \dots, \epsilon_{TJ}$.
- (A10) The functions m_0, \dots, m_L are linearly independent. In particular, no function is equal to 0.
- (A11) It holds that $[\log(KT)^2 \{ (KM_T/J)^{1/2} + T^{1/2} M_T^4 J^{-2} + K^{3/2} J^{-1} + K^{4/3} J^{-2/3} T^{-1/6} \} + 1] T^{1/2} (\rho^2 + \delta_K^2) = o(1)$.

Assumption (A11) poses very weak conditions on the growth of J , K , and T . Suppose, for example, that M_T is of logarithmic

order and that K is of order $(TJ)^{1/5}$ so that the variance and the bias are balanced for twice differentiable functions. In this setting, (A11) only requires that T/J^2 times a logarithmic factor converges to zero. Define $\tilde{Z}_t = B^\top \hat{Z}_t$,

$$\begin{aligned} \tilde{Z}_{c,t} &= \tilde{Z}_t - T^{-1} \sum_{s=1}^T \tilde{Z}_s, \\ Z_{c,t} &= Z_t - T^{-1} \sum_{s=1}^T Z_s, \\ \tilde{Z}_{n,t} &= (T^{-1} \sum_{s=1}^T \tilde{Z}_{c,s} \tilde{Z}_{c,s}^\top)^{-1/2} \tilde{Z}_{c,t}, \end{aligned}$$

$$\text{and } Z_{n,t} = (T^{-1} \sum_{s=1}^T Z_{c,s} Z_{c,s}^\top)^{-1/2} Z_{c,t}.$$

Theorem 3. Suppose that model (4) holds and that (\hat{Z}_t, \hat{A}) is defined by the minimization problem (5). Make the Assumptions (A1)–(A11). Then there exists a random matrix B such that for $h \neq 0$

$$\begin{aligned} \frac{1}{T} \sum_{t=\max[1,-h+1]}^{\min[T,T-h]} \tilde{Z}_{c,t} (\tilde{Z}_{c,t+h} - \tilde{Z}_{c,t})^\top - \frac{1}{T} \sum_{t=\max[1,-h+1]}^{\min[T,T-h]} Z_{c,t} \\ (Z_{c,t+h} - Z_{c,t})^\top &= o_P(T^{-1/2}), \\ \frac{1}{T} \sum_{t=\max[1,-h+1]}^{\min[T,T-h]} \tilde{Z}_{n,t} \tilde{Z}_{n,t+h}^\top - \frac{1}{T} \sum_{t=\max[1,-h+1]}^{\min[T,T-h]} Z_{n,t} Z_{n,t+h}^\top &= o_P(T^{-1/2}). \end{aligned}$$

To illustrate an implication of Theorem 3, suppose that the factor process Z_t in (4) is a stationary VAR(p) process in a mean adjusted form:

$$Z_t - \mu = \Theta_1(Z_{t-1} - \mu) + \dots + \Theta_p(Z_{t-p} - \mu) + U_t, \quad (14)$$

where $\mu = E(Z_t)$, Θ_j is a $L \times L$ matrix of coefficients and U_t is a white noise with a nonsingular covariance matrix. Let Γ_h be the autocovariance matrix of the process Z_t with the lag $h \geq 0$, which is estimated by $\hat{\Gamma}_h = T^{-1} \sum_{t=h+1}^T (Z_t - \bar{Z})(Z_{t-h} - \bar{Z})^\top$. Let $Y = (Z_{p+1} - \mu, \dots, Z_T - \mu)$, $\Theta = (\Theta_1, \dots, \Theta_p)$, and $U = (U_{p+1}, \dots, U_T)$. Define $W_t = ((Z_t - \mu)^\top, \dots, (Z_{t-p+1} - \mu)^\top)^\top$ and $W = (W_p, \dots, W_{T-1})$. Then, the model (14) can be rewritten as $Y = \Theta W + U$ and the least squares estimator of Θ is given by $\hat{\Theta} = \hat{Y} \hat{W}^\top (\hat{W} \hat{W}^\top)^{-1}$, where \hat{Y} and \hat{W} are the same as Y and W , respectively, except that μ is replaced by \bar{Z} . Likewise, fitting a VAR(p) model with the estimated factor process \tilde{Z}_t yields $\hat{\Theta} = \hat{Y} \hat{W}^\top (\hat{W} \hat{W}^\top)^{-1}$, where \hat{Y} and \hat{W} are defined as \hat{Y} and \hat{W} with Z_t being replaced by \tilde{Z}_t . Both \hat{Y} and \hat{W} are matrices composed of $\hat{\Gamma}_h$ for various h . The matrices \hat{Y} and \hat{W} have the same forms as \hat{Y} and \hat{W} , respectively, but with $\hat{\Gamma}_h$ being replaced by $\hat{\Gamma}_h = T^{-1} \sum_{t=h+1}^T (\tilde{Z}_t - \bar{\tilde{Z}})(\tilde{Z}_{t-h} - \bar{\tilde{Z}})^\top$. It is well known that $\sqrt{T}(\hat{\Theta} - \Theta) = \mathcal{O}_P(1)$, see Lütkepohl (1993). By Theorem 3, we have $\sqrt{T}(\hat{\Theta} - \Theta) = o_P(1)$.

APPENDIX: PROOFS OF THEOREMS

A.1 Proof of Theorem 1

We use the Newton-Kantorovich theorem to prove the theorem. The statement of the theorem may be found in Kantorovich and Akilov (1982), for example.

Suppose that $\sum_{t=1}^T \|\mathcal{Z}_t^{(0)\top} \mathcal{A}^{(0)} - \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}}\|^2 \leq r$ for some $r > 0$, which will be chosen later. With the Frobenius norm $\|M\|$ for a matrix M , we get

$$\begin{aligned} \|\mathcal{A}^{(0)} - \hat{\mathcal{A}}\|^2 &\leq \left\| \left(\sum_{t=1}^T \mathcal{Z}_t^{(0)} \mathcal{Z}_t^{(0)\top} \right)^{-1} \right\|^2 \\ &\quad \left\| \sum_{t=1}^T \mathcal{Z}_t^{(0)} \mathcal{Z}_t^{(0)\top} (\mathcal{A}^{(0)} - \hat{\mathcal{A}}) \right\|^2 \\ &= \left\| \left(\sum_{t=1}^T \mathcal{Z}_t^{(0)} \mathcal{Z}_t^{(0)\top} \right)^{-1} \right\|^2 \cdot \left\| \sum_{t=1}^T \mathcal{Z}_t^{(0)} \mathcal{Z}_t^{(0)\top} \mathcal{A}^{(0)} \right. \\ &\quad \left. - \sum_{t=1}^T \mathcal{Z}_t^{(0)} \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}} \right\|^2 \leq \left\| \left(\sum_{t=1}^T \mathcal{Z}_t^{(0)} \mathcal{Z}_t^{(0)\top} \right)^{-1} \right\|^2 \\ &\quad \times \left(\sum_{t=1}^T \left\| \mathcal{Z}_t^{(0)} \mathcal{Z}_t^{(0)\top} \mathcal{A}^{(0)} - \mathcal{Z}_t^{(0)} \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}} \right\|^2 \right) \leq \\ &\quad r \left\| \left(\sum_{t=1}^T \mathcal{Z}_t^{(0)} \mathcal{Z}_t^{(0)\top} \right)^{-1} \right\|^2 \left(\sum_{t=1}^T \|\mathcal{Z}_t^{(0)}\|^2 \right) \\ &= rc_1^2. \end{aligned} \quad (A.1)$$

For a matrix M , define $\|M\|_2 = \sup_{\|x\|=1} \|Mx\|$. It is known that $\|M\|_2 \leq \|M\|$. We get

$$\begin{aligned} \|\hat{\mathcal{A}}^\top (\mathcal{Z}_t^{(0)} - \hat{\mathcal{Z}}_t)\| &\geq \|\hat{\mathcal{A}}\|_2^{-1} \\ &\quad \cdot \|(\hat{\mathcal{A}} \hat{\mathcal{A}}^\top)^{-1}\|^{-1} \cdot \|\mathcal{Z}_t^{(0)} - \hat{\mathcal{Z}}_t\|, \end{aligned} \quad (A.2)$$

$$\begin{aligned} \|(\mathcal{Z}_t^{(0)} - \hat{\mathcal{Z}}_t)^\top \hat{\mathcal{A}}\| &\leq \|\mathcal{Z}_t^{(0)\top} (\hat{\mathcal{A}} - \mathcal{A}^{(0)})\| + \|\mathcal{Z}_t^{(0)\top} \mathcal{A}^{(0)} - \\ \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}}\| &\leq \|\mathcal{Z}_t^{(0)}\| \cdot \|\hat{\mathcal{A}} - \mathcal{A}^{(0)}\| + \|\mathcal{Z}_t^{(0)\top} \mathcal{A}^{(0)} - \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}}\|. \end{aligned} \quad (A.3)$$

The two inequalities (A.2) and (A.3) together with (A.1) give

$$\begin{aligned} \|\mathcal{Z}^{(0)} - \hat{\mathcal{Z}}\|^2 &\leq 2r \|\hat{\mathcal{A}}\|_2^2 \cdot \|(\hat{\mathcal{A}} \hat{\mathcal{A}}^\top)^{-1}\|^2 \\ &\quad \times \left(1 + c_1 \sum_{t=1}^T \|\mathcal{Z}_t^{(0)}\|^2 \right) = rc_2^2. \end{aligned} \quad (A.4)$$

Because $F'(\alpha, z)$ is quadratic in (α, z) , there exists $0 < c_3 < \infty$ for any compact set D in $\mathbb{R}^{K(L+1)+TL}$ such that $\|F'(\alpha', z') - F'(\alpha, z)\|_2 \leq c_3 \|(\alpha'^\top, z'^\top)^\top - (\alpha^\top, z^\top)^\top\|$ for all $(\alpha'^\top, z'^\top)^\top, (\alpha^\top, z^\top)^\top \in D$. Let $c_4 = \|F'_*(\alpha^{(0)}, Z^{(0)})^{-1}\|_2 < \infty$. Because F is continuous and $F(\hat{\alpha}, \hat{Z}) = 0$, there exists $r' > 0$ such that, if $\|\alpha^{(0)} - \hat{\alpha}\| + \|Z^{(0)} - \hat{Z}\| \leq r'$, then

$$\|F'_*(\alpha^{(0)}, Z^{(0)})^{-1} F(\alpha^{(0)}, Z^{(0)})\| \leq \frac{\gamma}{2c_3 c_4}.$$

By the Newton-Kantorovich theorem,

$$\|\alpha^{(k)} - \hat{\alpha}\| + \|Z^{(k)} - \hat{Z}\| \leq C_1 2^{-(k-1)} \gamma^{2^{k-1}} \quad (A.5)$$

for some $C_1 > 0$. This gives that if $\|\alpha^{(0)} - \hat{\alpha}\| + \|Z^{(0)} - \hat{Z}\| \leq r'$, then

$$\begin{aligned} \sum_{t=1}^T \|\mathcal{Z}_t^{(k)\top} \mathcal{A}^{(k)} - \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}}\|^2 &\leq C_2 (\|\alpha^{(k)} - \hat{\alpha}\|^2 + \\ \|\mathcal{Z}^{(k)} - \hat{\mathcal{Z}}\|^2) &\leq C_2^{-2(k-1)} \gamma^{2^{2(k-1)}} \end{aligned}$$

for some $C, C_2 > 0$. We take $r = (c_1 + c_2)^{-2} r'^2$. Then, by (A.1) and (A.4), $\|\alpha^{(0)} - \hat{\alpha}\| + \|Z^{(0)} - \hat{Z}\| \leq r'$ if $\sum_{t=1}^T \|\mathcal{Z}_t^{(0)\top} \mathcal{A}^{(0)} - \hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}}\|^2 \leq r$. This completes the proof of the theorem.

A.2 Proof of Theorem 2

For functions $g(t, x)$ we define the norms $\|g\|_1^2 = (1/TJ) \sum_{t=1}^T \sum_{j=1}^J g(t, X_{t,j})^2$, $\|g\|_2^2 = (1/T) \sum_{t=1}^T \int g(t, x)^2 f_t(x) dx$, and $\|g\|_3^2 = (1/T) \sum_{t=1}^T \int g(t, x)^2 dx$. Note that because of Assumption (A2) the last two norms are equivalent. Thus, for the statement of the theorem we have to show for $\Delta(t, x) = (\hat{\mathcal{Z}}_t^\top \hat{\mathcal{A}} - \mathcal{Z}_t^\top \mathcal{A}^*) \psi(x)$ that

$$\|\Delta\|_2^2 = \mathcal{O}_P(\rho^2 + \delta_K^2). \quad (A.6)$$

We start by showing that

$$\|\Delta\|_1^2 = \mathcal{O}_P([(K+T) \log(JTM_T)]/(JT) + \delta_K^2). \quad (A.7)$$

For this aim we apply Theorem 10.11 in Van de Geer (2000) that treats rates of convergence for least squares estimators on sieves. In our case we have the following sieve: $\mathcal{G}_T^* = \{g: \{1, \dots, T\} \times [0, 1]^d \rightarrow \mathbb{R}, g(t, x) = (1, z_t^\top) \mathcal{A} \psi(x) \text{ for an } (L+1) \times K \text{ matrix } \mathcal{A} \text{ and } z_t \in \mathbb{R}^L \text{ with the following properties: } |(1, z_t^\top) \mathcal{A} \psi(x)| \leq M_T \text{ for } 1 \leq t \leq T \text{ and } x \in [0, 1]^d\}$. With a constant C the δ -entropy $H_T(\delta, \mathcal{G}_T^*)$ of \mathcal{G}_T^* with respect to the empirical norm $\|g\|_1$ is bounded by

$$H_T(\delta, \mathcal{G}_T^*) \leq CT \log(M_T/\delta) + CK \log(KM_T/\delta). \quad (A.8)$$

For the proof of (A.8) note first that each element $g(t, x) = (1, z_t^\top) \mathcal{A}\psi(x)$ of \mathcal{G}_T^* can be chosen such that $T^{-1} \sum_{t=1}^T z_t z_t^\top$ is equal to the $L \times L$ identity matrix I_L . Then the bound $|(1, z_t^\top) \mathcal{A}\psi(x)| \leq M_T$ implies that $\|\mathcal{A}\psi(x)\| \leq M_T$. For the proof of (A.8) we use that the (δ/M_T) -entropy of a unit ball in \mathbb{R}^T is of order $\mathcal{O}(T \log(M_T/\delta))$ and that the δ -entropy with respect to the sup-norm for functions $\mathcal{A}\psi(x)$ with $\|\mathcal{A}\psi(x)\| \leq M_T$ is of order $\mathcal{O}(K \log(KM_T/\delta))$. In the last entropy bound we used that for each x it holds that $\|\psi(x)\| \leq K^{1/2}$. These two entropy bounds imply (A.8). Application of Theorem 10.11 in Van de Geer (2000) gives (A.7).

We now show that (A.7) implies (A.6). For this aim note first that by Bernstein's inequality for $a, d > 0, g \in \mathcal{G}_T^*$ with $\|g\|_2^2 \leq d$

$$P(|\|g\|_1^2 - \|g\|_2^2| \geq a) \leq 2 \exp\left(-\frac{a^2 JT}{2(a+d)M_T^2}\right). \quad (\text{A.9})$$

Furthermore, for $g, h \in \mathcal{G}_T^*$ it holds with constants C, C' that

$$\begin{aligned} \|g\|_1^2 - \|h\|_1^2 &\leq CK \left(T^{-1} \sum_{t=1}^T \|e_t - f_t\|^2 \right)^{1/2} \\ \left(T^{-1} \sum_{t=1}^T \|e_t + f_t\|^2 \right)^{1/2} &\leq C'K \|g - h\|_2 (\|g\|_2 + \|h\|_2), \end{aligned} \quad (\text{A.10})$$

where e_t and f_t are chosen such that $g(x, t) = e_t^\top \psi(x)$ and $h(x, t) = f_t^\top \psi(x)$. From (A.9) and (A.10) we get with a constant $C > 0$ that for $d = 1, 2, \dots$

$$\begin{aligned} P\left(\sup_{g \in \mathcal{G}_T^*, d\rho^2 \leq \|g\|_2^2 \leq (d+1)\rho^2} |\|g\|_1^2 - \|g\|_2^2| \geq d\rho^2/2\right) \\ \leq C \exp\left((C + K + T) \log(dKM_T) - d\rho^2 JT/[20M_T^2]\right). \end{aligned}$$

By summing these inequalities over $d \geq 1$ we get $\|\Delta\|_2^2 \leq \rho^2$ or

$$\|\Delta\|_2^2 \leq \|\|\Delta\|_1^2 - \|\Delta\|_2^2\| + \|\Delta\|_1^2 \leq \|\Delta\|_2^2/2 + \|\Delta\|_1^2$$

with probability tending to one. This shows Equation (A.6) and concludes the proof of Theorem 2.

A.3 Proof of Theorem 3

We will prove the first equation of the theorem for $h \neq 0$. The second equation follows from the first equation. We first prove that the matrix $T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \widehat{\mathcal{Z}}_{c,t}^\top$ is invertible, where $\mathcal{Z}_{c,t}^\top = (1, \mathcal{Z}_{c,t}^\top)$, $\widehat{\mathcal{Z}}_{c,t}^\top = (1, \widehat{\mathcal{Z}}_{c,t}^\top)$, and $\widehat{\mathcal{Z}}_{c,t} = \widehat{\mathcal{Z}}_t - T^{-1} \sum_{s=1}^T \widehat{\mathcal{Z}}_s$. This implies that $T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \widehat{\mathcal{Z}}_{c,t}^\top$ is invertible. Suppose that the assertion is not true. We can choose a random vector e such that $\|e\| = 1$ and $e^\top \sum_{t=1}^T \mathcal{Z}_{c,t} \widehat{\mathcal{Z}}_{c,t}^\top = 0$. Let \widehat{A} and A^* be the $L \times K$ matrices that are obtained by deleting the first rows of \widehat{A} and A^* , respectively. Let \widehat{A}_c and A_c^* be the matrices obtained from \widehat{A} and A^* by replacing their first rows by $\widehat{\alpha}_0^\top + (T^{-1} \sum_{t=1}^T \widehat{\mathcal{Z}}_t)^\top \widehat{A}$ and $\alpha_0^{\ast\top} + (T^{-1} \sum_{t=1}^T \mathcal{Z}_t)^\top A^*$, respectively. By definition, it follows that

$$\widehat{\mathcal{Z}}_t^\top \widehat{A} = \widehat{\mathcal{Z}}_{c,t}^\top \widehat{A}_c, \quad \mathcal{Z}_t^\top A^* = \mathcal{Z}_{c,t}^\top A_c^*. \quad (\text{A.11})$$

Note that

$$\begin{aligned} &\left\| T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \widehat{\mathcal{Z}}_{c,t}^\top \widehat{A}_c - T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \mathcal{Z}_{c,t}^\top A_c^* \right\| \\ &\leq T^{-1} \sum_{t=1}^T \left\| \mathcal{Z}_{c,t} \left(\widehat{\mathcal{Z}}_{c,t}^\top \widehat{A}_c - \mathcal{Z}_{c,t}^\top A_c^* \right) \right\| \\ &\leq \left(T^{-1} \sum_{t=1}^T \left\| \mathcal{Z}_{c,t} \right\|^2 \right)^{1/2} \left(T^{-1} \sum_{t=1}^T \left\| \widehat{\mathcal{Z}}_t^\top \widehat{A} - \mathcal{Z}_t^\top A^* \right\|^2 \right)^{1/2} \\ &= \mathcal{O}_P(\rho + \delta_K), \end{aligned} \quad (\text{A.12})$$

because of Assumption (A6) and Theorem 2. Thus with $f = T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \mathcal{Z}_{c,t}^\top e$, we obtain

$$\begin{aligned} \|f^\top m\| &= \|f^\top (\mathcal{A}_c^* \psi)\| + \mathcal{O}_P(T^{-1/2} + \delta_K) \\ &= \left\| e^\top T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \widehat{\mathcal{Z}}_{c,t}^\top \widehat{A}_c \psi \right\| + \mathcal{O}_P(T^{-1/2} + \rho + \delta_K) \\ &= \mathcal{O}_P(T^{-1/2} + \rho + \delta_K). \end{aligned}$$

This implies that m_0, \dots, m_d are linearly dependent, contradicting to Assumption (A10).

Let \widetilde{B} be the matrix given at (8) with B defined as in (9). Define $\widetilde{\mathcal{Z}}_{c,t} = \widetilde{B}^\top \widehat{\mathcal{Z}}_{c,t}$ and $\widetilde{A}_c = \widetilde{B}^{-1} \widehat{A}_c$. Then $\widetilde{\mathcal{Z}}_{c,t}^\top \widetilde{A}_c = \widehat{\mathcal{Z}}_{c,t}^\top \widehat{A}_c$ and $T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \widetilde{\mathcal{Z}}_{c,t}^\top = T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \mathcal{Z}_{c,t}^\top$. This gives with (A.12)

$$\begin{aligned} \|\widetilde{A}_c - A_c^*\| &= \left\| T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \mathcal{Z}_{c,t}^\top (\widetilde{A}_c - A_c^*) \right\| \mathcal{O}_P(1) \\ &= \left\| T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \widetilde{\mathcal{Z}}_{c,t}^\top \widetilde{A}_c - T^{-1} \sum_{t=1}^T \mathcal{Z}_{c,t} \mathcal{Z}_{c,t}^\top A_c^* \right\| \mathcal{O}_P(1) \\ &= \mathcal{O}_P(\rho + \delta_K). \end{aligned} \quad (\text{A.13})$$

Because of Theorem 2 this implies

$$\|\widetilde{A} - A^*\| = \mathcal{O}_P(\rho + \delta_K). \quad (\text{A.14})$$

Define $\widetilde{\mathcal{Z}}_{c,t}$ by $\widetilde{\mathcal{Z}}_{c,t}^\top = (1, \widetilde{\mathcal{Z}}_{c,t}^\top)$. Note that $\widetilde{\mathcal{Z}}_{c,t} = B^\top \widehat{\mathcal{Z}}_{c,t}$. Also, define $\widetilde{A} = B^{-1} \widehat{A}$, which equals \widehat{A}_c without the first row. From (A10), (A5), (A.14), and Theorem 2, we get

$$\begin{aligned} T^{-1} \sum_{t=1}^T \|\widetilde{\mathcal{Z}}_t - \mathcal{Z}_t\|^2 &= T^{-1} \sum_{t=1}^T \|\widetilde{\mathcal{Z}}_t - \mathcal{Z}_t\|^2 \\ &= T^{-1} \sum_{t=1}^T \left\| \widetilde{\mathcal{Z}}_t^\top (m_0, \dots, m_L)^\top - \mathcal{Z}_t^\top (m_0, \dots, m_L)^\top \right\|^2 \mathcal{O}_P(1) \\ &= T^{-1} \sum_{t=1}^T \left\| \widetilde{\mathcal{Z}}_t^\top A^* - \mathcal{Z}_t^\top A^* \right\|^2 \mathcal{O}_P(1) + T^{-1} \sum_{t=1}^T \left\| \widetilde{\mathcal{Z}}_t^\top \widetilde{A} - \mathcal{Z}_t^\top A^* \right\|^2 \\ &\quad \times \mathcal{O}_P(1) + \mathcal{O}_P(\delta_K^2) \\ &\leq T^{-1} \sum_{t=1}^T \|\widetilde{\mathcal{Z}}_t - \mathcal{Z}_t\|^2 \|\widetilde{A} - A^*\|^2 \mathcal{O}_P(1) + T^{-1} \sum_{t=1}^T \|\mathcal{Z}_t\|^2 \\ &\quad \times \|\widetilde{A} - A^*\|^2 \mathcal{O}_P(1) \\ &+ T^{-1} \sum_{t=1}^T \|\widetilde{\mathcal{Z}}_t^\top \widetilde{A} - \mathcal{Z}_t^\top A^*\|^2 \mathcal{O}_P(1) + \mathcal{O}_P(\rho^2 + \delta_K^2) \\ &= \mathcal{O}_P(\rho^2 + \delta_K^2). \end{aligned} \quad (\text{A.15})$$

From Equation (A.15) one gets

$$T^{-1} \sum_{t=1}^T \|\tilde{Z}_{c,t} - Z_{c,t}\|^2 = \mathcal{O}_P(\rho^2 + \delta_K^2). \quad (\text{A.16})$$

We will show that for $h \neq 0$

$$T^{-1} \sum_{t=h+1}^T \{(\tilde{Z}_{c,t+h} - Z_{c,t+h}) - (\tilde{Z}_{c,t} - Z_{c,t})\} Z_{c,t}^\top = \mathcal{O}_P(T^{-1/2}). \quad (\text{A.17})$$

This implies the first statement of Theorem 3, because by (A.16)

$$\begin{aligned} T^{-1} \sum_{t=-h+1}^T (\tilde{Z}_{c,t} - Z_{c,t})(\tilde{Z}_{c,t+h}^\top - Z_{c,t+h}^\top) &= \mathcal{O}_P(\rho^2 + \delta_K^2) \\ &= \mathcal{O}_P(T^{-1/2}). \end{aligned}$$

For the proof of (A.17), let $\tilde{\alpha}_c$ be the stack form of $\tilde{\mathcal{A}}_c$ and $\tilde{\alpha}_{c,0}^\top$ be its first row. Using the representation (6) and the first identity of (A.11), it can be verified that

$$\tilde{Z}_{c,t} = \tilde{S}_{t,Z}^{-1} J^{-1} \sum_{j=1}^J \{Y_{t,j} \tilde{A} \psi(X_{t,j}) - \tilde{A} \psi(X_{t,j}) \psi(X_{t,j})^\top \tilde{\alpha}_{c,0}\}, \quad (\text{A.18})$$

$$\tilde{\alpha}_c = \tilde{S}_\alpha^{-1} T^{-1} J^{-1} \sum_{t=1}^T \sum_{j=1}^J \{\psi(X_{t,j}) \otimes \tilde{Z}_{c,t}\} Y_{t,j}, \quad (\text{A.19})$$

where $\tilde{S}_{t,Z} = J^{-1} \sum_{j=1}^J \tilde{A} \psi(X_{t,j}) \psi(X_{t,j})^\top \tilde{A}^\top$ and $\tilde{S}_\alpha = T^{-1} J^{-1} \sum_{t=1}^T \sum_{j=1}^J \{\psi(X_{t,j}) \otimes \tilde{Z}_{c,t}\} \{\psi(X_{t,j}) \otimes \tilde{Z}_{c,t}\}^\top$. Define $\tilde{S}_{t,Z}$ as $\tilde{S}_{t,Z}$ with $\tilde{\mathcal{A}}_c$ replacing \tilde{A} . Also, define $\mathcal{S}_{t,Z} = \mathcal{A}_c^* E\{\psi(X_{t,j}) \psi(X_{t,j})^\top\} \mathcal{A}_c^{*\top}$ and

$$\mathcal{S}_\alpha = T^{-1} \sum_{t=1}^T E\{\{\psi(X_{t,j}) \otimes Z_{c,t}\} \{\psi(X_{t,j}) \otimes Z_{c,t}\}^\top | Z_t\}.$$

Let $\gamma = T^{-1/2}(\rho + \delta_K)^{-1}$. We argue that

$$\sup_{1 \leq t \leq T} \|\tilde{S}_{t,Z} - \mathcal{S}_{t,Z}\| = \mathcal{O}_P(\gamma), \quad \|\tilde{S}_\alpha - \mathcal{S}_\alpha\| = \mathcal{O}_P(\gamma). \quad (\text{A.20})$$

We show the first part of (A.20). The second part can be shown similarly. To prove the first part it suffices to show that, uniformly for $1 \leq t \leq T$,

$$\begin{aligned} J^{-1} \sum_{j=1}^J \mathcal{A}_c^* [\psi(X_{t,j}) \psi(X_{t,j})^\top - E\{\psi(X_{t,j}) \psi(X_{t,j})^\top\}] (\tilde{\mathcal{A}}_c - \mathcal{A}_c^*)^\top \\ = \mathcal{O}_P(\gamma), \end{aligned} \quad (\text{A.21})$$

$$\begin{aligned} J^{-1} \sum_{j=1}^J (\tilde{\mathcal{A}}_c - \mathcal{A}_c^*) [\psi(X_{t,j}) \psi(X_{t,j})^\top - E\{\psi(X_{t,j}) \psi(X_{t,j})^\top\}] \\ (\tilde{\mathcal{A}}_c - \mathcal{A}_c^*)^\top = \mathcal{O}_P(\gamma), \end{aligned} \quad (\text{A.22})$$

$$\begin{aligned} J^{-1} \sum_{j=1}^J \mathcal{A}_c^* [\psi(X_{t,j}) \psi(X_{t,j})^\top - E\{\psi(X_{t,j}) \psi(X_{t,j})^\top\}] \mathcal{A}_c^{*\top} \\ = \mathcal{O}_P(\gamma), \end{aligned} \quad (\text{A.23})$$

$$\begin{aligned} J^{-1} \sum_{j=1}^J \mathcal{A}_c^* E\{\psi(X_{t,j}) \psi(X_{t,j})^\top\} (\tilde{\mathcal{A}}_c - \mathcal{A}_c^*)^\top = \mathcal{O}_P(\gamma), \quad (\text{A.24}) \\ J^{-1} \sum_{j=1}^J (\tilde{\mathcal{A}}_c - \mathcal{A}_c^*) E\{\psi(X_{t,j}) \psi(X_{t,j})^\top\} (\tilde{\mathcal{A}}_c - \mathcal{A}_c^*)^\top = \mathcal{O}_P(\gamma). \end{aligned} \quad (\text{A.25})$$

The proof of (A.23)–(A.25) follows by simple arguments. We now show (A.21). Claim (A.22) can be shown similarly. For the proof of (A.21) we use Bernstein's inequality for the following sum:

$$P\left(\left|\sum_{j=1}^J W_j\right| > x\right) \leq 2 \exp\left(-\frac{1}{2} \frac{x^2}{V + Mx/3}\right). \quad (\text{A.26})$$

Here for a value of t with $1 \leq t \leq T$, the random variable W_j is an element of the $(L+1) \times 1$ -matrix $S = J^{-1} \mathcal{A}_c^* [\psi(X_{t,j}) \psi(X_{t,j})^\top e - E\{\psi(X_{t,j}) \psi(X_{t,j})^\top e\}]$ where $e \in \mathbb{R}^K$ with $\|e\| = 1$. In (A.26), V is an upper bound for the variance of $\sum_{j=1}^J W_j$ and M is a bound for the absolute values of W_j (i.e. $|W_j| \leq M$ for $1 \leq j \leq J$, a.s.). With some constants C_1 and C_2 that do not depend on t and the row number we get $V \leq C_1 J^{-1}$ and $M \leq C_2 K^{1/2} J^{-1}$. Application of Bernstein's inequality gives that, uniformly for $1 \leq t \leq T$ and $e \in \mathbb{R}^K$ with $\|e\| = 1$, all $(L+1)$ elements of S are of order $\mathcal{O}_P(\gamma)$. This shows claim (A.21).

From (A.13), (A.15), (A.18), (A.19), and (A.20) it follows that uniformly for $1 \leq t \leq T$,

$$\begin{aligned} \tilde{Z}_{c,t} - Z_{c,t} &= S_{t,Z}^{-1} J^{-1} \sum_{j=1}^J \varepsilon_{t,j} \mathcal{A}_c^* \psi(X_{t,j}) + S_{t,Z}^{-1} J^{-1} \sum_{j=1}^J \varepsilon_{t,j} \\ &\quad \times (\tilde{\mathcal{A}}_c - \mathcal{A}_c^*) \psi(X_{t,j}) \\ &+ S_{t,Z}^{-1} J^{-1} \sum_{j=1}^J (\tilde{\mathcal{A}}_c - \mathcal{A}_c^*) \psi(X_{t,j}) \psi(X_{t,j})^\top \mathcal{A}_c^{*\top} Z_{c,t} + \mathcal{O}_P(T^{-1/2}) \\ &\equiv \Delta_{t,1,Z} + \Delta_{t,2,Z} + \Delta_{t,3,Z} + \mathcal{O}_P(T^{-1/2}). \end{aligned} \quad (\text{A.27})$$

For the proof of the theorem it remains to show that for $1 \leq j \leq 3$

$$T^{-1} \sum_{t=-h+1}^T (\Delta_{t+h,j,Z} - \Delta_{t,j,Z}) Z_{c,t}^\top = \mathcal{O}_P(T^{-1/2}). \quad (\text{A.28})$$

This can be easily checked for $j = 1$. For $j = 2$ it follows from $\|\tilde{\mathcal{A}}_c - \mathcal{A}_c^*\| = \mathcal{O}(\rho + \delta_k)$ and

$$E\left\{\left\|T^{-1} J^{-1} \sum_{t=1}^T \sum_{j=1}^J \varepsilon_{t,j} S_{t,Z}^{-1} \mathcal{M} \psi(X_{t,j})\right\|^2\right\} = \mathcal{O}(KJ^{-1}T^{-1})$$

for any $L \times K$ matrix \mathcal{M} with $\|\mathcal{M}\| = 1$. For the proof of (A.28) for $j = 3$, it suffices to show that

$$T^{-1} \sum_{t=1}^{T+h} \Delta_{t,j,Z} (Z_{c,t-h} - Z_{c,t})^\top = \mathcal{O}_P(T^{-1/2}). \quad (\text{A.29})$$

We note first that for $1 \leq l \leq L$

$$\begin{aligned} T^{-1} \sum_{t=1}^{T+h} \Delta_{t,3,Z} (Z_{c,t-h,l} - Z_{c,t,l}) \\ = T^{-1} J^{-1} \sum_{t=1}^{T+h} \sum_{j=1}^J \left\{ \left(V_{h,t}^\top \mathcal{A}_c^* \psi(X_{t,j}) \psi(X_{t,j})^\top \right) \otimes S_{t,Z}^{-1} \right\} (\tilde{\alpha} - \alpha^*), \end{aligned}$$

where $V_{h,t} = (Z_{c,t-h,t} - Z_{c,t,t})Z_{c,t,t}$, and $\tilde{\alpha}$ and α^* denote the stack forms of A and A^* , respectively. For the proof of (A.29) it suffices to show

$$T^{-1}J^{-1} \sum_{t=1}^{T+h} \sum_{j=1}^J \{ (E[V_{h,t}]^\top \mathcal{A}_c^* \psi(X_{t,j}) \psi(X_{t,j})^\top) \otimes S_{t,Z}^{-1} \} \times (\tilde{\alpha} - \alpha^*) = o_P(T^{-1/2}), \tag{A.30}$$

$$\left\| T^{-1}J^{-1} \sum_{t=1}^{T+h} \sum_{j=1}^J \{ (\{V_{h,t} - E[V_{h,t}]\}^\top \mathcal{A}_c^* \psi(X_{t,j}) \psi(X_{t,j})^\top) \otimes S_{t,Z}^{-1} \} \right\|^2 = \mathcal{O}_P(KJ^{-1}T^{-1}). \tag{A.31}$$

Claim (A.31) can be easily shown by calculating the expectation of the left hand side of (A.31) and by using the mixing condition at Assumption (A9). For a proof of (A.30) we remark first that by construction

$$0 = T^{-1} \sum_{t=1}^T (\tilde{Z}_{c,t} - Z_{c,t})Z_{c,t}^\top.$$

Using (A.27) and similar arguments as in the proof of (A.28) for $j = 1, 2$ we get that

$$T^{-1} \sum_{t=1}^T \Delta_{t,3,Z} Z_{c,t}^\top = T^{-1}J^{-1} \sum_{t=1}^T \sum_{j=1}^J \{ (Z_{c,t} Z_{c,t}^\top \mathcal{A}_c^* \psi(X_{t,j}) \psi(X_{t,j})^\top) \otimes S_{t,Z}^{-1} \} (\tilde{\alpha} - \alpha^*) = o_P(T^{-1/2}).$$

As in the proof of (A.31) one can show that

$$\left\| T^{-1}J^{-1} \sum_{t=1}^{T+h} \sum_{j=1}^J \{ (\{Z_{c,t} Z_{c,t}^\top - E[Z_{c,t} Z_{c,t}^\top]\} \mathcal{A}_c^* \psi(X_{t,j}) \psi(X_{t,j})^\top) \otimes S_{t,Z}^{-1} \} \right\|^2 = \mathcal{O}_P(KJ^{-1}T^{-1}).$$

The last two equalities imply that

$$T^{-1}J^{-1} \sum_{t=1}^T \sum_{j=1}^J \{ (E[Z_{c,t} Z_{c,t}^\top] \mathcal{A}_c^* \psi(X_{t,j}) \psi(X_{t,j})^\top) \otimes S_{t,Z}^{-1} \} \times (\tilde{\alpha} - \alpha^*) = o_P(T^{-1/2}).$$

Because of Assumption (A9) this implies claim (A.29) and concludes the proof of Theorem 3.

[Received June 2007. Revised August 2008.]

REFERENCES

Biswal, B., Yetkin, F., Haughton, V., and Hyde, J. (1995), "Functional Connectivity in the Motor Cortex of Resting Human Brain Using Echo-Planar MRI," *Magnetic Resonance in Medicine*, 34, 537–541.
 Black, F., and Scholes, M. (1973), "The Pricing of Options and Corporate Liabilities," *The Journal of Political Economy*, 81, 637–654.

Brüggemann, R., Lütkepohl, H., and Saikkonen, P. (2006), "Residual Autocorrelation Testing for Vector Error Correction Models," *Journal of Econometrics*, 134, 579–604.
 Brumback, B., and Rice, J. A. (1998), "Smoothing Spline Models for the Analysis of Nested and Crossed Samples of Curves," *Journal of the American Statistical Association*, 93, 961–994.
 Connor, G., Hagmann, M., and Linton, O. (2007). *Efficient Semiparametric Estimation of the Fama-French Model and Extensions*, Preprint.
 Connor, G., and Linton, O. (2007), "Semiparametric Estimation of a Characteristic-based Factor Model of Stock Returns," *Journal of Empirical Finance*, 14, 694–717.
 Cont, R., and da Fonseca, J. (2002), "The Dynamics of Implied Volatility Surfaces," *Quantitative Finance*, 2, 45–60.
 de Boor, C. (2001). *A Practical Guide to Splines*, Berlin, Heidelberg: Springer-Verlag.
 Diebold, F. X., and Li, C. (2006), "Forecasting the Term Structure of Government Bond Yields," *Journal of Econometrics*, 130, 337–364.
 Fama, E. F., and French, K. R. (1992), "The Cross-Section of Expected Stock Returns," *The Journal of Finance*, 47, 427–465.
 Fan, J., Yao, Q., and Cai, Z. (2003), "Adaptive Varying-Coefficient Linear Models," *Journal of the Royal Statistical Society: Series B*, 65, 57–80.
 Fengler, M. R., Härdle, W., and Mammen, E. (2007), "A Semiparametric Factor Model for Implied Volatility Surface Dynamics," *Journal of Financial Econometrics*, 5, 189–218.
 Forni, M., Hallin, M., Lippi, M., and Reichlin, L. (2000), "The Generalized Dynamic Factor Model: Identification and Estimation," *The Review of Economics and Statistics*, 82, 540–554.
 Forni, M., and Lippi, M. (2001), "The Generalized Factor Model: Representation Theory," *Econometric Theory*, 17, 1113–1141.
 Gasser, T., Möcks, R., and Verleger, R. (1983), "Selavco: A Method to Deal With Trial-to-Trial Variability of Evoked Potential," *Electroencephalography and Clinical Neurophysiology*, 55, 717–723.
 Hafner, R. (2004). *Stochastic Implied Volatility*, Berlin: Springer.
 Hallin, M., and Liska, R. (2007), "Determining the Number of Factors in the Generalized Dynamic Factor Model," *Journal of the American Statistical Association*, 102, 603–617.
 Hansen, L.H., Nielsen, B., and Nielsen, J.P. (2004). "Two Sided Analysis of the Variance With a Latent Time Series," Nuffield College Economic Working Paper 2004-W25, University of Oxford.
 Hosking, J. R. M. (1980), "The Multivariate Portmanteau Statistic," *Journal of the American Statistical Association*, 75, 602–608.
 Hosking, J. R. M. (1981), "Lagrange-Multiplier Tests of Multivariate Time-Series Models," *Journal of the Royal Statistical Society, Series B*, 43, 219–230.
 Kantorovich, L. V., and Akilov, G. P. (1982). *Functional Analysis* (2nd ed.), Oxford, U.K.: Pergamon Press.
 Kauermann, G. (2000), "Modeling Longitudinal Data With Ordinal Response by Varying Coefficients," *Biometrics*, 56, 1692–1698.
 Lee, R. D., and Carter, L. (1992), "Modeling and Forecasting the Time Series of U.S. Mortality," *Journal of the American Statistical Association*, 87, 659–671.
 Logothetis, N., and Wandell, B. (2004), "Interpreting the Bold Signal," *Annual Review of Physiology*, 66, 735–769.
 Lütkepohl, H. (1993). *Introduction to Multiple Time Series Analysis*, Berlin, Heidelberg: Springer-Verlag.
 Martinussen, T., and Scheike, T. (2000), "A Nonparametric Dynamic Additive Regression Model for Longitudinal Data," *Annals of Statistics*, 28, 1000–1025.
 Nelson, C. R., and Siegel, A. F. (1987), "Parsimonious Modeling of Yield Curves," *Journal of Business*, 60, 473–489.
 Peña, D., and Box, E. P. (1987), "Identifying a Simplifying Structure in Time Series," *Journal of the American Statistical Association*, 82, 836–843.
 Stock, J. H., and Watson, M.W. (2005). "Implications of Dynamic Factor Models for VAR Analysis," NBER Working Papers 11467, National Bureau of Economic Research, Inc., available at <http://ideas.repec.org/p/nbr/nberwo/11467.html>.
 Van de Geer, S. (2000). *Empirical Processes in M-Estimation*, Cambridge, U.K.: Cambridge University Press.
 Vrahatis, M. N. (1989), "A Short Proof and a Generalization of Miranda's Existence Theorem," *Proceedings of the American Mathematical Society*, 107, 701–703.
 Worsley, K., Liao, C., Aston, J., Petre, V., Duncan, G., Morales, F., and Evans, A. (2002), "A General Statistical Analysis for fMRI Data," *NeuroImage*, 15, 1–15.
 Yang, L., Park, B. U., Xue, L., and Härdle, W. (2006), "Estimation and Testing for Varying Coefficients in Additive Models With Marginal Integration," *Journal of the American Statistical Association*, 101, 1212–1227.

A Generalized ARFIMA Process with Markov-Switching Fractional Differencing Parameter

Wen-Jen Tsay¹ and Wolfgang Karl Härdle^{2*}

April 24, 2007

¹ The Institute of Economics, Academia Sinica, Taiwan

² CASE – Center for Applied Statistics and Economics

Humboldt-Universität zu Berlin,

Spandauer Straße 1, 10178 Berlin, Germany

Abstract

We propose a general class of Markov-switching-ARFIMA processes in order to combine strands of long memory and Markov-switching literature. Although the coverage of this class of models is broad, we show that these models can be easily estimated with the DLV algorithm proposed. This algorithm combines the Durbin-Levinson and Viterbi procedures. A Monte Carlo experiment reveals that the finite sample performance of the proposed algorithm for a simple mixture model of Markov-switching mean and ARFIMA(1, d , 1) process is satisfactory. We apply the Markov-switching-ARFIMA models to the U.S. real interest rates, the Nile river level, and the U.S. unemployment rates, respectively. The results are all highly consistent with the conjectures made or empirical results found in the literature. Particularly, we confirm the conjecture in Beran and Terrin (1996) that the observations 1 to about 100 of the Nile river data seem to be more independent than the subsequent observations, and the value of differencing parameter is lower for the first 100 observations than for the subsequent data.

Key words: Markov chain; ARFIMA process; Viterbi algorithm; Long memory

JEL classification: C14, C22, C32, C52, C53, G12

*This research was supported by the Deutsche Forschungsgemeinschaft through the SFB 649 'Economic Risk'. Correspondence to Wen-Jen Tsay. The Institute of Economics, Academia Sinica, Taipei, Taiwan, R.O.C. Tel: (886-2) 2782-2791 ext. 296. Fax: (886-2) 2785-3946. E-Mail: wtsay@ieas.econ.sinica.edu.tw

1 Introduction

It is well known that many time series data exhibit long memory, or long-range dependence, including the Nile river level, *ex post* real interest rate, forward premium, and the dynamics of aggregate partisanship and macroideology. Among the many other examples that Beran (1994) gives the Nile river data has been known for its long memory behavior since ancient times, and this is one of the time series that led to the discovery of the Hurst effect (Hurst, 1951) and motivated Mandelbrot and his co-workers (Mandelbrot and van Ness, 1968; Mandelbrot and Wallis, 1969) to introduce fractional Gaussian noise to model long memory phenomenon.

Long range dependence also has been observed in financial data. As demonstrated by Ding et al. (1993), de Lima and Crato (1993) and Bollerslev and Mikkelsen (1996) that the volatility of most financial time series exhibits strong persistency and can be well described as a long memory process. Evidence of financial market volatility's strong persistency inspired Breidt et al. (1998) to propose a class of long memory stochastic volatility (LMSV) models. Deo et al. (2006) also show that the LMSV model is useful for forecasting realized volatility (RV) which is an important quantity in finance.

Figure 1 displays the yearly Nile river minima based on measurements at the Roda gauge near Cairo during the years 622-1284. Beran (1994, p.33) documents that "When one only looks at short time periods, then there seem to be cycles or local trend. However, looking at the whole series, there is no apparent persisting cycle." The changing pattern of the Nile river data leads Bhattacharya et al. (1983) to argue that the so-called Hurst effect can also be explained as if the observations are composed as the sum of a weakly dependent stationary process and a deterministic function. As a consequence it is important to distinguish between a long memory time series and a weakly dependent time series with change-points in the mean. This question has been intensively considered in the literature, including Künsch (1986) and Heyde and Dai (1996). Berkes et al. (2006) presents an overview about this strand of literature. Similarly, Diebold and Inoue (2001) shows that long memory also may be easily confused with a Markov-switching mean. Thus, most of the existing literature considers long memory as a competing modeling framework against the structural change and Markov-switching models.

The Nile river level time series is far more complicated than a pure long memory or

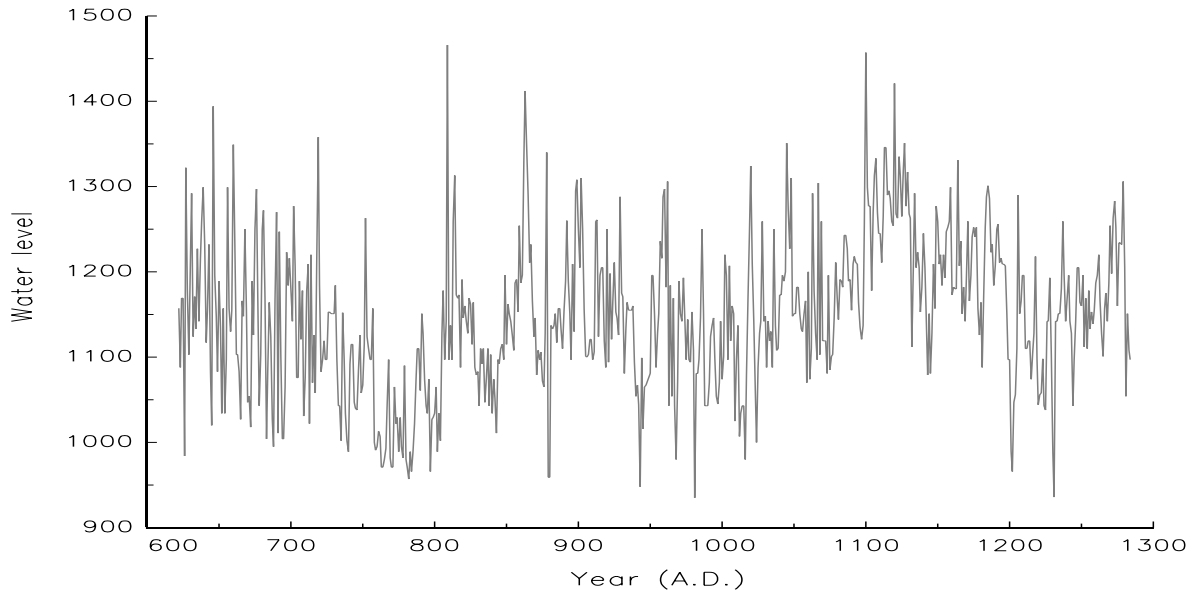


Figure 1: Yearly Nile river minima based on measurements at the Roda gauge near Cairo.

a weakly dependent time series with change-points in the mean to describe. Beran and Terrin (BT) (1996) suggest therefore that the Hurst parameter characterizing the yearly Nile river might change over time. When estimating the Nile river data with the autoregressive fractionally-integrated moving-average (ARFIMA) or $I(d)$ process introduced by Granger (1980), Granger and Joyeux (1980) and Hosking (1981), Beran and Terrin (1996, p.629) show that the data can be well fitted with an ARFIMA(0, d , 0) model with $d = 0.4$, where the fractional differencing parameter d of ARFIMA process acts like the Hurst parameter H of fractional Gaussian noise in characterizing the hyperbolic decay of the autocovariance function of a long memory process. BT further claim that the observations 1 to about 100 seem to be more independent than the subsequent observations, and the value of the fractional differencing parameter might be *lower* for the first 100 observations than for the subsequent data. If this claim is right, then there should be a structural change in the long range persistence of the Nile river data around the year 720, and the Nile river data neither can be described with a pure long memory nor a weakly dependent time series with change-points in the mean.

The possible change of the differencing parameter stimulate BT to propose a statistic for testing the stability of the fractional differencing parameter. This testing statistic has been further discussed and extended in Horváth and Shao (1999) and Horváth (2001). However,

their methods can not identify the change points of the fractional differencing parameter. A Bayesian random persistent-shift (RPS) method for detecting structural change in the differencing parameter and the process level has been considered in Ray and Tsay (2002). Nevertheless, the RPS method is not built on the Markov-switching framework, thus may not fully characterize the cycling behavior of the data series, i.e., “seven years of great abundance” and “seven years of famine” — the Joseph effect named by Mandelbrot and van Ness (1968) and Mandelbrot and Wallis (1969).

The above considerations lead us to combine the long memory and Markov-switching literature into a unified framework. We introduce a Markov-switching-ARFIMA (MS-ARFIMA) process by extending the hidden Markov model. Given that the hidden Markov model has become extremely popular in speech recognition as shown in Juang and Rabiner (1991) and Qian and Titterton (1991), and in econometrics, finance, genetics, and neurophysiology as outlined in Robert et al. (2000), the MS-ARFIMA model provides a flexible modeling framework for many applications to these fields. Moreover, the research conducted in this paper also solve the puzzle raised by Diebold and Inoue (2001) by estimating the differencing parameter allowing for the parameters of interest are Markov-switching.

The remaining parts of this paper are arranged as follows: Section 2 presents the MS-ARFIMA process and the algorithms for estimating the parameters of interest. In Section 3 we consider the finite sample performance of the proposed algorithm under the simple mixture of a Markov-switching mean and an ARFIMA(1, d , 1) process. We then apply the proposed methodology to the U.S. real interest rates, the Nile river data, and the U.S. unemployment rates in Section 4. Section 5 provides a conclusion.

2 Models and Main Results

The objective of this paper is to propose a general class of Markov-switching-ARFIMA processes in order to combine strands of long memory and Markov-switching literature. This class of models offers a rich dynamic mixture of a Markov chain and an $I(d)$ process.

Let $\{s_t\}_{t=1}^T$ be the latent sample path of an N -state Markov chain. At each time s_t can

assume only an integer value of $1, 2, \dots, N$, and its transition probability matrix is

$$\mathcal{P} \equiv \begin{bmatrix} p_{11} & p_{21} & \cdots & p_{N1} \\ p_{12} & p_{22} & \cdots & p_{N2} \\ \vdots & \vdots & \ddots & \vdots \\ p_{1N} & p_{2N} & \cdots & p_{NN} \end{bmatrix},$$

where $p_{ij} = P(s_t = j \mid s_{t-1} = i)$ and $\sum_{j=1}^N p_{ij} = 1$ for all i .

An $I(d)$ process, x_t , is defined as:

$$(1 - L)^d x_t = h_t,$$

where L is the lag operator ($Lk_t = k_{t-1}$) and h_t is a short memory process. When $d > 0$, the $I(d)$ process is often called the long memory process, because its autocovariance function is not summable so as to capture the long range dependence of a time series. In addition, the $I(d)$ process is nonstationary when $d \geq \frac{1}{2}$, otherwise, it is covariance stationary.

Combining the defining feature of a Markov chain and that of an $I(d)$ process, we propose the following MS-ARFIMA(p, d, q) process:

$$w_t = \mu_{s_t} I\{t \geq 1\} + (1 - L)^{-d_{s_t}} \sigma_{s_t} z_t I\{t \geq 1\} = \mu_{s_t} I\{t \geq 1\} + y_{s_t}, \quad (1)$$

where $I\{\cdot\}$ is the indicator function and z_t is stationary process with mean zero and bounded positive spectral density $f_u(\lambda) \sim G_0$ as $\lambda \rightarrow 0$ at each possible regime, thus including stationary and invertible ARMA process as its special case. The most distinguished feature of the process is that the fractional differencing parameter d_{s_t} well known in the long memory literature is allowed to be a Markov chain satisfying the following Assumption A:

Assumption A. s_t is independent of z_τ for all t and τ .

The model in (1) subsumes many interesting models in the literature. When $N = 1$, w_t reduces to the specification in (7) of Shimotsu and Phillips (2005):

$$w_t = \mu_0 + (1 - L)^{-d_0} \sigma_0 z_t I\{t \geq 1\} \quad (2)$$

which also can be represented as:

$$w_t = \mu_0 + \sum_{k=0}^{t-1} \frac{(d_0)_k}{k} \sigma_0 z_{t-k}, \quad (3)$$

where

$$(d_0)_k = \frac{\Gamma(d_0 + k)}{\Gamma(d_0)} = (d_0)(d_0 + 1) \dots (d_0 + k - 1) \quad (4)$$

is Pochhammer's symbol for the forward factorial and $\Gamma(\cdot)$ is the gamma function. Moreover, under the model in (1) and $d_{s_t} = 0$, w_t still includes the Markov-switching AR model considered in Hamilton (1989) as one of its special cases. We will show that the estimation of the model in (1) can be easily implemented with the algorithm proposed in this paper, even though the parameter estimation from a noisy version of realizations of Markov models is extremely difficult in all but very simple examples as well documented in Qian and Titterton (1991).

Let the total sample size be T , and denote $\mathcal{W}_t \equiv (w_1, w_2, \dots, w_t)^\top$ the column vector containing the observations from time 1 to time t , while $\mathcal{S}_t = (s_1, s_2, \dots, s_t)^\top$ represents the corresponding states, and $\mathcal{Y}_t = (y_1, y_2, \dots, y_t)^\top$ in (1) is similarly defined. The column vector $\alpha = (\mu_1, \dots, \mu_N, \sigma_1, \dots, \sigma_N, \phi_{11}, \dots, \phi_{1p}, \phi_{21}, \dots, \phi_{Np}, d_1, \dots, d_N, \theta_{11}, \dots, \theta_{Nq})^\top$ and p_{ij} (transition probabilities) consist of the parameters characterizing the conditional density function (cdf) of w_t . After stacking the parameter vector α and the transition probabilities p_{ij} into one column vector ξ , we can represent the cdf of w_t as $f(w_t | \mathcal{S}_t, \mathcal{W}_{t-1}; \xi)$, clearly showing that the cdf of w_t depends on the entire past routes of states (in general). Indeed, there are N^T possible paths of states running throughout the observations \mathcal{W}_T .

To illustrate the proposed algorithm for the model in (1), we first consider the simplest case where w_t in (1) is generated as:

$$w_t = \mu_{s_t} I\{t \geq 1\} + (1 - L)^{-d_0} \sigma_0 \varepsilon_t I\{t \geq 1\} = \mu_{s_t} I\{t \geq 1\} + y_t, \quad (5)$$

where $d < \frac{1}{2}$ and ε_t is a zero mean normally, independently and identically distributed white noise (i.i.d.) with $E(\varepsilon_t^2) = 1$. That is, w_t in (5) is a special type of MS-ARFIMA(0, d , 0) process whose differencing parameter is fixed across different regimes. Under Assumption A and $\varepsilon_t \sim N(0, 1)$ i.i.d. process, the likelihood function of \mathcal{W}_T , $L(\mathcal{S}_T, \mathcal{W}_T; \xi)$ hereafter, for the hidden Markov model in (5) equals

$$L(\mathcal{S}_T, \mathcal{W}_T; \xi) = (2\pi)^{-T/2} |\Lambda|^{-1/2} \exp\left(-\frac{1}{2} \mathcal{Y}_T^\top \Lambda^{-1} \mathcal{Y}_T\right) \prod_{t=1}^T \Pr(s_t | s_{t-1}), \quad (6)$$

where $\Lambda = E(\mathcal{Y}_T \mathcal{Y}_T^\top)$, and $\Pr(s_1 | s_0)$ is evaluated with the unconditional probability that the process will be in regime s_1 . Given that y_t in (5) is a simple ARFIMA(0, d , 0) process,

we can use the Durbin-Levinson algorithm to derive

$$(2\pi)^{-T/2} |\Lambda|^{-1/2} \exp\left(-\frac{1}{2} \mathcal{Y}_T^\top \Lambda^{-1} \mathcal{Y}_T\right) = \prod_{t=1}^T (2\pi)^{-1/2} v_{t-1}^{-1/2} \exp\left\{-\frac{(y_t - \hat{y}_t)^2}{2v_{t-1}}\right\}, \quad (7)$$

where \hat{y}_t denotes the one-step ahead predictor of y_t with the observation \mathcal{Y}_{t-1} as $j \geq 2$, and v_{t-1} is the corresponding one-step ahead prediction variance. Deriche and Tewfik (1993) also have employed the Durbin-Levinson algorithm to estimate a univariate ARFIMA(0, d , 0) processes without Markov-switching characteristic. Note that as $t = 1$, $\hat{y}_1 = 0$, and $v_0 = \gamma_0$ corresponds to the variance of y_t . As a result, the likelihood function in (6) can be rewritten as:

$$L(\mathcal{S}_T, \mathcal{W}_T; \xi) = \prod_{t=1}^T (2\pi)^{-1/2} v_{t-1}^{-1/2} \exp\left\{-\frac{(y_t - \hat{y}_t)^2}{2v_{t-1}}\right\} \Pr(s_t | s_{t-1}), \quad (8)$$

indicating that the *unconditional* likelihood function of the mixture model in (5) can be exactly and recursively evaluated provided that we can identify the true path of s_t , \mathcal{S}_T^* .

We do not know in reality the value of \mathcal{S}_T^* . However, the recursive structure shown in (8) is especially suitable for implementing the Viterbi (1967) algorithm in the digital communication literature to identify the most likely path of states among the N^T possible routes within \mathcal{W}_T . We thus combine the Durbin-Levinson algorithm and the Viterbi algorithm to suggest a *Durbin-Levinson-Viterbi* (DLV) algorithm for the model in (5). When compared to the original Viterbi algorithm designed for solving the problem of maximum a posteriori probability estimate of the state sequence of a finite-state discrete-time Markov process observed in white noise, the DLV algorithm proposed in this paper is concerned with the hidden Markov process observed in a much more general ARFIMA noise. Since the DLV algorithm can estimate the differencing parameter of a time series allowing for the presence of a Markov-switching mean, the puzzle raised by Diebold and Inoue (2001) that long memory can be easily confused with a Markov-switching mean is thus resolved by using this DLV algorithm.

To locate the most likely path running through the data \mathcal{W}_T with the idea of Viterbi (1967), we note first that, for each time t , there are N possible states ending at time t , i.e., $(s_t = i)$, $i = 1, \dots, N$. For a particular node of these N end points at time t , say $(s_t = j)$, there exists a corresponding most likely path:

$$(\mathcal{S}_{t-1}(s_t = j), s_t = j) = (s_1(s_t = j), s_2(s_t = j), \dots, s_{t-1}(s_t = j), s_t = j), \quad (9)$$

which ends at this particular node ($s_t = j$). We refer to the path $(\mathcal{S}_{t-1}(s_t = j), s_t = j)$ in (9) as the *survivor* associated with the node ($s_t = j$). Note that, with little loss of clarity, we do not explicitly specify that the path depends on the parameter ξ and the observations \mathcal{W}_t in order to simplify the notation. The likelihood function generated from this survivor $(\mathcal{S}_{t-1}(s_t = j), s_t = j)$ and the formula in (8) is recorded as $L(\mathcal{S}_{t-1}(s_t = j), s_t = j, \mathcal{W}_t; \xi)$ and is crucial for locating the most likely path running from time 1 to time T . In short, for each node ($s_t = j$) at time t , there exists a most likely path, survivor $(\mathcal{S}_{t-1}(s_t = j), s_t = j)$, and its associated likelihood function $L(\mathcal{S}_{t-1}(s_t = j), s_t = j, \mathcal{W}_t; \xi)$. Most importantly, the number of survivors at each time t is always equal to N .

Given the N survivors at time t and in order to locate the survivor $(\mathcal{S}_t(s_{t+1} = i), s_{t+1} = i)$ for a particular node ($s_{t+1} = i$) at time $t + 1$, among the N segments connecting the node ($s_{t+1} = i$) and the N time- t survivors $(\mathcal{S}_{t-1}(s_t = j), s_t = j)$ recorded at time t , we select the one producing the largest likelihood function $L(\mathcal{S}_t(s_{t+1} = i), s_{t+1} = i, \mathcal{W}_{t+1}; \xi)$ among these N possible candidates, and name it as the survivor $(\mathcal{S}_t(s_{t+1} = i), s_{t+1} = i)$ for this particular node ($s_{t+1} = i$). The computation of the aforementioned likelihoods is simple, because we record the likelihood functions of the N time- t survivors at each time t .

This recursive updating process proceeds from time 1 to time T and results in N time- T survivors $(\mathcal{S}_{T-1}(s_T = i), s_T = i)$ and their associated likelihood function $L(\mathcal{S}_{T-1}(s_T = i), s_T = i, \mathcal{W}_T; \xi)$, for each $i = 1, \dots, N$. From these N time- T survivors we select the one producing the largest likelihood function, say $L(\mathcal{S}_{T-1}(s_T = g), s_T = g, \mathcal{W}_T; \xi)$, as the most likely path running from time 1 to time T . Combining a numerical optimization procedure and this chosen likelihood function $L(\mathcal{S}_{T-1}(s_T = g), s_T = g, \mathcal{W}_T; \xi)$ generated from the Viterbi algorithm and the Durbin-Levinson algorithm displayed in (7), we can estimate the parameters ξ and identify the states \mathcal{S}_T hidden in the observations \mathcal{W}_T .

We now consider another special type of MS-ARFIMA(p, d, q) process:

$$w_t = \mu_{s_t} I\{t \geq 1\} + y_t = \mu_{s_t} I\{t \geq 1\} + (1 - L)^{-d_0} \sigma_0 z_t I\{t \geq 1\}, \quad \phi(L)z_t = \theta(L)\varepsilon_t, \quad (10)$$

where

$$\phi(L) = 1 - \phi_1 L - \dots - \phi_p L^p, \quad \theta(L) = 1 + \theta_1 L + \dots + \theta_q L^q, \quad (11)$$

and the roots of the polynomial $\phi(L)$ and those of $\theta(L)$ in (11) are all outside the unit circle and share no common roots. The model in (10) is much more general than that in (5), but

still can be estimated with the preceding Viterbi algorithm after some modifications. Please note that the value of fractional differencing parameter is unchanged across different regimes as that imposed in (5).

Note that the term y_t in (10) can be rearranged as

$$y_t = (1 - L)^{-d_0} \sigma_0 \phi(L)^{-1} \theta(L) \varepsilon_t, \quad t = 1, 2, \dots \quad (12)$$

We then have

$$\phi(L)y_t = (1 - L)^{-d_0} \sigma_0 \theta(L) \varepsilon_t = \sigma_0 \theta(L) (1 - L)^{-d_0} \varepsilon_t = \sigma_0 \theta(L) \tilde{y}_t, \quad t = 1, 2, \dots, \quad (13)$$

where $\tilde{y}_t = (1 - L)^{-d_0} \varepsilon_t$ is an ARFIMA(0, d , 0) process. Dueker and Serletis (2000) use the same transformation method for estimating an ARFIMA(p , d , q) process. Conditional on a set of $\phi(L)$ and $\theta(L)$ and a suitable starting value, the *conditional* likelihood function of y_t in (12) can still be evaluated exactly with the transformed ARFIMA(0, d , 0) \tilde{y}_t in (13) and the Durbin-Levinson algorithm defined in (7). For example, conditional on y_0 being equal to 0, we can extract an ARFIMA(0, d , 0) process from an ARFIMA(1, d , 1) process as follows:

$$\sigma_0 \tilde{y}_t = y_t - \phi_1 y_{t-1} - \sigma_0 \theta_1 \tilde{y}_{t-1}, \quad t = 1, \dots, T. \quad (14)$$

Conditional on a set of $\phi(L)$ and $\theta(L)$ and a suitable starting value for the parameter ξ , we can recursively and exactly evaluate the conditional likelihood function of the hidden Markov model using the DLV algorithm proposed previously.

The same idea also applies to the class of MS-ARFIMA(p , d , q) processes in (1) where d can be Markov-switching. However, we cannot use the Durbin-Levinson algorithm when the fractional differencing parameter is allowed to be Markov-switching. Nevertheless, the Viterbi algorithm is still powerful enough to locate the most likely path under this circumstance. That is, conditional on a suitable starting value for the parameter ξ , we employ the recursive structure inherent in Viterbi algorithms to identify the most likely path running through the data set.

3 Monte Carlo Experiment

In this section we consider a Monte Carlo experiment to demonstrate the finite sample performance of the proposed DLV algorithm on a special version of the model in (1):

$$w_t = \mu_{s_t} I\{t \geq 1\} + (1 - L)^{-d_0} \sigma_0 (1 - \phi_1 L)^{-1} (1 + \theta_1 L) \varepsilon_t I\{t \geq 1\}. \quad (15)$$

We employ three different values of the fractional differencing parameter:

$$d_0 = \{0.2, 0.3, 0.4\}, \tag{16}$$

along with the following parameters:

$$\mu_1 = 4, \mu_2 = 1, \phi_1 = 0.5, \theta_1 = 0.5, p_{11} = p_{22} = 0.95, \tag{17}$$

and σ_0 is chosen to ensure that the variance of the ARFIMA(1, d , 1) noise in (15) is equal to 1 across different configurations. Note that the positive values of d_0 in (16) are chosen to reflect the variations used in the long memory literature.

All the computations are performed with GAUSS. Two hundred replications are conducted for each specification at 3 different sample sizes ($T = 100, 200, 400$) usually encountered in the empirical applications. For each sample size T , 200 additional values are generated in order to obtain random starting values. The optimization algorithm used to implement the DLV algorithm is the quasi-Newton algorithm of Broyden, Fletcher, Goldfarb, and Shanno (BFGS) contained in the GAUSS MAXLIK library. The maximum number of iterations for each replication is 100.

Table 1 contains the simulation results when the true value of parameters are used as the initial values for estimation procedure. The results reveal that the bias performance from the DLV algorithm is satisfactory (especially when the sample size is larger) for all configurations considered. Moreover, the associated root-mean-squared error (RMSE) almost always decreases with the increasing sample size. We find only two cases where the pattern of RMSE change is not what we expect, i.e., when $d_0 = 0.4$, the RMSE of estimating the parameters μ_1 and μ_2 as $T = 400$ is found to be a little higher than that of estimating the parameters μ_1 and μ_2 as $T = 200$. These two observations demonstrate the ability of the DLV algorithm to deal with the mixture model considered in this section. The performance of DLV algorithm for estimating the fractional differencing parameter is particularly displayed with the box-plots in Figure 2. The above-mentioned observations are clearly borne out in this figure.

We also check the robustness of the preceding simulation results by changing the choice of initial values for estimation. The simulations in Table 1 are replicated by setting the initial values for parameters at the true values except that of d_0 is set at zero. The results

**Table 1. Finite sample performance of the DLV algorithm:
Initial values of parameters are set at the true values of parameters**

Parameter		μ_1	μ_2	p_{11}	p_{22}	σ_0	d_0	ϕ_1	θ_1
$d_0 = 0.4$									
$T = 100$	Bias	-0.010	-0.106	0.010	0.016	0.008	0.183	-0.128	-0.040
	RMSE	1.008	0.991	0.039	0.060	0.022	0.294	0.233	0.138
$T = 200$	Bias	-0.094	-0.098	0.006	0.006	0.003	0.135	-0.101	-0.026
	RMSE	0.978	0.978	0.028	0.025	0.015	0.233	0.191	0.086
$T = 400$	Bias	-0.074	-0.076	0.004	0.004	0.001	0.096	-0.073	-0.013
	RMSE	0.990	0.990	0.019	0.017	0.010	0.192	0.163	0.060
$d_0 = 0.3$									
$T = 100$	Bias	-0.057	-0.070	0.009	0.017	0.012	0.175	-0.109	-0.041
	RMSE	1.042	1.024	0.037	0.060	0.030	0.319	0.245	0.131
$T = 200$	Bias	-0.058	-0.055	0.006	0.006	0.005	0.122	-0.079	-0.030
	RMSE	0.947	0.944	0.027	0.025	0.020	0.260	0.212	0.086
$T = 400$	Bias	-0.038	-0.043	0.004	0.004	0.001	0.090	-0.061	-0.016
	RMSE	0.885	0.883	0.019	0.017	0.015	0.217	0.185	0.061
$d_0 = 0.2$									
$T = 100$	Bias	-0.017	-0.041	0.009	0.017	0.014	0.201	-0.115	-0.044
	RMSE	0.874	0.853	0.037	0.060	0.037	0.341	0.258	0.128
$T = 200$	Bias	-0.042	-0.047	0.006	0.006	0.006	0.167	-0.106	-0.037
	RMSE	0.795	0.792	0.028	0.025	0.024	0.297	0.239	0.088
$T = 400$	Bias	-0.038	-0.046	0.004	0.004	0.002	0.122	-0.085	-0.019
	RMSE	0.670	0.669	0.019	0.017	0.018	0.239	0.203	0.061

Notes: Simulations are based on 200 replications. The data is generated from the mixture model defined in (15), (16) and (17). DLV algorithm is the Durbin-Levinson-Viterbi algorithm proposed in this paper. Bias is computed as the true parameter minus the corresponding average estimated values.

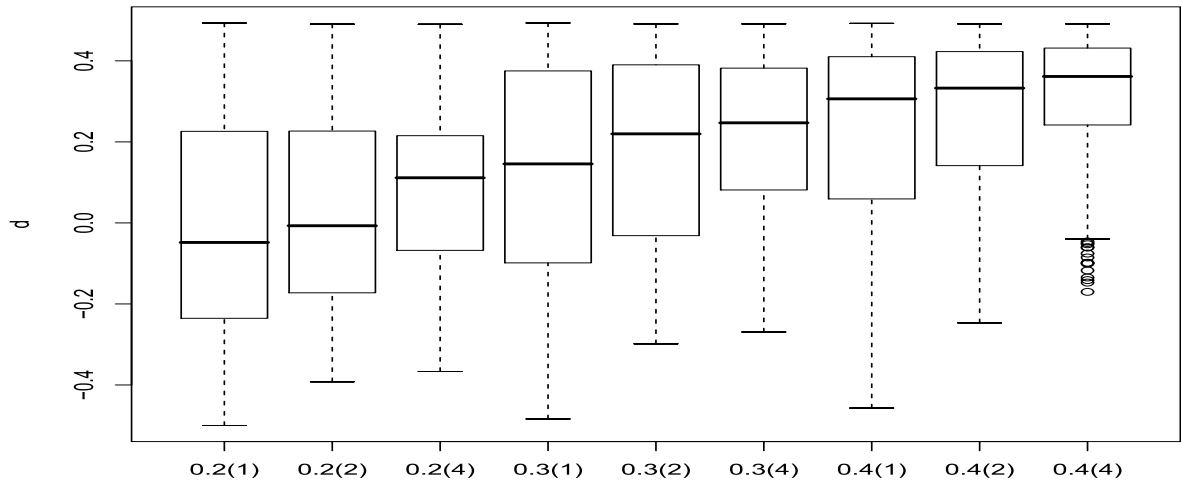


Figure 2: Box-plots of the estimated d from the model defined in (15), (16) and (17) with 200 realizations. The initial values of parameters are set at the true values of parameters. The value $f(g)$ denotes the model specification where $d = f$ and $T = 100 \times g$.

contained Table 2 and Figure 3 indicate that the finite sample performance of our procedure is not sensitive to the initial values used for estimation.

4 Empirical Applications

The methodology developed in this paper is motivated by the dynamic pattern of long memory behavior. Evidence has been given by many methods for such a changing covariance behavior of the Nile river. The applications of the proposed MS-ARFIMA model to actual data are far reaching. For that reason, we consider three data set. The first one is the U.S. real interest rates, the second one is the Nile river data, and the third one is the U.S. unemployment rates.

4.1 Example with real interest rates

In this subsection we first consider the U.S. ex post monthly real interest rate constructed from monthly inflation and Treasury bill rates from January 1953 to December 1990 in Mishkin (1990). The reason we use the original dataset of Mishkin (1990) is to employ it as a benchmark for a clear comparison between the results from the MS-ARFIMA model and

**Table 2. Finite sample performance of the DLV algorithm:
Initial values of parameters are set at the true values of parameters
except that of d_0 is set at zero**

Parameter		μ_1	μ_2	p_{11}	p_{22}	σ_0	d_0	ϕ_1	θ_1
$d_0 = 0.4$									
$T = 100$	Bias	-0.116	-0.122	0.010	0.017	0.009	0.188	-0.130	-0.041
	RMSE	1.030	1.017	0.039	0.060	0.021	0.298	0.235	0.137
$T = 200$	Bias	-0.093	-0.096	0.006	0.006	0.003	0.138	-0.103	-0.027
	RMSE	0.979	0.979	0.028	0.025	0.015	0.238	0.193	0.087
$T = 400$	Bias	-0.074	-0.076	0.004	0.004	0.001	0.096	-0.073	-0.013
	RMSE	0.990	0.990	0.019	0.017	0.010	0.192	0.163	0.060
$d_0 = 0.3$									
$T = 100$	Bias	-0.021	-0.034	0.010	0.017	0.012	0.186	-0.115	-0.040
	RMSE	0.972	0.949	0.039	0.060	0.030	0.325	0.241	0.127
$T = 200$	Bias	-0.046	-0.049	0.006	0.006	0.005	0.126	-0.081	-0.030
	RMSE	0.936	0.937	0.028	0.025	0.021	0.261	0.212	0.086
$T = 400$	Bias	-0.040	-0.044	0.004	0.004	0.002	0.088	-0.059	-0.016
	RMSE	0.912	0.912	0.019	0.017	0.015	0.217	0.184	0.060
$d_0 = 0.2$									
$T = 100$	Bias	-0.018	-0.038	0.009	0.016	0.014	0.195	-0.110	-0.045
	RMSE	0.892	0.864	0.037	0.060	0.037	0.340	0.260	0.130
$T = 200$	Bias	-0.036	-0.040	0.006	0.006	0.006	0.160	-0.100	-0.037
	RMSE	0.804	0.801	0.028	0.025	0.024	0.294	0.238	0.087
$T = 400$	Bias	-0.044	-0.051	0.004	0.004	0.002	0.117	-0.082	-0.018
	RMSE	0.674	0.673	0.019	0.017	0.018	0.235	0.200	0.060

Notes: Simulations are based on 200 replications. The data is generated from the mixture model defined in (15), (16) and (17). DLV algorithm is the Durbin-Levinson-Viterbi algorithm proposed in this paper. Bias is computed as the true parameter minus the corresponding average estimated values.

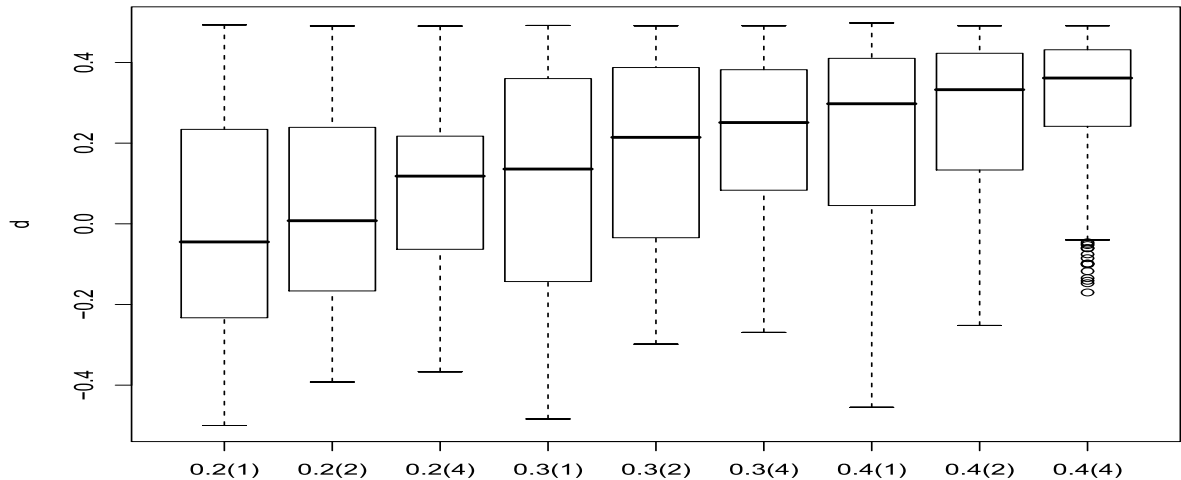


Figure 3: Box-plots of the estimated d from the model defined in (15), (16) and (17) with 200 realizations. The initial values of parameters are set at the true values except that of d_0 is set at zero. The value $f(g)$ denotes the model specification where $d = f$ and $T = 100 \times g$.

those generated from the methodology employed in earlier papers.

The main feature of the real interest rate is that the whole dataset can be split into three subperiods, January 1953-October 1979, November 1979-October 1982, and November 1982-December 1990, because the operating procedure of the monetary authority changed in October 1979 and October 1982 as argued in Mishkin (1990). Another interesting feature of the real interest rate is that the data of these three subperiods can be well described with the ARFIMA models as shown in Tsay (2000). The simultaneous presence of structural break and long memory within the real interest rate allows itself to be an ideal subject to be investigated with the MS-ARFIMA model.

Allowing the break points to be endogenously determined, Table 3 contains the parameter estimates from the following mixture model with a 2-state Markov chain and an ARFIMA(1, d , 1) noise:

$$w_t = \mu_{s_t} I\{t \geq 1\} + (1 - L)^{-d_0} \sigma_0 z_t I\{t \geq 1\}, \quad (1 - \phi_1 L) z_t = (1 + \theta_1 L) \varepsilon_t, \quad (18)$$

where ϕ_1 or θ_1 is assumed to be zero depending on the noise specification. Following Hamilton (1989), asymptotic standard errors are calculated numerically.

Table 3 shows that the estimates of μ_1 , μ_2 , p_{11} , p_{22} , σ_0 , and d_0 from the DLV algorithm

Table 3. Estimates of Parameters Based on Data for U.S. Monthly Real Interest Rate and the DLV Algorithm

	ARFIMA(0, d , 0)		ARFIMA(0, d , 1)		ARFIMA(1, d , 0)		ARFIMA(1, d , 1)	
	Estimate	S.E.	Estimate	S.E.	Estimate	S.E.	Estimate	S.E.
μ_1	5.3455	0.7494	5.3168	0.7162	5.3116	0.7124	5.3626	0.7706
μ_2	0.7226	0.4814	0.7194	0.4383	0.7184	0.4322	0.7352	0.4958
p_{11}	0.9833	0.0150	0.9833	0.0150	0.9833	0.0150	0.9833	0.0150
p_{22}	0.9977	0.0023	0.9977	0.0023	0.9977	0.0023	0.9977	0.0023
σ_0	2.5094	0.0831	2.5091	0.0831	2.5091	0.0831	2.4979	0.0827
d_0	0.2225	0.0367	0.2062	0.0520	0.2034	0.0653	0.2337	0.0376
ϕ_1	-	-	-	-	0.0324	0.0946	-0.9847	0.0155
θ_1	-	-	0.0279	0.0663	-	-	0.9675	0.0200
L^*	1079.0875		1079.0009		1078.9918		1077.0173	

Notes: The results are based on the MS-ARFIMA model defined in (18). S.E. stands for the standard error of the estimate. L^* represents the negative of the log-likelihood function of the switching model. DLV algorithm is the Durbin-Levinson-Viterbi algorithm proposed in this paper.

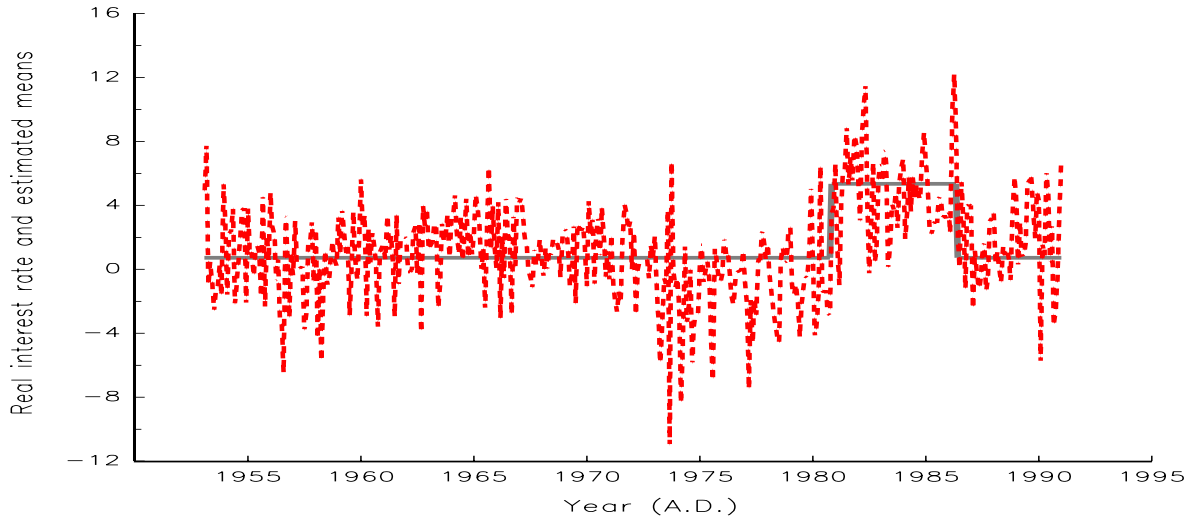


Figure 4: US monthly ex post real interest rates, January 1953-December 1990. Solid line denotes the path of estimated switching means from the specification $ARFIMA(0, d, 0)$ in Table 3, while dotted line denotes the observed monthly ex post real interest rates.

are quite robust across all 4 different configurations. More importantly, two identical break points are identified with these four models, thus divide the whole data into three subperiods as suggested in Mishkin (1990). The endogenous break points identified are November 1980 and May 1986, respectively.

Figure 4 displays the U.S. monthly ex post real interest rates and the path of estimated switching means generated from the DLV algorithm. Without loss of generality, only the path of the estimated switching means from the specification $ARFIMA(0, d, 0)$ in Table 3 is reported. Figure 4 shows that the model in (18) provides a satisfactory fitting of the U.S. monthly real interest rates. Although the endogenously identified break points are later than the well-known monetary operating procedure change points (October 1979 and October 1982), this finding is quite reasonable, because it takes some time for the ex post real interest rate to adjust its path after new information arrives. This argument is buttressed with the findings in Figure 4 that the endogenously identified break points are more closely connected to the observed path of the U.S. monthly ex post real interest rates than the monetary operating procedure change points are.

Table 3 also shows that a long memory phenomenon is found in the real interest rate as has been documented in Tsay (2000). Nevertheless, the estimate of the fractional differencing parameter in Table 3 is much lower than that of 0.666 in Table 3 of Tsay (2000) where the

change points are exogenously determined, and it is more in line with the estimates of 0.204, 0.275, and 0.193 from the individual subperiod data presented in Table 3 of Tsay (2000). This implies that the persistence of long memory in the real interest rate is much more mitigated, once we take the potentially switching mean of the data into account, thus confirming the arguments of Diebold and Inoue (2001) that the presence of Markov-switching level might increase the persistence of the data under investigation.

4.2 Example with Nile river data

In this subsection we apply the Viterbi algorithm to the Nile river data with the following model:

$$w_t = \mu_{s_t} I\{t \geq 1\} + (1 - L)^{-d_{s_t}} \sigma_{s_t} \varepsilon_t I\{t \geq 1\}, \quad (19)$$

where N is assumed to be 2. For the purpose of comparison, we estimate a fixed regime ARFIMA(0, d , 0) model for the Nile river data, i.e., $N = 1$ is imposed on this model. The estimated value of d from such a fixed regime ARFIMA(0, d , 0) model is 0.3986 and is almost identical to the finding in Beran and Terrin (1996).

When estimating the model in (19) with the Viterbi algorithm, we find that the value of the differencing parameter in Table 4 is 0.5770 (nonstationary) for one state, and is 0.2143 (stationary) for the other one. In addition, we identify 5 transitions within the Nile river data in the year 720, 805, 815, 878, and 1070. The estimated path of d_{s_t} from the MS-ARFIMA(0, d , 0) model in Table 4 is graphed in Figure 5.

Most impressively, the first transition data occurs in the year of 720, and the associated estimated value of d_{s_t} within the period 622 to 719 is 0.2143 which is lower than the 0.5770 observed in the other regime. These two findings correspond closely to the conjectures in Beran and Terrin (1996) that the observations 1 to about 100 seem to be more independent than the subsequent observations and the value of differencing parameter might be lower for the first 100 observations than for the subsequent data.

In Figures 6 and 7 we present the observations and the fitted values generated from the estimated parameters displayed in Table 4. It is clear that the fitted value from the MS-ARFIMA(0, d , 0) model is much closer to the real data than that generated from the model whose differencing parameter is not Markov switching. Combining the findings of the likelihood values in Table 4, we find that the MS-ARFIMA(0, d , 0) model is a promising

Table 4. Estimates of MS-ARFIMA(0, d , 0) Model based on the Nile River Data

	MS-ARFIMA(0, d , 0)		ARFIMA(0, d , 0)	
	Estimate	S.E.	Estimate	S.E.
μ_1	10.8593	0.6903	11.4847	0.2607
μ_2	11.4939	0.0917	-	-
p_{11}	0.9930	0.0042	-	-
p_{22}	0.9918	0.0050	-	-
σ_1	0.5430	0.0202	0.6995	0.0192
σ_2	0.8143	0.0332	-	-
d_1	0.5770	0.0430	0.3986	0.0309
d_2	0.2143	0.0510	-	-
L^*	687.5642		703.8541	

Notes: The MS-ARFIMA(0, d , 0) model is defined in (19). S.E. stands for the standard error of the estimate based on numerical derivative. L^* represents the negative of the log-likelihood function of the estimated model.

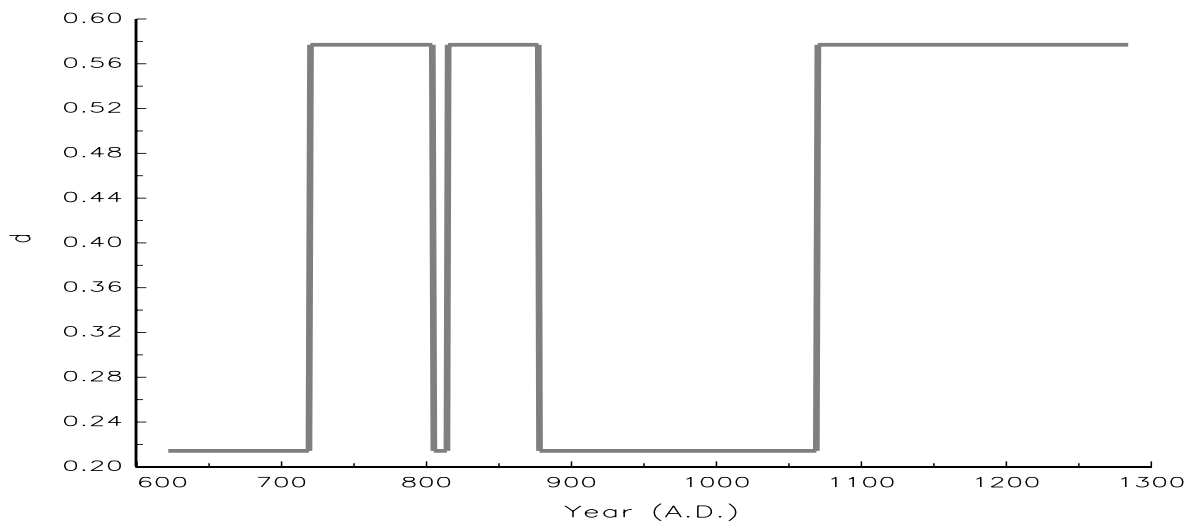


Figure 5: Estimated d_{st} from the MS-ARFIMA(0, d , 0) model in Table 4.

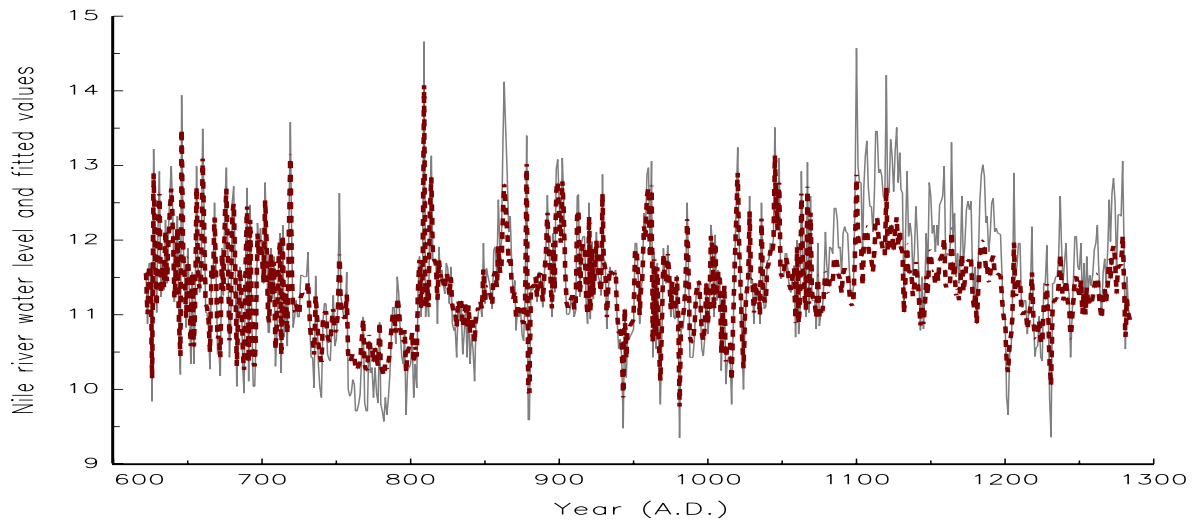


Figure 6: Solid line denotes the Nile river water level divided by 100, while dotted line denotes the corresponding fitted values from the MS-ARFIMA(0, d , 0) model in Table 4.

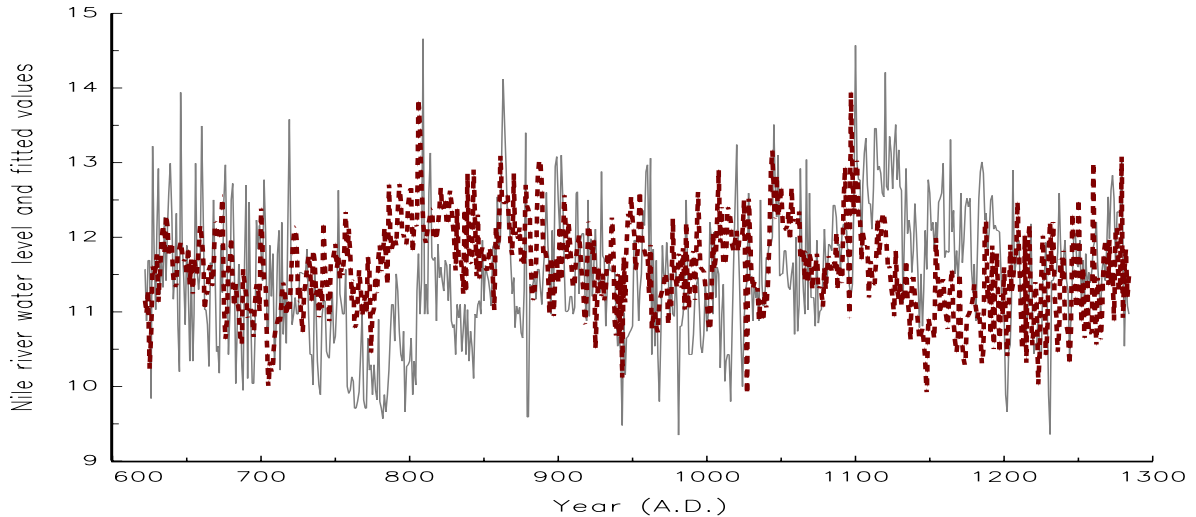


Figure 7: Solid line denotes the Nile river water level divided by 100, while dotted line denotes the corresponding fitted values from the ARFIMA(0, d , 0) model in Table 4.

alternative to describe the Nile river data.

4.3 Example with unemployment rates

In this subsection we apply the Viterbi algorithm to the U.S. quarterly unemployment rates from 1948 to 2006. This data is based on the monthly unemployment rates contained in *Bureau of Labour Statistics* as those employed in van Dijk et al. (2002) for estimating a fractionally integrated smooth transition autoregressive (FI-STAR) model. However, van Dijk et al (2002) employ the original monthly unemployment rates ranging from July 1986 to December 1999, while we use all the data contained in *Bureau of Labour Statistics*, but focusing on the quarterly frequency usually considered in the business cycle related studies.

As clearly argued in van Dijk et al. (2002) and shown in Figure 8, there are two important empirical features of U.S. unemployment rates, i.e., the shocks to the series is quite persistent and the series seem to rise faster during recessions than it falls during expansions. van Dijk et al. (2002) find that the estimated d is 0.43 from a FI-STAR model presented in their Table 1. This implies that a time series model describing long memory and nonlinearity simultaneously may be useful for modeling U.S. unemployment rates and many other applications.

The aforementioned two features contained in U.S. unemployment also provide another good opportunity to test the applicability of the MS-ARFIMA model. As a consequence we estimate the U.S. quarterly unemployment rates with the following MS-ARFIMA($p, d, 0$)

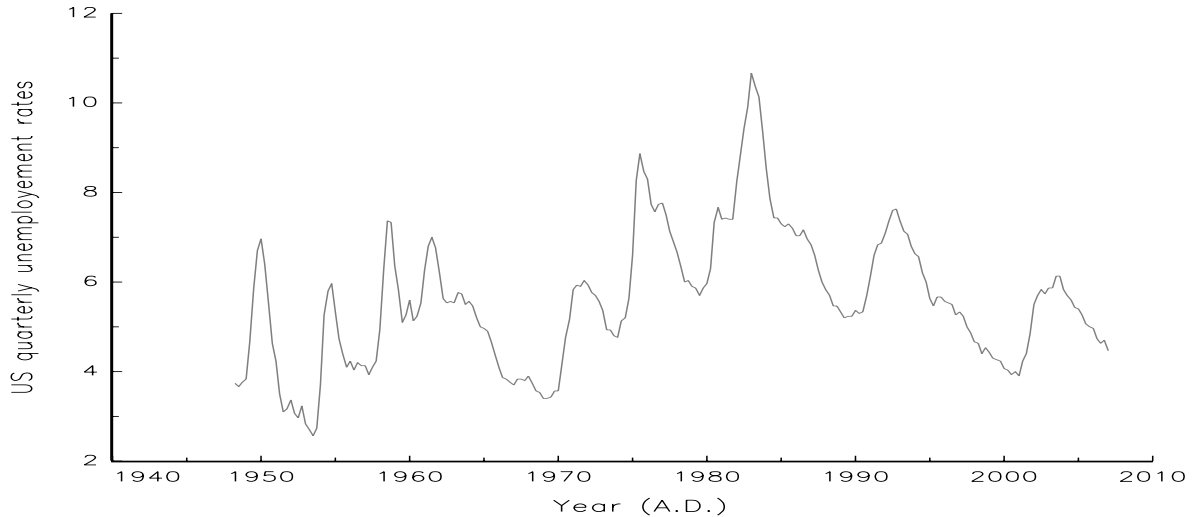


Figure 8: U.S. quarterly seasonally adjusted unemployment rates, 1948-2006.

model:

$$w_t = \mu_{s_t} I\{t \geq 1\} + (1 - L)^{-d_{s_t}} \sigma_{s_t} \phi(B)^{-1} \varepsilon_t I\{t \geq 1\}, \quad (20)$$

where N is assumed to be 2, and $p = \{3, 4\}$. The choice of $p = 4$ is adopted by following the model specification in (30) of van Dijk et al (2002), while $p = 3$ is chosen to check the robustness of the estimation results from the specification $p = 4$. The major objective of this subsection is to investigate whether the long memory observed in van Dijk et al. (2002) can also be retained from the MS-ARFIMA methodology.

When estimating the model in (20) with the Viterbi algorithm, we find that the values of the estimated fractional differencing parameter from both MS-ARFIMA(3, d , 0) and MS-ARFIMA(4, d , 0) models in Table 5 are very close to that found in van Dijk et al. (2002), thus confirming that long memory phenomenon seems to be present in the U.S. unemployment rates. For clarity of exposition, the estimated path of d_{s_t} from the MS-ARFIMA(3, d , 0) model and that of d_{s_t} from the MS-ARFIMA(4, d , 0) one are graphed in Figure 9 and Figure 10, respectively. These figures clearly show that d_{s_t} are around 0.4-0.5 for both regimes estimated in each MS-ARFIMA(p , d , 0) model in Table 5.

We also check to what extent the fitted values generated from the models in Table 5 can capture the feature of U.S. unemployment rates. This task is not taken in van Dijk et al. (2002) when estimating their FI-STAR model for the U.S. monthly unemployment rates. It is interesting to find in Figure 11 and Figure 12 that the MS-ARFIMA(p , d , 0) model in (20)

Table 5. Estimates of MS-ARFIMA($p, d, 0$) Model based on the U.S. quarterly unemployment rates

	MS-ARFIMA(3, d , 0)		MS-ARFIMA(4, d , 0)	
	Estimate	S.E.	Estimate	S.E.
μ_1	3.8080	0.1552	3.4572	0.3711
μ_2	5.1358	0.4403	3.8254	0.3334
p_{11}	0.9939	0.0067	0.9877	0.0093
p_{22}	0.9896	0.0083	0.9867	0.0101
σ_1	0.1973	0.0135	0.1535	0.0101
σ_2	0.3380	0.0206	0.3921	0.0274
d_1	0.4919	0.1215	0.4429	0.0987
d_2	0.4143	0.1337	0.4342	0.1058
ϕ_1	1.2570	0.1415	1.1325	0.1215
ϕ_2	-0.3822	0.1510	-0.2301	0.1239
ϕ_3	-0.0666	0.0788	-0.0141	0.1053
ϕ_4	-	-	-0.0495	0.0712
L^*	36.1377		14.7662	

Notes: The results are based on the MS-ARFIMA($p, d, 0$) model defined in (20). S.E. stands for the standard error of the estimate based on numerical derivative. L^* represents the negative of the log-likelihood function of the estimated model.

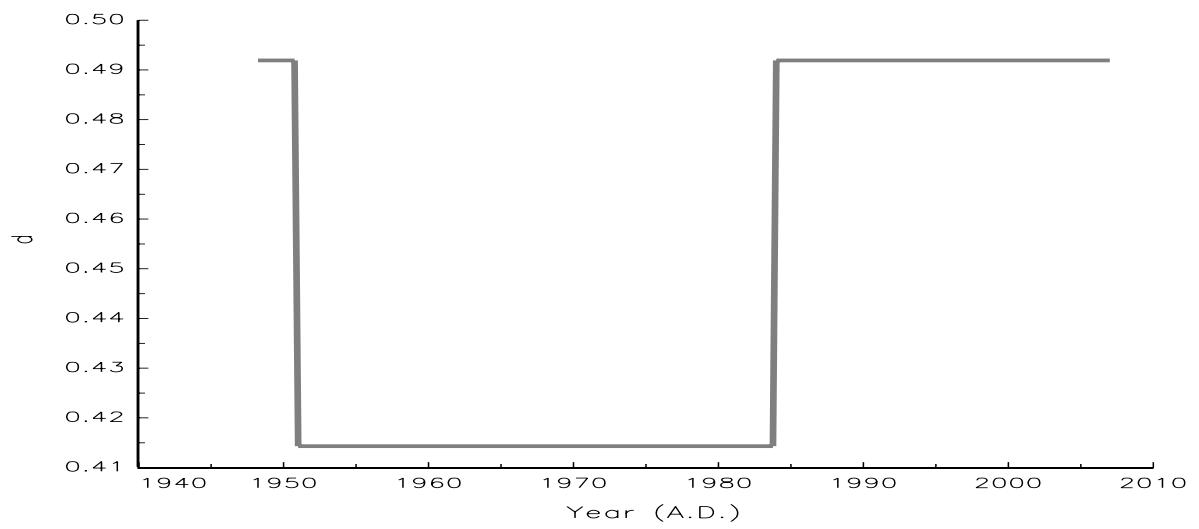


Figure 9: Estimated d_{st} from the MS-ARFIMA(3, d , 0) model in Table 5.

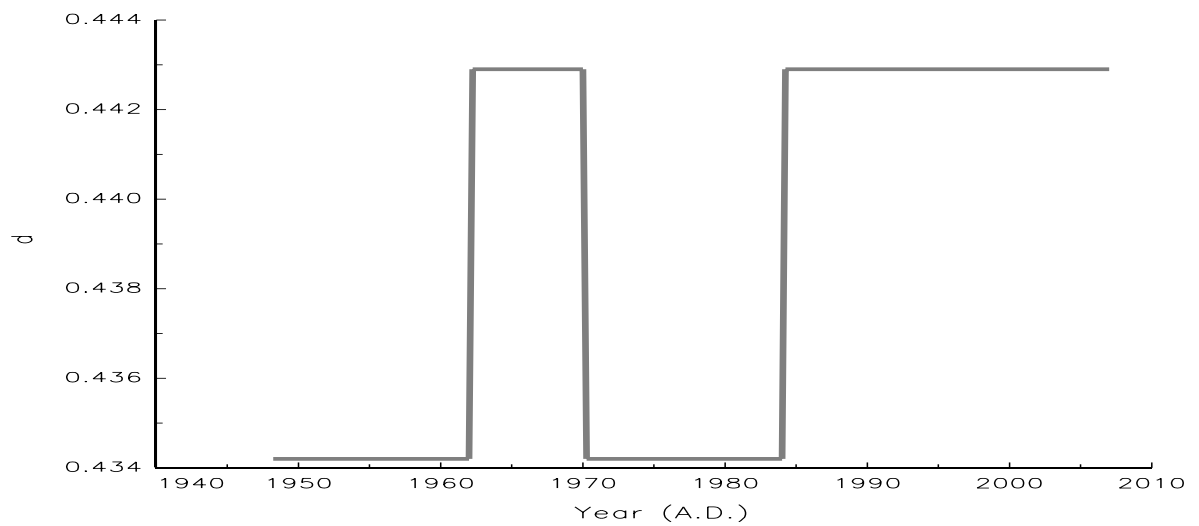


Figure 10: Estimated d_{st} from the MS-ARFIMA(4, d , 0) model in Table 5.

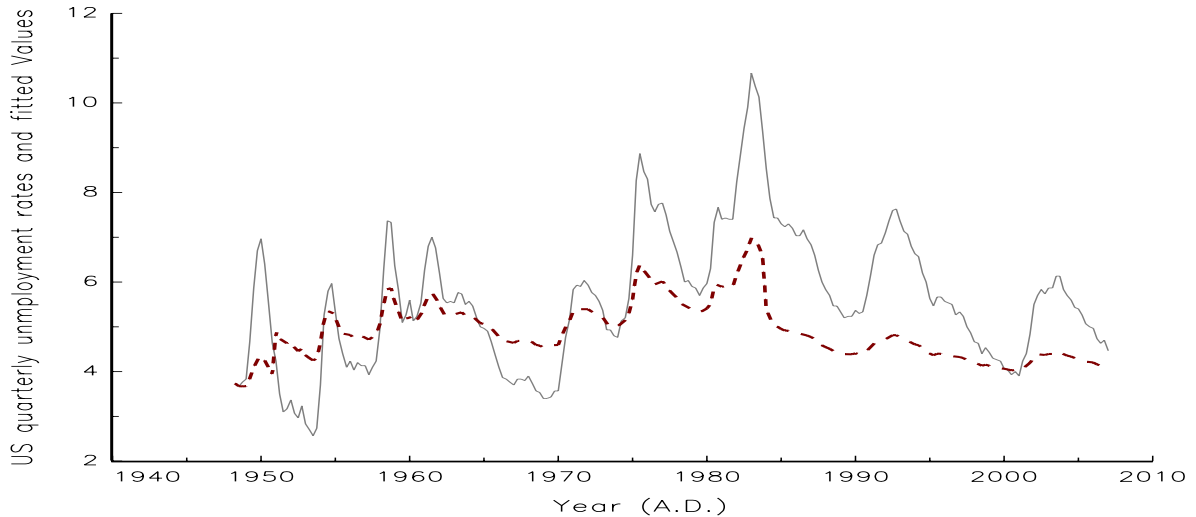


Figure 11: Solid line denotes the U.S. quarterly seasonally adjusted unemployment rates (1948-2006), while dotted line denotes the corresponding fitted values from the MS-ARFIMA(3, d , 0) model in Table 5.

provides a reasonable fit to the data, even though we do not include some seasonal control variables, like seasonal difference operator, as van Dijk et al. (2002) have done for their empirical studies.

5 Conclusions

A general class of MS-ARFIMA processes is suggested to combine long memory and Markov-switching models into one unified framework. The coverage of this class of MS-ARFIMA models is far-reaching, but we show that they still can be easily estimated with the original Viterbi algorithm or the DLV algorithm proposed in this paper. In addition, the simulation reveals that the finite sample performance of the DLV algorithm for a simple mixture model of Markov-switching mean and ARFIMA(1, d , 1) process is satisfactory. When applying the MS-ARFIMA models to the U.S. real interest rates, the Nile river level, and the U.S. unemployment rates, the estimation results are both highly compatible with the conjectures made in the literature. Accordingly, the MS-ARFIMA model considered in this paper not only can be used for solving the puzzle raised by Diebold and Inoue (2001), but can also find many potential applications in several scientific research fields.

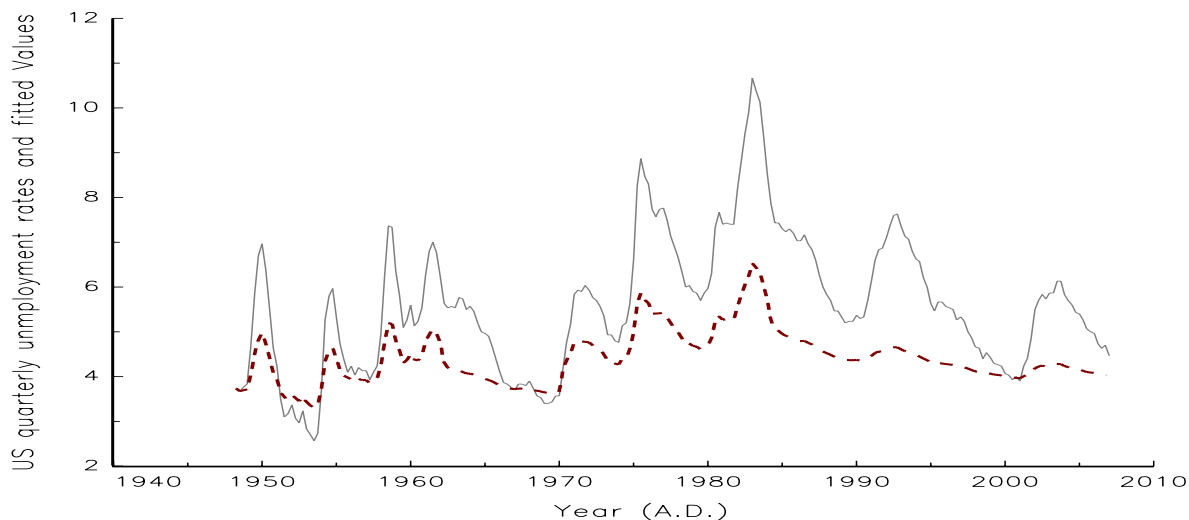


Figure 12: Solid line denotes the U.S. quarterly seasonally adjusted unemployment rates (1948-2006), while dotted line denotes the corresponding fitted values from the MS-ARFIMA(4, d , 0) model in Table 5.

References

- Beran, J. (1994), *Statistics for Long-Memory Processes*. New York: Chapman and Hall.
- Beran, J. and Terrin, N. (1996), “Testing for a Change of the Long-Memory Parameter”, *Biometrika*, 83, 627-638.
- Berkes, I., Horváth, L., Kokoszka, P. and Shao, Q.M. (2006), “On Discriminating between Long-Range Dependence and Changes in Mean”, *The Annals of Statistics*, 34, 1140-1165.
- Bhattacharya, R.N., Gupta, V.K. and Waymire, E. (1983), “The Hurst Effect under Trends”, *Journal of Applied Probability*, 20, 649-662.
- Bollerslev, T. and Mikkelsen, H.O.A. (1996), “Modeling and Pricing Long-Memory in Stock Market Volatility”, *Journal of Econometrics*, 73, 151-184.
- Breidt, F.J., Crato, N. and de Lima, P. (1998), “The Detection and Estimation of Long Memory in Stochastic Volatility”, *Journal of Econometrics*, 83, 325-348.
- de Lima, P. and Crato, N. (1993), “Long-Range Dependence in the Conditional Variance of Stock Returns”, *Proceedings of the Business and Economic Statistics Section*, August 1993 Joint Statistical Meetings, San Francisco.

- Deo, R., Hurvich, C. and Lu, Y. (2006), “Forecasting Realized Volatility Using a Long-Memory Stochastic Volatility Model: Estimation, Prediction and Seasonal Adjustment”, *Journal of Econometrics* 131, 29-58.
- Deriche, J.A. and Tewfik, A.H. (1993), “Maximum Likelihood Estimation of the Parameters of Discrete Fractionally Differenced Gaussian Noise Process,” *IEEE Transactions on Signal Processing*, 41, 2977-2989.
- Diebold, F.X. and Inoue, A. (2001), “Long Memory and Regime Switching”, *Journal of Econometrics*, 105, 131-159.
- Ding, Z., Granger, C.W.J. and Engle, R.F. (1993), “A Long Memory Property of Stock Market Returns and a New Model”, *Journal of Empirical Finance*, 1, 83-106.
- Dueker, M. and Serletis, A. (2000), “Do Real Exchange Rates Have Autoregressive Unit Roots? A Test under the Alternative of Long Memory and Breaks”, Working Paper 2000-016A, Federal Reserve Bank of St. Louis.
- Granger, C.W.J. (1980), “Long Memory Relationships and the Aggregation of Dynamic Models”, *Journal of Econometrics*, 14, 227-238.
- Granger, C.W.J. and Joyeux, R. (1980), “An Introduction to Long-Memory Time Series Models and Fractional Differencing”, *Journal of Time Series Analysis*, 1, 15-29.
- Hamilton, J.D. (1989), “A New Approach to the Economic Analysis of Nonstationary Time Series and the Business Cycle”, *Econometrica*, 57, 357-384.
- Heyde, C.C. and Dai, W. (1996), “On the Robustness to Small Trends of Estimation based on the Smoothed Periodogram”, *Journal of Time Series Analysis*, 17, 141-150.
- Horváth, L. (2001), “Change-Point Detection in Long-Memory Processes. *Journal of Multivariate Analysis*, 78, 218-234.
- Horváth, L. and Shao, Q.M. (1999), “Limit Theorems for the Union-Intersection Test”, *Journal of Statistical Planning and Inference*, 44, 133-148.
- Hosking, J.R.M. (1981), “Fractional Differencing”, *Biometrika*, 68, 165-176.

- Hurst, H.E. (1951), “Long-Term Storage Capacity of Reservoirs”, *Transactions of the American Society of Civil Engineers*, 116, 770-799.
- Juang, B.H. and Rabiner, L.R. (1991), “Hidden Markov Models for Speech Recognition”, *Technometrics*, 33, 251-272.
- Künsch, H. (1986), “Discrimination between Monotonic Trends and Long-Range Dependence”, *Journal of Applied Probability*, 23, 1025-1030.
- Mandelbrot, B.B. and van Ness, J.W. (1968), “Fractional Brownian Motions, Fractional Noises and Applications”, *SIAM Review*, 10, 422-437.
- Mandelbrot, B.B. and Wallis, J.R. (1969), “Some Long-Run Properties of Geophysical Records”, *Water Resources Research*, 5, 321-340.
- Mishkin, F.S. (1990), “What does the Term Structure of Interest Rate Tell Us about Future Inflation? *Journal of Monetary Economics*, 25, 77-95.
- Qian, W. and Titterton, D.M. (1991), “Estimation of Parameters in Hidden Markov Models”, *Philosophical Transactions: Physical Sciences and Engineering*, 337, 407-428.
- Ray, B.K. and Tsay, R.S. (2002), “Bayesian Methods for Change-Point Detection in Long-Range Dependent Process”, *Journal of Time Series Analysis*, 23, 687-705.
- Robert, C.P., Rydén, R. and Titterton, D.M. (2000), “Bayesian Inference in Hidden Markov Models through the Reversible Jump Markov Chain Monte Carlo Method”, *Journal of the Royal Statistical Society B*, 62, 57-75.
- Shimotsu, K. and Phillips, P.C.B. (2005), “Exact Local Whittle Estimation of Fractional Integration”, *The Annals of Statistics*, 33, 1890-1933.
- Tsay, W.J. (2000), “Long Memory Story of the Real Interest Rate”, *Economic Letters*, 67, 325-330.
- van Dijk, D., Franses, P.H. and Raap, R. (2002), “A Nonlinear Long Memory Model, with an Application to US Unemployment”, *Journal of Econometrics*, 110, 135-165.
- Viterbi, A.J. (1967), “Error Bounds for Convolutional Codes and an Asymptotic Optimum Decoding Algorithm”, *IEEE Transactions on Signal Processing*, IT-13, 260-269.

SFB 649 Discussion Paper Series 2007

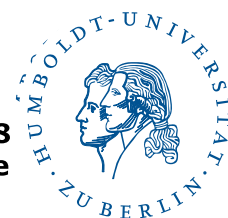
For a complete list of Discussion Papers published by the SFB 649, please visit <http://sfb649.wiwi.hu-berlin.de>.

- 001 "Trade Liberalisation, Process and Product Innovation, and Relative Skill Demand" by Sebastian Braun, January 2007.
- 002 "Robust Risk Management. Accounting for Nonstationarity and Heavy Tails" by Ying Chen and Vladimir Spokoiny, January 2007.
- 003 "Explaining Asset Prices with External Habits and Wage Rigidities in a DSGE Model." by Harald Uhlig, January 2007.
- 004 "Volatility and Causality in Asia Pacific Financial Markets" by Enzo Weber, January 2007.
- 005 "Quantile Sieve Estimates For Time Series" by Jürgen Franke, Jean-Pierre Stockis and Joseph Tadjuidje, February 2007.
- 006 "Real Origins of the Great Depression: Monopolistic Competition, Union Power, and the American Business Cycle in the 1920s" by Monique Ebell and Albrecht Ritschl, February 2007.
- 007 "Rules, Discretion or Reputation? Monetary Policies and the Efficiency of Financial Markets in Germany, 14th to 16th Centuries" by Oliver Volckart, February 2007.
- 008 "Sectoral Transformation, Turbulence, and Labour Market Dynamics in Germany" by Ronald Bachmann and Michael C. Burda, February 2007.
- 009 "Union Wage Compression in a Right-to-Manage Model" by Thorsten Vogel, February 2007.
- 010 "On σ -additive robust representation of convex risk measures for unbounded financial positions in the presence of uncertainty about the market model" by Volker Krätschmer, March 2007.
- 011 "Media Coverage and Macroeconomic Information Processing" by Alexandra Niessen, March 2007.
- 012 "Are Correlations Constant Over Time? Application of the CC-TRIG_t-test to Return Series from Different Asset Classes." by Matthias Fischer, March 2007.
- 013 "Uncertain Paternity, Mating Market Failure, and the Institution of Marriage" by Dirk Bethmann and Michael Kvasnicka, March 2007.
- 014 "What Happened to the Transatlantic Capital Market Relations?" by Enzo Weber, March 2007.
- 015 "Who Leads Financial Markets?" by Enzo Weber, April 2007.
- 016 "Fiscal Policy Rules in Practice" by Andreas Thams, April 2007.
- 017 "Empirical Pricing Kernels and Investor Preferences" by Kai Detlefsen, Wolfgang Härdle and Rouslan Moro, April 2007.
- 018 "Simultaneous Causality in International Trade" by Enzo Weber, April 2007.
- 019 "Regional and Outward Economic Integration in South-East Asia" by Enzo Weber, April 2007.
- 020 "Computational Statistics and Data Visualization" by Antony Unwin, Chun-houh Chen and Wolfgang Härdle, April 2007.
- 021 "Ideology Without Ideologists" by Lydia Mechtenberg, April 2007.
- 022 "A Generalized ARFIMA Process with Markov-Switching Fractional Differencing Parameter" by Wen-Jen Tsay and Wolfgang Härdle, April 2007.

SFB 649, Spandauer Straße 1, D-10178 Berlin
<http://sfb649.wiwi.hu-berlin.de>

This research was supported by the Deutsche
Forschungsgemeinschaft through the SFB 649 "Economic Risk".

SFB 649, Spandauer Straße 1, D-10178
<http://sfb649.wiwi.hu-berlin.de>



This research was supported by the Deutsche
Forschungsgemeinschaft through the SFB 649 "Economic Risk".

COMMON FUNCTIONAL PRINCIPAL COMPONENTS¹

BY MICHAL BENKO, WOLFGANG HÄRDLE AND ALOIS KNEIP

Humboldt-Universität, Humboldt-Universität and Bonn Universität

Functional principal component analysis (FPCA) based on the Karhunen–Loève decomposition has been successfully applied in many applications, mainly for one sample problems. In this paper we consider common functional principal components for two sample problems. Our research is motivated not only by the theoretical challenge of this data situation, but also by the actual question of dynamics of implied volatility (IV) functions. For different maturities the log-returns of IVs are samples of (smooth) random functions and the methods proposed here study the similarities of their stochastic behavior. First we present a new method for estimation of functional principal components from discrete noisy data. Next we present the two sample inference for FPCA and develop the two sample theory. We propose bootstrap tests for testing the equality of eigenvalues, eigenfunctions, and mean functions of two functional samples, illustrate the test-properties by simulation study and apply the method to the IV analysis.

1. Introduction. In many applications in biometrics, chemometrics, econometrics, etc., the data come from the observation of continuous phenomena of time or space and can be assumed to represent a sample of i.i.d. smooth random functions $X_1(t), \dots, X_n(t) \in L^2[0, 1]$. Functional data analysis has received considerable attention in the statistical literature during the last decade. In this context functional principal component analysis (FPCA) has proved to be a key technique. An early reference is Rao (1958), and important methodological contributions have been given by various authors. Case studies and references, as well as methodological and algorithmical details, can be found in the books by Ramsay and Silverman (2002, 2005) or Ferraty and Vieu (2006).

Received January 2006; revised February 2007.

¹Supported by the Deutsche Forschungsgemeinschaft and the Sonderforschungsbereich 649 “Ökonomisches Risiko.”

AMS 2000 subject classifications. Primary 62H25, 62G08; secondary 62P05.

Key words and phrases. Functional principal components, nonparametric regression, bootstrap, two sample problem.

<p>This is an electronic reprint of the original article published by the Institute of Mathematical Statistics in <i>The Annals of Statistics</i>, 2009, Vol. 37, No. 1, 1–34. This reprint differs from the original in pagination and typographic detail.</p>

The well-known Karhunen–Loève (KL) expansion provides a basic tool to describe the distribution of the random functions X_i and can be seen as the theoretical basis of FPCA. For $v, w \in L^2[0, 1]$, let $\langle v, w \rangle = \int_0^1 v(t)w(t) dt$, and let $\|\cdot\| = \langle \cdot, \cdot \rangle^{1/2}$ denote the usual L^2 -norm. With $\lambda_1 \geq \lambda_2 \geq \dots$ and $\gamma_1, \gamma_2, \dots$ denoting eigenvalues and corresponding orthonormal eigenfunctions of the covariance operator Γ of X_i , we obtain $X_i = \mu + \sum_{r=1}^{\infty} \beta_{ri} \gamma_r$, $i = 1, \dots, n$, where $\mu = E(X_i)$ is the mean function and $\beta_{ri} = \langle X_i - \mu, \gamma_r \rangle$ are (scalar) factor loadings with $E(\beta_{ri}^2) = \lambda_r$. Structure and dynamics of the random functions can be assessed by analyzing the “functional principal components” γ_r , as well as the distribution of the factor loadings. For a given functional sample, the unknown characteristics λ_r, γ_r are estimated by the eigenvalues and eigenfunctions of the empirical covariance operator $\hat{\Gamma}_n$ of X_1, \dots, X_n . Note that an eigenfunction γ_r is identified (up to sign) only if the corresponding eigenvalue λ_r has multiplicity one. This therefore establishes a necessary regularity condition for any inference based on an estimated functional principal component $\hat{\gamma}_r$ in FPCA. Signs are arbitrary (γ_r and β_{ri} can be replaced by $-\gamma_r$ and $-\beta_{ri}$) and may be fixed by a suitable standardization. More detailed discussion on this topic and precise assumptions can be found in Section 2.

In many important applications a small number of functional principal components will suffice to approximate the functions X_i with a high degree of accuracy. Indeed, FPCA plays a much more central role in functional data analysis than its well-known analogue in multivariate analysis. There are two major reasons. First, distributions on function spaces are complex objects, and the Karhunen–Loève expansion seems to be the only practically feasible way to access their structure. Second, in multivariate analysis a substantial interpretation of principal components is often difficult and has to be based on vague arguments concerning the correlation of principal components with original variables. Such a problem does not at all exist in the functional context, where $\gamma_1(t), \gamma_2(t), \dots$ are functions representing the major modes of variation of $X_i(t)$ over t .

In this paper we consider inference and tests of hypotheses on the structure of functional principal components. Motivated by an application to implied volatility analysis, we will concentrate on the two sample case. A central point is the use of bootstrap procedures. We will show that the bootstrap methodology can also be applied to functional data.

In Section 2 we start by discussing one-sample inference for FPCA. Basic results on asymptotic distributions have already been derived by [Dauxois, Pousse and Romain \(1982\)](#) in situations where the functions are directly observable. [Hall and Hosseini-Nasab \(2006\)](#) develop asymptotic Taylor expansions of estimated eigenfunctions in terms of the difference $\hat{\Gamma}_n - \Gamma$.

Without deriving rigorous theoretical results, they also provide some qualitative arguments as well as simulation results motivating the use of bootstrap in order to construct confidence regions for principal components.

In practice, the functions of interest are often not directly observed, but are regression curves which have to be reconstructed from discrete, noisy data. In this context the standard approach is to first estimate individual functions nonparametrically (e.g., by B-splines) and then to determine principal components of the resulting estimated empirical covariance operator—see [Besse and Ramsay \(1986\)](#), [Ramsay and Dalzell \(1991\)](#), among others. Approaches incorporating a smoothing step into the eigenanalysis have been proposed by [Rice and Silverman \(1991\)](#), [Pezzulli and Silverman \(1993\)](#) or [Silverman \(1996\)](#). Robust estimation of principal components has been considered by [Lacontore et al. \(1999\)](#). [Yao, Müller and Wang \(2005\)](#) and [Hall, Müller and Wang \(2006\)](#) propose techniques based on nonparametric estimation of the covariance function $E[\{X_i(t) - \mu(t)\}\{X_i(s) - \mu(s)\}]$ which can also be applied if there are only a few scattered observations per curve.

Section 2.1 presents a new method for estimation of functional principal components. It consists in an adaptation of a technique introduced by [Kneip and Utikal \(2001\)](#) for the case of density functions. The key-idea is to represent the components of the Karhunen–Loève expansion in terms of an (L^2) scalar-product matrix of the sample. We investigate the asymptotic properties of the proposed method. It is shown that under mild conditions the additional error caused by estimation from discrete, noisy data is first-order asymptotically negligible, and inference may proceed “as if” the functions were directly observed. Generalizing the results of [Dauxois, Pousse and Romain \(1982\)](#), we then present a theorem on the asymptotic distributions of the empirical eigenvalues and eigenfunctions. The structure of the asymptotic expansion derived in the theorem provides a basis to show consistency of bootstrap procedures.

Section 3 deals with two-sample inference. We consider two independent samples of functions $\{X_i^{(1)}\}_{i=1}^{n_1}$ and $\{X_i^{(2)}\}_{i=1}^{n_2}$. The problem of interest is to test in how far the distributions of these random functions coincide. The structure of the different distributions in function space can be accessed by means of the respective Karhunen–Loève expansions

$$X_i^{(p)} = \mu^{(p)} + \sum_{r=1}^{\infty} \beta_{ri}^{(p)} \gamma_r^{(p)}, \quad p = 1, 2.$$

Differences in the distribution of these random functions will correspond to differences in the components of the respective KL expansions above. Without restriction, one may require that signs are such that $\langle \gamma_r^{(1)}, \gamma_r^{(2)} \rangle \geq 0$. Two sample inference for FPCA in general has not been considered in the literature so far. In Section 3 we define bootstrap procedures for testing

the equality of mean functions, eigenvalues, eigenfunctions and eigenspaces. Consistency of the bootstrap is derived in Section 3.1, while Section 3.2 contains a simulation study providing insight into the finite sample performance of our tests.

It is of particular interest to compare the functional components characterizing the two samples. If these factors are “common,” this means $\gamma_r := \gamma_r^{(1)} = \gamma_r^{(2)}$, then only the factor loadings $\beta_{ri}^{(p)}$ may vary across samples. This situation may be seen as a functional generalization of the concept of “common principal components” as introduced by Flury (1988) in multivariate analysis. A weaker hypothesis may only require equality of the eigenspaces spanned by the first $L \in \mathbb{N}$ functional principal components. [\mathbb{N} denotes the set of all natural numbers $1, 2, \dots$ ($0 \notin \mathbb{N}$)]. If for both samples the common L -dimensional eigenspaces suffice to approximate the functions with high accuracy, then the distributions in function space are well represented by a low-dimensional factor model, and subsequent analysis may rely on comparing the multivariate distributions of the random vectors $(\beta_{r1}^{(p)}, \dots, \beta_{rL}^{(p)})^\top$.

The idea of “common functional principal components” is of considerable importance in implied volatility (IV) dynamics. This application is discussed in detail in Section 4. Implied volatility is obtained from the pricing model proposed by Black and Scholes (1973) and is a key parameter for quoting options prices. Our aim is to construct low-dimensional factor models for the log-returns of the IV functions of options with different maturities. In our application the first group of functional observations— $\{X_i^{(1)}\}_{i=1}^{n_1}$, are log-returns on the maturity “1 month” (1M group) and second group— $\{X_i^{(2)}\}_{i=1}^{n_2}$, are log-returns on the maturity “3 months” (3M group).

The first three eigenfunctions (ordered with respect to the corresponding eigenvalues), estimated by the method described in Section 2.1, are plotted in Figure 1. The estimated eigenfunctions for both groups are of similar structure, which motivates a common FPCA approach. Based on discretized vectors of functional values, a (multivariate) common principal components analysis of implied volatilities has already been considered by Fengler, Härdle and Villa (2003). They rely on the methodology introduced by Flury (1988) which is based on maximum likelihood estimation under the assumption of multivariate normality. Our analysis overcomes the limitations of this approach by providing specific hypothesis tests in a fully functional setup. It will be shown in Section 4 that for both groups $L = 3$ components suffice to explain 98.2% of the variability of the sample functions. An application of the tests developed in Section 3 does not reject the equality of the corresponding eigenspaces.

2. Functional principal components and one sample inference. In this section we will focus on one sample of i.i.d. smooth random functions $X_1, \dots,$

$X_n \in L^2[0, 1]$. We will assume a well-defined mean function $\mu = E(X_i)$, as well as the existence of a continuous covariance function $\sigma(t, s) = E[\{X_i(t) - \mu(t)\}\{X_i(s) - \mu(s)\}]$. Then $E(\|X_i - \mu\|^2) = \int \sigma(t, t) dt < \infty$, and the covariance operator Γ of X_i is given by

$$(\Gamma v)(t) = \int \sigma(t, s)v(s) ds, \quad v \in L^2[0, 1].$$

The Karhunen–Loève decomposition provides a basic tool to describe the distribution of the random functions X_i . With $\lambda_1 \geq \lambda_2 \geq \dots$ and $\gamma_1, \gamma_2, \dots$ denoting eigenvalues and a corresponding complete orthonormal basis of eigenfunctions of Γ , we obtain

$$(1) \quad X_i = \mu + \sum_{r=1}^{\infty} \beta_{ri} \gamma_r, \quad i = 1, \dots, n,$$

where $\beta_{ri} = \langle X_i - \mu, \gamma_r \rangle$ are uncorrelated (scalar) factor loadings with $E(\beta_{ri}) = 0$, $E(\beta_{ri}^2) = \lambda_r$ and $E(\beta_{ri}\beta_{ki}) = 0$ for $r \neq k$. Structure and dynamics of the random functions can be assessed by analyzing the “functional principal components” γ_r , as well as the distribution of the factor loadings.

A discussion of basic properties of (1) can, for example, be found in [Gihman and Skorohod \(1973\)](#). Under our assumptions, the infinite sums in (1) converge with probability 1, and $\sum_{r=1}^{\infty} \lambda_r = E(\|X_i - \mu\|^2) < \infty$. Smoothness of X_i carries over to a corresponding degree of smoothness of $\sigma(t, s)$ and γ_r . If, with probability 1, $X_i(t)$ is twice continuously differentiable, then σ as well as γ_r are also twice continuously differentiable. The particular case of a Gaussian random function X_i implies that the β_{ri} are independent $N(0, \lambda_r)$ -distributed random variables.

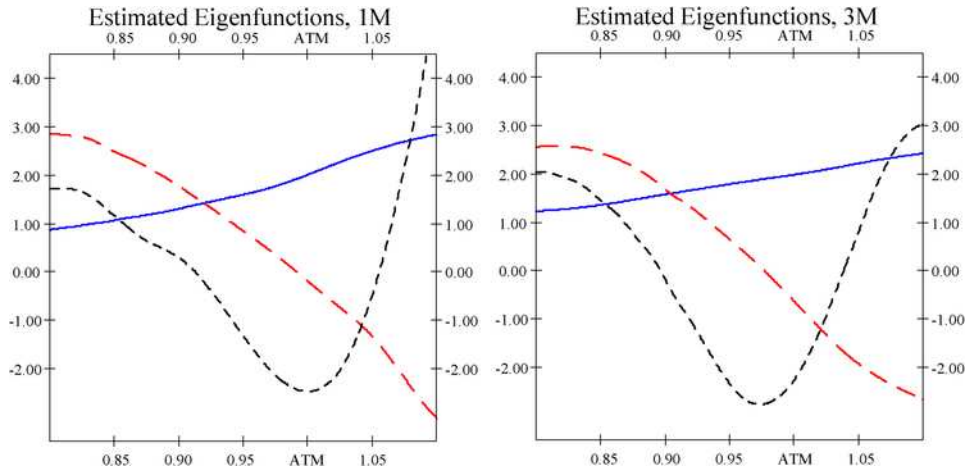


FIG. 1. *Estimated eigenfunctions for 1M group in the left plot and 3M group in the right plot: solid—first function, dashed—second function, finely dashed—third function.*

An important property of (1) consists in the known fact that the first L principal components provide a “best basis” for approximating the sample functions in terms of the integrated square error; see [Ramsay and Silverman \(2005\)](#), Section 6.2.3, among others. For any choice of L orthonormal basis functions v_1, \dots, v_L , the mean integrated square error

$$(2) \quad \rho(v_1, \dots, v_L) = \mathbb{E} \left(\left\| X_i - \mu - \sum_{r=1}^L \langle X_i - \mu, v_r \rangle v_r \right\|^2 \right)$$

is minimized by $v_r = \gamma_r$.

2.1. Estimation of functional principal components. For a given sample an empirical analog of (1) can be constructed by using eigenvalues $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots$ and orthonormal eigenfunctions $\hat{\gamma}_1, \hat{\gamma}_2, \dots$ of the empirical covariance operator $\hat{\Gamma}_n$, where

$$(\hat{\Gamma}_n v)(t) = \int \hat{\sigma}(t, s) v(s) ds,$$

with $\bar{X} = n^{-1} \sum_{i=1}^n X_i$ and $\hat{\sigma}(t, s) = n^{-1} \sum_{i=1}^n \{X_i(t) - \bar{X}(t)\} \{X_i(s) - \bar{X}(s)\}$ denoting sample mean and covariance function. Then

$$(3) \quad X_i = \bar{X} + \sum_{r=1}^n \hat{\beta}_{ri} \hat{\gamma}_r, \quad i = 1, \dots, n,$$

where $\hat{\beta}_{ri} = \langle \hat{\gamma}_r, X_i - \bar{X} \rangle$. We necessarily obtain $n^{-1} \sum_i \hat{\beta}_{ri} = 0$, $n^{-1} \sum_i \hat{\beta}_{ri} \hat{\beta}_{si} = 0$ for $r \neq s$, and $n^{-1} \sum_i \hat{\beta}_{ri}^2 = \hat{\lambda}_r$.

Analysis will have to concentrate on the leading principal components explaining the major part of the variance. In the following we will assume that $\lambda_1 > \lambda_2 > \dots > \lambda_{r_0} > \lambda_{r_0+1}$, where r_0 denotes the maximal number of components to be considered. For all $r = 1, \dots, r_0$, the corresponding eigenfunction γ_r is then uniquely defined up to sign. Signs are arbitrary, decompositions (1) or (3) may just as well be written in terms of $-\gamma_r, -\beta_{ri}$ or $-\hat{\gamma}_r, -\hat{\beta}_{ri}$, and any suitable standardization may be applied by the statistician. In order to ensure that $\hat{\gamma}_r$ may be viewed as an estimator of γ_r rather than of $-\gamma_r$, we will in the following only assume that signs are such that $\langle \gamma_r, \hat{\gamma}_r \rangle \geq 0$. More generally, any subsequent statement concerning differences of two eigenfunctions will be based on the condition of a nonnegative inner product. This does not impose any restriction and will go without saying.

The results of [Dauxois, Pousse and Romain \(1982\)](#) imply that, under regularity conditions, $\|\hat{\gamma}_r - \gamma_r\| = \mathcal{O}_p(n^{-1/2})$, $|\hat{\lambda}_r - \lambda_r| = \mathcal{O}_p(n^{-1/2})$, as well as $|\hat{\beta}_{ri} - \beta_{ri}| = \mathcal{O}_p(n^{-1/2})$ for all $r \leq r_0$.

However, in practice, the sample functions X_i are often not directly observed, but have to be reconstructed from noisy observations Y_{ij} at discrete

design points t_{ik} :

$$(4) \quad Y_{ik} = X_i(t_{ik}) + \varepsilon_{ik}, \quad k = 1, \dots, T_i,$$

where ε_{ik} are independent noise terms with $E(\varepsilon_{ik}) = 0$, $\text{Var}(\varepsilon_{ik}) = \sigma_i^2$.

Our approach for estimating principal components is motivated by the well-known duality relation between row and column spaces of a data matrix; see [Härdle and Simar \(2003\)](#), Chapter 8, among others. In a first step this approach relies on estimating the elements of the matrix:

$$(5) \quad M_{lk} = \langle X_l - \bar{X}, X_k - \bar{X} \rangle, \quad l, k = 1, \dots, n.$$

Some simple linear algebra shows that all nonzero eigenvalues $\hat{\lambda}_1 \geq \hat{\lambda}_2 \cdots$ of $\hat{\Gamma}_n$ and $l_1 \geq l_2 \cdots$ of M are related by $\hat{\lambda}_r = l_r/n$, $r = 1, 2, \dots$. When using the corresponding orthonormal eigenvectors p_1, p_2, \dots of M , the empirical scores $\hat{\beta}_{ri}$, as well as the empirical eigenfunctions $\hat{\gamma}_r$, are obtained by $\hat{\beta}_{ri} = \sqrt{l_r} p_{ir}$ and

$$(6) \quad \hat{\gamma}_r = \frac{1}{\sqrt{l_r}} \sum_{i=1}^n p_{ir} (X_i - \bar{X}) = \frac{1}{\sqrt{l_r}} \sum_{i=1}^n p_{ir} X_i.$$

The elements of M are functionals which can be estimated with asymptotically negligible bias and a parametric rate of convergence $T_i^{-1/2}$. If the data in (4) is generated from a balanced, equidistant design, then it is easily seen that for $i \neq j$ this rate of convergence is achieved by the estimator

$$\widehat{M}_{ij} = T^{-1} \sum_{k=1}^T (Y_{ik} - \bar{Y}_{\cdot k})(Y_{jk} - \bar{Y}_{\cdot k}), \quad i \neq j,$$

and

$$\widehat{M}_{ii} = T^{-1} \sum_{k=1}^T (Y_{ik} - \bar{Y}_{\cdot k})^2 - \hat{\sigma}_i^2,$$

where $\hat{\sigma}_i^2$ denotes some nonparametric estimator of variance and $\bar{Y}_{\cdot k} = n^{-1} \times \sum_{j=1}^n Y_{jk}$.

In the case of a random design some adjustment is necessary: Define the ordered sample $t_{i(1)} \leq t_{i(2)} \leq \dots \leq t_{i(T_i)}$ of design points, and for $j = 1, \dots, T_i$, let $Y_{i(j)}$ denote the observation belonging to $t_{i(j)}$. With $t_{i(0)} = -t_{i(1)}$ and $t_{i(T_i+1)} = 2 - t_{i(T_i)}$, set

$$\chi_i(t) = \sum_{j=1}^{T_i} Y_{i(j)} I \left(t \in \left[\frac{t_{i(j-1)} + t_{i(j)}}{2}, \frac{t_{i(j)} + t_{i(j+1)}}{2} \right) \right), \quad t \in [0, 1],$$

where $I(\cdot)$ denotes the indicator function, and for $i \neq j$, define the estimate of M_{ij} by

$$\widehat{M}_{ij} = \int_0^1 \{\chi_i(t) - \bar{\chi}(t)\} \{\chi_j(t) - \bar{\chi}(t)\} dt,$$

where $\bar{\chi}(t) = n^{-1} \sum_{i=1}^n \chi_i(t)$. Finally, by redefining $t_{i(1)} = -t_{i(2)}$ and $t_{i(T_i+1)} = 2 - t_{i(T_i)}$, set $\chi_i^*(t) = \sum_{j=2}^{T_i} Y_{i(j-1)} I(t \in [\frac{t_{i(j-1)}+t_{i(j)}}{2}, \frac{t_{i(j)}+t_{i(j+1)}}{2}))$, $t \in [0, 1]$. Then construct estimators of the diagonal terms M_{ii} by

$$(7) \quad \widehat{M}_{ii} = \int_0^1 \{\chi_i(t) - \bar{\chi}(t)\} \{\chi_i^*(t) - \bar{\chi}(t)\} dt.$$

The aim of using the estimator (7) for the diagonal terms is to avoid the additional bias implied by $E_\varepsilon(Y_{ik}^2) = X_i(t_{ij})^2 + \sigma_i^2$. Here E_ε denotes conditional expectation given t_{ij} , X_i . Alternatively, we can construct a bias corrected estimator using some nonparametric estimation of variance σ_i^2 , for example, the difference based model-free variance estimators studied in Hall, Kay and Titterton (1990) can be employed.

The eigenvalues $\hat{l}_1 \geq \hat{l}_2 \cdots$ and eigenvectors $\hat{p}_1, \hat{p}_2, \dots$ of the resulting matrix \widehat{M} then provide estimates $\hat{\lambda}_{r;T} = \hat{l}_r/n$ and $\hat{\beta}_{ri;T} = \sqrt{\hat{l}_r} \hat{p}_{ir}$ of λ_r and β_{ri} . Estimates $\hat{\gamma}_{r;T}$ of the empirical functional principal component $\hat{\gamma}_r$ can be determined from (6) when replacing the unknown true functions X_i by nonparametric estimates \hat{X}_i (as, for example, local polynomial estimates) with smoothing parameter (bandwidth) b :

$$(8) \quad \hat{\gamma}_{r;T} = \frac{1}{\sqrt{\hat{l}_r}} \sum_{i=1}^n \hat{p}_{ir} \hat{X}_i.$$

When considering (8), it is important to note that $\hat{\gamma}_{r;T}$ is defined as a *weighted average* of all estimated sample functions. Averaging reduces variance, and efficient estimation of $\hat{\gamma}_r$ therefore requires *undersmoothing* of individual function estimates \hat{X}_i . Theoretical results are given in Theorem 1 below. Indeed, if, for example, n and $T = \min_i T_i$ are of the same order of magnitude, then under suitable additional regularity conditions it will be shown that for an optimal choice of a smoothing parameter $b \sim (nT)^{-1/5}$ and twice continuously differentiable X_i , we obtain the rate of convergence $\|\hat{\gamma}_r - \hat{\gamma}_{r;T}\| = \mathcal{O}_p\{(nT)^{-2/5}\}$. Note, however, that the bias corrected estimator (7) may yield negative eigenvalues. In practice, these values will be small and will have to be interpreted as zero. Furthermore, the eigenfunctions determined by (8) may not be exactly orthogonal. Again, when using reasonable bandwidths, this effect will be small, but of course (8) may be followed by a suitable orthogonalization procedure.

It is of interest to compare our procedure to more standard methods for estimating λ_r and $\hat{\gamma}_r$ as mentioned above. When evaluating eigenvalues and eigenfunctions of the empirical covariance operator of nonparametrically estimated curves \hat{X}_i , then for fixed $r \leq r_0$ the above rate of convergence for the estimated eigenfunctions may well be achieved for a suitable choice of smoothing parameters (e.g., number of basis functions). But as will be seen

from Theorem 1, our approach also implies that $|\hat{\lambda}_r - \frac{\hat{t}_r}{n}| = \mathcal{O}_p(T^{-1} + n^{-1})$. When using standard methods it does not seem to be possible to obtain a corresponding rate of convergence, since any smoothing bias $|\mathbb{E}[\hat{X}_i(t)] - X_i(t)|$ will invariably affect the quality of the corresponding estimate of $\hat{\lambda}_r$.

We want to emphasize that any finite sample interpretation will require that T is sufficiently large such that our nonparametric reconstructions of individual curves can be assumed to possess a fairly small bias. The above arguments do not apply to extremely sparse designs with very few observations per curve [see Hall, Müller and Wang (2006) for an FPCA methodology focusing on sparse data].

Note that, in addition to (8), our final estimate of the empirical mean function $\hat{\mu} = \bar{X}$ will be given by $\hat{\mu}_T = n^{-1} \sum_i \hat{X}_i$. A straightforward approach to determine a suitable bandwidth b consists in a “leave-one-individual-out” cross-validation. For the maximal number r_0 of components to be considered, let $\hat{\mu}_{T,-i}$ and $\hat{\gamma}_{r;T,-i}$, $r = 1, \dots, r_0$, denote the estimates of $\hat{\mu}$ and $\hat{\gamma}_r$ obtained from the data (Y_{lj}, t_{lj}) , $l = 1, \dots, i-1, i+1, \dots, n$, $j = 1, \dots, T_k$. By (8), these estimates depend on b , and one may approximate an optimal smoothing parameter by minimizing

$$\sum_i \sum_j \left\{ Y_{ij} - \hat{\mu}_{T,-i}(t_{ij}) - \sum_{r=1}^{r_0} \hat{\vartheta}_{ri} \hat{\gamma}_{r;T,-i}(t_{ij}) \right\}^2$$

over b , where $\hat{\vartheta}_{ri}$ denote ordinary least squares estimates of $\hat{\beta}_{ri}$. A more sophisticated version of this method may even allow to select different bandwidths b_r when estimating different functional principal components by (8). Although, under certain regularity conditions, the same qualitative rates of convergence hold for any arbitrary *fixed* $r \leq r_0$, the quality of estimates decreases when r becomes large. Due to $\langle \gamma_s, \gamma_r \rangle = 0$ for $s < r$, the number of zero crossings, peaks and valleys of γ_r has to increase with r . Hence, in tendency γ_r will be less and less smooth as r increases. At the same time, $\lambda_r \rightarrow 0$, which means that for large r the r th eigenfunctions will only possess a very small influence on the structure of X_i . This in turn means that the relative importance of the error terms ε_{ik} in (4) on the structure of $\hat{\gamma}_{r;T}$ will increase with r .

2.2. One sample inference. Clearly, in the framework described by (1)–(4) we are faced with two sources of variability of estimated functional principal components. Due to sampling variation, $\hat{\gamma}_r$ will differ from the true component γ_r , and due to (4), there will exist an additional estimation error when approximating $\hat{\gamma}_r$ by $\hat{\gamma}_{r;T}$.

The following theorems quantify the order of magnitude of these different types of error. Our theoretical results are based on the following assumptions on the structure of the random functions X_i .

ASSUMPTION 1. $X_1, \dots, X_n \in L^2[0, 1]$ is an i.i.d. sample of random functions with mean μ and continuous covariance function $\sigma(t, s)$, and (1) holds for a system of eigenfunctions satisfying $\sup_{s \in \mathbb{N}} \sup_{t \in [0, 1]} |\gamma_s(t)| < \infty$. Furthermore, $\sum_{r=1}^{\infty} \sum_{s=1}^{\infty} \mathbb{E}[\beta_{ri}^2 \beta_{si}^2] < \infty$ and $\sum_{q=1}^{\infty} \sum_{s=1}^{\infty} \mathbb{E}[\beta_{ri}^2 \beta_{qi} \beta_{si}] < \infty$ for all $r \in \mathbb{N}$.

Recall that $\mathbb{E}[\beta_{ri}] = 0$ and $\mathbb{E}[\beta_{ri} \beta_{si}] = 0$ for $r \neq s$. Note that the assumption on the factor loadings is necessarily fulfilled if X_i are Gaussian random functions. Then β_{ri} and β_{si} are independent for $r \neq s$, all moments of β_{ri} are finite, and hence $\mathbb{E}[\beta_{ri}^2 \beta_{qi} \beta_{si}] = 0$ for $q \neq s$, as well as $\mathbb{E}[\beta_{ri}^2 \beta_{si}^2] = \lambda_r \lambda_s$ for $r \neq s$; see [Gihman and Skorohod \(1973\)](#).

We need some further assumptions concerning smoothness of X_i and the structure of the discrete model (4).

ASSUMPTION 2. (a) X_i is a.s. twice continuously differentiable. There exists a constant $D_1 < \infty$ such that the derivatives are bounded by $\sup_t \mathbb{E}[X_i'(t)^4] \leq D_1$, as well as $\sup_t \mathbb{E}[X_i''(t)^4] \leq D_1$.

(b) The design points t_{ik} , $i = 1, \dots, n$, $k = 1, \dots, T_i$, are i.i.d. random variables which are independent of X_i and ε_{ik} . The corresponding design density f is continuous on $[0, 1]$ and satisfies $\inf_{t \in [0, 1]} f(t) > 0$.

(c) For any i , the error terms ε_{ik} are i.i.d. zero mean random variables with $\text{Var}(\varepsilon_{ik}) = \sigma_i^2$. Furthermore, ε_{ik} is independent of X_i , and there exists a constant D_2 such that $\mathbb{E}(\varepsilon_{ik}^8) < D_2$ for all i, k .

(d) The estimates \hat{X}_i used in (8) are determined by either a local linear or a Nadaraya–Watson kernel estimator with smoothing parameter b and kernel function K . K is a continuous probability density which is symmetric at 0.

The following theorems provide asymptotic results as $n, T \rightarrow \infty$, where $T = \min_{i=1}^n \{T_i\}$.

THEOREM 1. *In addition to Assumptions 1 and 2, assume that $\inf_{s \neq r} |\lambda_r - \lambda_s| > 0$ holds for some $r = 1, 2, \dots$. Then we have the following:*

(i) $n^{-1} \sum_{i=1}^n (\hat{\beta}_{ri} - \hat{\beta}_{ri;T})^2 = \mathcal{O}_p(T^{-1})$ and

$$(9) \quad \left| \hat{\lambda}_r - \frac{\hat{l}_r}{n} \right| = \mathcal{O}_p(T^{-1} + n^{-1}).$$

(ii) *If additionally $b \rightarrow 0$ and $(Tb)^{-1} \rightarrow 0$ as $n, T \rightarrow \infty$, then for all $t \in [0, 1]$,*

$$(10) \quad |\hat{\gamma}_r(t) - \hat{\gamma}_{r;T}(t)| = \mathcal{O}_p\{b^2 + (nTb)^{-1/2} + (Tb^{1/2})^{-1} + n^{-1}\}.$$

A proof is given in the [Appendix](#).

THEOREM 2. Under Assumption 1 we obtain the following:

(i) For all $t \in [0, 1]$,

$$\sqrt{n}\{\bar{X}(t) - \mu(t)\} = \sum_r \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^n \beta_{ri} \right\} \gamma_r(t) \xrightarrow{\mathcal{L}} N\left(0, \sum_r \lambda_r \gamma_r(t)^2\right).$$

If, furthermore, $\lambda_{r-1} > \lambda_r > \lambda_{r+1}$ holds for some fixed $r \in \{1, 2, \dots\}$, then

(ii)

$$(11) \quad \sqrt{n}(\hat{\lambda}_r - \lambda_r) = \frac{1}{\sqrt{n}} \sum_{i=1}^n (\beta_{ri}^2 - \lambda_r) + \mathcal{O}_p(n^{-1/2}) \xrightarrow{\mathcal{L}} N(0, \Lambda_r),$$

where $\Lambda_r = \mathbb{E}[(\beta_{ri}^2 - \lambda_r)^2]$,

(iii) and for all $t \in [0, 1]$

$$(12) \quad \hat{\gamma}_r(t) - \gamma_r(t) = \sum_{s \neq r} \left\{ \frac{1}{n(\lambda_r - \lambda_s)} \sum_{i=1}^n \beta_{si} \beta_{ri} \right\} \gamma_s(t) + R_r(t),$$

where $\|R_r\| = \mathcal{O}_p(n^{-1})$.

Moreover,

$$\begin{aligned} & \sqrt{n} \sum_{s \neq r} \left\{ \frac{1}{n(\lambda_r - \lambda_s)} \sum_{i=1}^n \beta_{si} \beta_{ri} \right\} \gamma_s(t) \\ & \xrightarrow{\mathcal{L}} N\left(0, \sum_{q \neq r} \sum_{s \neq r} \frac{\mathbb{E}[\beta_{ri}^2 \beta_{qi} \beta_{si}]}{(\lambda_q - \lambda_r)(\lambda_s - \lambda_r)} \gamma_q(t) \gamma_s(t)\right). \end{aligned}$$

A proof can be found in the [Appendix](#). The theorem provides a generalization of the results of [Dauxois, Pousse and Romain \(1982\)](#) who derive explicit asymptotic distributions by assuming Gaussian random functions X_i . Note that in this case $\Lambda_r = 2\lambda_r^2$ and $\sum_{q \neq r} \sum_{s \neq r} \frac{\mathbb{E}[\beta_{ri}^2 \beta_{qi} \beta_{si}]}{(\lambda_q - \lambda_r)(\lambda_s - \lambda_r)} \gamma_q(t) \gamma_s(t) = \sum_{s \neq r} \frac{\lambda_r \lambda_s}{(\lambda_s - \lambda_r)^2} \gamma_s(t)^2$.

When evaluating the bandwidth-dependent terms in (10), best rates of convergence $|\hat{\gamma}_r(t) - \hat{\gamma}_{r;T}(t)| = \mathcal{O}_p\{(nT)^{-2/5} + T^{-4/5} + n^{-1}\}$ are achieved when choosing an undersmoothing bandwidth $b \sim \max\{(nT)^{-1/5}, T^{-2/5}\}$. Theoretical work in functional data analysis is usually based on the implicit assumption that the additional error due to (4) is negligible, and that one can proceed “as if” the functions X_i were directly observed. In view of Theorems 1 and 2, this approach is justified in the following situations:

(1) T is much larger than n , that is, $n/T^{4/5} \rightarrow 0$, and the smoothing parameter b in (8) is of order $T^{-1/5}$ (optimal smoothing of individual functions).

(2) An undersmoothing bandwidth $b \sim \max\{(nT)^{-1/5}, T^{-2/5}\}$ is used and $n/T^{8/5} \rightarrow 0$. This means that T may be smaller than n , but T must be at least of order of magnitude larger than $n^{5/8}$.

In both cases (1) and (2) the above theorems imply that $|\hat{\lambda}_r - \frac{\hat{L}_r}{n}| = \mathcal{O}_p(|\hat{\lambda}_r - \lambda_r|)$, as well as $\|\hat{\gamma}_r - \hat{\gamma}_{r;T}\| = \mathcal{O}_p(\|\hat{\gamma}_r - \gamma_r\|)$. Inference about functional principal components will then be first-order equivalent to an inference based on known functions X_i .

In such situations Theorem 2 suggests bootstrap procedures as tools for one sample inference. For example, the distribution of $\|\hat{\gamma}_r - \gamma_r\|$ may be approximated by the bootstrap distribution of $\|\hat{\gamma}_r^* - \hat{\gamma}_r\|$, where $\hat{\gamma}_r^*$ are estimates to be obtained from i.i.d. bootstrap resamples $X_1^*, X_2^*, \dots, X_n^*$ of $\{X_1, X_2, \dots, X_n\}$. This means that $X_1^* = X_{i_1}, \dots, X_n^* = X_{i_n}$ for some indices i_1, \dots, i_n drawn independently and with replacement from $\{1, \dots, n\}$ and, in practice, $\hat{\gamma}_r^*$ may thus be approximated from corresponding discrete data $(Y_{i_1 j}, t_{i_1 j})_{j=1, \dots, T_{i_1}}, \dots, (Y_{i_n j}, t_{i_n j})_{j=1, \dots, T_{i_n}}$. The additional error is negligible if either (1) or (2) is satisfied.

One may wonder about the validity of such a bootstrap. Functions are complex objects and there is no established result in bootstrap theory which readily generalizes to samples of random functions. But by (1), i.i.d. bootstrap resamples $\{X_i^*\}_{i=1, \dots, n}$ may be equivalently represented by corresponding, i.i.d. resamples $\{\beta_{1i}^*, \beta_{2i}^*, \dots\}_{i=1, \dots, n}$ of factor loadings. Standard multivariate bootstrap theorems imply that for any $q \in \mathbb{N}$ the distribution of moments of the random vectors $(\beta_{1i}, \dots, \beta_{qi})$ may be consistently approximated by the bootstrap distribution of corresponding moments of $(\beta_{1i}^*, \dots, \beta_{qi}^*)$. Together with some straightforward limit arguments as $q \rightarrow \infty$, the structure of the first-order terms in the asymptotic expansions (11) and (12) then allows to establish consistency of the functional bootstrap. These arguments will be made precise in the proof of Theorem 3 below, which concerns related bootstrap statistics in two sample problems.

REMARK. Theorem 2(iii) implies that the variance of $\hat{\gamma}_r$ is large if one of the differences $\lambda_{r-1} - \lambda_r$ or $\lambda_r - \lambda_{r+1}$ is small. In the limit case of eigenvalues of multiplicity $m > 1$ our theory does not apply. Note that then only the m -dimensional eigenspace is identified, but not a particular basis (eigenfunctions). In multivariate PCA Tyler (1981) provides some inference results on corresponding projection matrices assuming that $\lambda_r > \lambda_{r+1} \geq \dots \geq \lambda_{r+m} > \lambda_{r+m+1}$ for known values of r and m .

Although the existence of eigenvalues λ_r , $r \leq r_0$, with multiplicity $m > 1$ may be considered as a degenerate case, it is immediately seen that $\lambda_r \rightarrow 0$ and, hence, $\lambda_r - \lambda_{r+1} \rightarrow 0$ as r increases. Even in the case of fully observed

functions X_i , estimates of eigenfunctions corresponding to very small eigenvalues will thus be poor. The problem of determining a sensible upper limit of the number r_0 of principal components to be analyzed is addressed in [Hall and Hosseini-Nasab \(2006\)](#).

3. Two sample inference. The comparison of functional components across groups leads naturally to two sample problems. Thus, let

$$X_1^{(1)}, X_2^{(1)}, \dots, X_{n_1}^{(1)} \quad \text{and} \quad X_1^{(2)}, X_2^{(2)}, \dots, X_{n_2}^{(2)}$$

denote two independent samples of smooth functions. The problem of interest is to test in how far the distributions of these random functions coincide. The structure of the different distributions in function space can be accessed by means of the respective Karhunen–Loève decompositions. The problem to be considered then translates into testing equality of the different components of these decompositions given by

$$(13) \quad X_i^{(p)} = \mu^{(p)} + \sum_{r=1}^{\infty} \beta_{ri}^{(p)} \gamma_r^{(p)}, \quad p = 1, 2,$$

where again $\gamma_r^{(p)}$ are the eigenfunctions of the respective covariance operator $\Gamma^{(p)}$ corresponding to the eigenvalues $\lambda_1^{(p)} = \text{E}\{(\beta_{1i}^{(p)})^2\} \geq \lambda_2^{(p)} = \text{E}\{(\beta_{2i}^{(p)})^2\} \geq \dots$. We will again suppose that $\lambda_{r-1}^{(p)} > \lambda_r^{(p)} > \lambda_{r+1}^{(p)}$, $p = 1, 2$, for all $r \leq r_0$ components to be considered. Without restriction, we will additionally assume that signs are such that $\langle \gamma_r^{(1)}, \gamma_r^{(2)} \rangle \geq 0$, as well as $\langle \hat{\gamma}_r^{(1)}, \hat{\gamma}_r^{(2)} \rangle \geq 0$.

It is of great interest to detect possible variations in the functional components characterizing the two samples in (13). Significant difference may give rise to substantial interpretation. Important hypotheses to be considered thus are as follows:

$$H_{01} : \mu^{(1)} = \mu^{(2)} \quad \text{and} \quad H_{02,r} : \gamma_r^{(1)} = \gamma_r^{(2)}, \quad r \leq r_0.$$

Hypothesis $H_{02,r}$ is of particular importance. Then $\gamma_r^{(1)} = \gamma_r^{(2)}$ and only the factor loadings β_{ri} may vary across samples. If, for example, $H_{02,r}$ is accepted, one may additionally want to test hypotheses about the distributions of $\beta_{ri}^{(p)}$, $p = 1, 2$. Recall that necessarily $\text{E}\{\beta_{ri}^{(p)}\} = 0$, $\text{E}\{\beta_{ri}^{(p)}\}^2 = \lambda_r^{(p)}$, and $\beta_{si}^{(p)}$ is uncorrelated with $\beta_{ri}^{(p)}$ if $r \neq s$. If the $X_i^{(p)}$ are Gaussian random variables, the $\beta_{ri}^{(p)}$ are independent $N(0, \lambda_r)$ random variables. A natural hypothesis to be tested then refers to the equality of variances:

$$H_{03,r} : \lambda_r^{(1)} = \lambda_r^{(2)}, \quad r = 1, 2, \dots$$

Let $\hat{\mu}^{(p)}(t) = \frac{1}{n_p} \sum_i X_i^{(p)}(t)$, and let $\hat{\lambda}_1^{(p)} \geq \hat{\lambda}_2^{(p)} \geq \dots$ and $\hat{\gamma}_1^{(p)}, \hat{\gamma}_2^{(p)}, \dots$ denote eigenvalues and corresponding eigenfunctions of the empirical covariance operator $\hat{\Gamma}_{n_p}^{(p)}$ of $X_1^{(p)}, X_2^{(p)}(t), \dots, X_{n_p}^{(p)}$. The following test statistics are

defined in terms of $\hat{\mu}^{(p)}$, $\hat{\lambda}_r^{(p)}$ and $\hat{\gamma}_r^{(p)}$. As discussed in the proceeding section, all curves in both samples are usually not directly observed, but have to be reconstructed from noisy observations according to (4). In this situation, the “true” empirical eigenvalues and eigenfunctions have to be replaced by their discrete sample estimates. Bootstrap estimates are obtained by resampling the observations corresponding to the unknown curves $X_i^{(p)}$. As discussed in Section 2.2, the validity of our test procedures is then based on the assumption that T is sufficiently large such that the additional estimation error is asymptotically negligible.

Our tests of the hypotheses H_{0_1} , $H_{0_{2,r}}$ and $H_{0_{3,r}}$ rely on the statistics

$$\begin{aligned} D_1 &\stackrel{\text{def}}{=} \|\hat{\mu}^{(1)} - \hat{\mu}^{(2)}\|^2, \\ D_{2,r} &\stackrel{\text{def}}{=} \|\hat{\gamma}_r^{(1)} - \hat{\gamma}_r^{(2)}\|^2, \\ D_{3,r} &\stackrel{\text{def}}{=} |\hat{\lambda}_r^{(1)} - \hat{\lambda}_r^{(2)}|^2. \end{aligned}$$

The respective null-hypothesis has to be rejected if $D_1 \geq \Delta_{1;1-\alpha}$, $D_{2,r} \geq \Delta_{2,r;1-\alpha}$ or $D_{3,r} \geq \Delta_{3,r;1-\alpha}$, where $\Delta_{1;1-\alpha}$, $\Delta_{2,r;1-\alpha}$ and $\Delta_{3,r;1-\alpha}$ denote the critical values of the distributions of

$$\begin{aligned} \Delta_1 &\stackrel{\text{def}}{=} \|\hat{\mu}^{(1)} - \mu^{(1)} - (\hat{\mu}^{(2)} - \mu^{(2)})\|^2, \\ \Delta_{2,r} &\stackrel{\text{def}}{=} \|\hat{\gamma}_r^{(1)} - \gamma_r^{(1)} - (\hat{\gamma}_r^{(2)} - \gamma_r^{(2)})\|^2, \\ \Delta_{3,r} &\stackrel{\text{def}}{=} |\hat{\lambda}_r^{(1)} - \lambda_r^{(1)} - (\hat{\lambda}_r^{(2)} - \lambda_r^{(2)})|^2. \end{aligned}$$

Of course, the distributions of the different Δ 's cannot be accessed directly, since they depend on the unknown true population mean, eigenvalues and eigenfunctions. However, it will be shown below that these distributions and, hence, their critical values are approximated by the bootstrap distribution of

$$\begin{aligned} \Delta_1^* &\stackrel{\text{def}}{=} \|\hat{\mu}^{(1)*} - \hat{\mu}^{(1)} - (\hat{\mu}^{(2)*} - \hat{\mu}^{(2)})\|^2, \\ \Delta_{2,r}^* &\stackrel{\text{def}}{=} \|\hat{\gamma}_r^{(1)*} - \hat{\gamma}_r^{(1)} - (\hat{\gamma}_r^{(2)*} - \hat{\gamma}_r^{(2)})\|^2, \\ \Delta_{3,r}^* &\stackrel{\text{def}}{=} |\hat{\lambda}_r^{(1)*} - \hat{\lambda}_r^{(1)} - (\hat{\lambda}_r^{(2)*} - \hat{\lambda}_r^{(2)})|^2, \end{aligned}$$

where $\hat{\mu}^{(1)*}$, $\hat{\gamma}_r^{(1)*}$, $\hat{\lambda}_r^{(1)*}$, as well as $\hat{\mu}^{(2)*}$, $\hat{\gamma}_r^{(2)*}$, $\hat{\lambda}_r^{(2)*}$, are estimates to be obtained from independent bootstrap samples $X_1^{1*}(t)$, $X_2^{1*}(t)$, \dots , $X_{n_1}^{1*}(t)$, as well as $X_1^{2*}(t)$, $X_2^{2*}(t)$, \dots , $X_{n_2}^{2*}(t)$.

This test procedure is motivated by the following insights:

(1) Under each of our null-hypotheses the respective test statistics D is equal to the corresponding Δ . The test will thus asymptotically possess the correct level: $P(D > \Delta_{1-\alpha}) \approx \alpha$.

(2) If the null hypothesis is false, then $D \neq \Delta$. Compared to the distribution of Δ , the distribution of D is shifted by the difference in the true means, eigenfunctions or eigenvalues. In tendency D will be larger than $\Delta_{1-\alpha}$.

Let $1 < L \leq r_0$. Even if for $r \leq L$ the equality of eigenfunctions is rejected, we may be interested in the question of whether at least the L -dimensional eigenspaces generated by the first L eigenfunctions are identical. Therefore, let $\mathcal{E}_L^{(1)}$, as well as $\mathcal{E}_L^{(2)}$, denote the L -dimensional linear function spaces generated by the eigenfunctions $\gamma_1^{(1)}, \dots, \gamma_L^{(1)}$ and $\gamma_1^{(2)}, \dots, \gamma_L^{(2)}$, respectively. We then aim to test the null hypothesis:

$$H_{0_{4,L}} : \mathcal{E}_L^{(1)} = \mathcal{E}_L^{(2)}.$$

Of course, $H_{0_{4,L}}$ corresponds to the hypothesis that the operators projecting into $\mathcal{E}_L^{(1)}$ and $\mathcal{E}_L^{(2)}$ are identical. This in turn translates into the condition that

$$\sum_{r=1}^L \gamma_r^{(1)}(t)\gamma_r^{(1)}(s) = \sum_{r=1}^L \gamma_r^{(2)}(t)\gamma_r^{(2)}(s) \quad \text{for all } t, s \in [0, 1].$$

Similar to above, a suitable test statistic is given by

$$D_{4,L} \stackrel{\text{def}}{=} \iint \left\{ \sum_{r=1}^L \hat{\gamma}_r^{(1)}(t)\hat{\gamma}_r^{(1)}(s) - \sum_{r=1}^L \hat{\gamma}_r^{(2)}(t)\hat{\gamma}_r^{(2)}(s) \right\}^2 dt ds$$

and the null hypothesis is rejected if $D_{4,L} \geq \Delta_{4,L;1-\alpha}$, where $\Delta_{4,L;1-\alpha}$ denotes the critical value of the distribution of

$$\Delta_{4,L} \stackrel{\text{def}}{=} \iint \left[\sum_{r=1}^L \{ \hat{\gamma}_r^{(1)}(t)\hat{\gamma}_r^{(1)}(s) - \gamma_r^{(1)}(t)\gamma_r^{(1)}(s) \} - \sum_{r=1}^L \{ \hat{\gamma}_r^{(2)}(t)\hat{\gamma}_r^{(2)}(s) - \gamma_r^{(2)}(t)\gamma_r^{(2)}(s) \} \right]^2 dt ds.$$

The distribution of $\Delta_{4,L}$ and, hence, its critical values are approximated by the bootstrap distribution of

$$\Delta_{4,L}^* \stackrel{\text{def}}{=} \iint \left[\sum_{r=1}^L \{ \hat{\gamma}_r^{(1)*}(t)\hat{\gamma}_r^{(1)*}(s) - \hat{\gamma}_r^{(1)}(t)\hat{\gamma}_r^{(1)}(s) \} - \sum_{r=1}^L \{ \hat{\gamma}_r^{(2)*}(t)\hat{\gamma}_r^{(2)*}(s) - \hat{\gamma}_r^{(2)}(t)\hat{\gamma}_r^{(2)}(s) \} \right]^2 dt ds.$$

It will be shown in Theorem 3 below that under the null hypothesis, as well as under the alternative, the distributions of $n\Delta_1, n\Delta_{2,r}, n\Delta_{3,r}, n\Delta_{4,L}$ converge to continuous limit distributions which can be consistently approximated by the bootstrap distributions of $n\Delta_1^*, n\Delta_{2,r}^*, n\Delta_{3,r}^*, n\Delta_{4,L}^*$.

3.1. *Theoretical results.* Let $n = (n_1 + n_2)/2$. We will assume that asymptotically $n_1 = n \cdot q_1$ and $n_2 = n \cdot q_2$ for some fixed proportions q_1 and q_2 . We will then study the asymptotic behavior of our statistics as $n \rightarrow \infty$.

We will use $\mathcal{X}_1 = \{X_1^{(1)}, \dots, X_{n_1}^{(1)}\}$ and $\mathcal{X}_2 = \{X_1^{(2)}, \dots, X_{n_2}^{(2)}\}$ to denote the observed samples of random functions.

THEOREM 3. *Assume that $\{X_1^{(1)}, \dots, X_{n_1}^{(1)}\}$ and $\{X_1^{(2)}, \dots, X_{n_2}^{(2)}\}$ are two independent samples of random functions, each of which satisfies Assumption 1. As $n \rightarrow \infty$ we then obtain the following:*

(i) *There exists a nondegenerated, continuous probability distribution F_1 such that $n\Delta_1 \xrightarrow{\mathcal{L}} F_1$, and for any $\delta > 0$,*

$$|P(n\Delta_1 \geq \delta) - P(n\Delta_1^* \geq \delta | \mathcal{X}_1, \mathcal{X}_2)| = \mathcal{O}_p(1).$$

(ii) *If, furthermore, $\lambda_{r-1}^{(1)} > \lambda_r^{(1)} > \lambda_{r+1}^{(1)}$ and $\lambda_{r-1}^{(2)} > \lambda_r^{(2)} > \lambda_{r+1}^{(2)}$ hold for some fixed $r = 1, 2, \dots$, there exist a nondegenerated, continuous probability distributions $F_{k,r}$ such that $n\Delta_{k,r} \xrightarrow{\mathcal{L}} F_{k,r}$, $k = 2, 3$, and for any $\delta > 0$,*

$$|P(n\Delta_{k,r} \geq \delta) - P(n\Delta_{k,r}^* \geq \delta | \mathcal{X}_1, \mathcal{X}_2)| = \mathcal{O}_p(1), \quad k = 2, 3.$$

(iii) *If $\lambda_r^{(1)} > \lambda_{r+1}^{(1)} > 0$ and $\lambda_r^{(2)} > \lambda_{r+1}^{(2)} > 0$ hold for all $r = 1, \dots, L$, there exists a nondegenerated, continuous probability distribution $F_{4,L}$ such that $n\Delta_{4,L} \xrightarrow{\mathcal{L}} F_{4,L}$, and for any $\delta > 0$,*

$$|P(n\Delta_{4,L} \geq \delta) - P(n\Delta_{4,L}^* \geq \delta | \mathcal{X}_1, \mathcal{X}_2)| = \mathcal{O}_p(1).$$

The structures of the distributions $F_1, F_{2,r}, F_{3,r}, F_{4,L}$ are derived in the proof of the theorem which can be found in the [Appendix](#). They are obtained as limits of distributions of quadratic forms.

3.2. *Simulation study.* In this paragraph we illustrate the finite behavior of the proposed test. The basic simulation-setup (setup ‘‘a’’) is established as follows: the first sample is generated by the random combination of orthonormalized sine and cosine functions (Fourier functions) and the second sample is generated by the random combination of the same but shifted factor functions:

$$\begin{aligned} X_i^{(1)}(t_{ik}) &= \beta_{1i}^{(1)} \sqrt{2} \sin(2\pi t_{ik}) + \beta_{2i}^{(1)} \sqrt{2} \cos(2\pi t_{ik}), \\ X_i^{(2)}(t_{ik}) &= \beta_{1i}^{(2)} \sqrt{2} \sin\{2\pi(t_{ik} + \delta)\} + \beta_{2i}^{(2)} \sqrt{2} \cos\{2\pi(t_{ik} + \delta)\}. \end{aligned}$$

The factor loadings are i.i.d. random variables with $\beta_{1i}^{(p)} \sim N(0, \lambda_1^{(p)})$ and $\beta_{2i}^{(p)} \sim N(0, \lambda_2^{(p)})$. The functions are generated on the equidistant grid $t_{ik} = t_k = k/T$, $k = 1, \dots, T = 100$, $i = 1, \dots, n = 70$. The simulation setup is based

TABLE 1

The results of the simulations for $\alpha = 0.1$, $n = 70$, $T = 100$, number of simulations 250

Setup/shift	0	0.05	0.1	0.15	0.2	0.25
(a) 10, 5, 8, 4	0.13	0.41	0.85	0.96	1	1
(a) 4, 2, 2, 1	0.12	0.48	0.87	0.96	1	1
(a) 2, 1, 1.5, 2	0.14	0.372	0.704	0.872	0.92	0.9
(b) 10, 5, 8, 4 D_1	0.10	0.44	0.86	0.95	1	1
(b) 10, 5, 8, 4 D_2	1	1	1	1	1	1

on the fact that the error of the estimation of the eigenfunctions simulated by sine and cosine functions is, in particular, manifested by some shift of the estimated eigenfunctions. The focus of this simulation study is the test of common eigenfunctions.

For the presentation of results in Table 1, we use the following notation: “(a) $\lambda_1^{(1)}, \lambda_2^{(1)}, \lambda_2^{(2)}, \lambda_2^{(2)}$.” The shift parameter δ is changing from 0 to 0.25 with the step 0.05. It should be mentioned that the shift $\delta = 0$ yields the simulation of level and setup with shift $\delta = 0.25$ yields the simulation of the alternative, where the two factor functions are exchanged.

In the second setup (setup “b”) the first factor functions are the same and the second factor functions differ:

$$X_i^{(1)}(t_{ik}) = \beta_{1i}^{(1)} \sqrt{2} \sin(2\pi t_{ik}) + \beta_{2i}^{(1)} \sqrt{2} \cos(2\pi t_{ik}),$$

$$X_i^{(2)}(t_{ik}) = \beta_{1i}^{(2)} \sqrt{2} \sin\{2\pi(t_{ik} + \delta)\} + \beta_{2i}^{(2)} \sqrt{2} \sin\{4\pi(t_{ik} + \delta)\}.$$

In Table 1 we use the notation “(b) $\lambda_1^{(1)}, \lambda_2^{(1)}, \lambda_2^{(2)}, \lambda_2^{(2)}, D_r$.” D_r means the test for the equality of the r th eigenfunction. In the bootstrap tests we used 500 bootstrap replications. The critical level in this simulation is $\alpha = 0.1$. The number of simulations is 250.

We can interpret Table 1 in the following way: In power simulations ($\delta \neq 0$) test behaves as expected: less powerful if the functions are “hardly distinguishable” (small shift, small difference in eigenvalues). The level approximation seems to be less precise if the difference in the eigenvalues ($\lambda_1^{(p)} - \lambda_2^{(p)}$) becomes smaller. This can be explained by relative small sample-size n , small number of bootstrap-replications and increasing estimation-error as argued in Theorem 2, assertion (iii).

In comparison to our general setup (4), we used an equidistant and common design for all functions. This simplification is necessary, it simplifies and speeds-up the simulations, in particular, using general random and observation-specific design makes the simulation computationally untractable.

Second, we omitted the additional observation error, this corresponds to the standard assumptions in the functional principal components theory. As

TABLE 2
The results of the simulation for $\alpha = 0.1$, $n = 70$, $T = 100$ with additional error in observation

Setup/shift	0	0.05	0.1	0.15	0.2	0.25
(a) 10, 5, 8, 4	0.09	0.35	0.64	0.92	0.94	0.97

argued in Section 2.2, the inference based on the directly observed functions and estimated functions X_i is first-order equivalent under mild conditions implied by Theorems 1 and 2. In order to illustrate this theoretical result in the simulation, we used the following setup:

$$X_i^{(1)}(t_{ik}) = \beta_{1i}^{(1)} \sqrt{2} \sin(2\pi t_{ik}) + \beta_{2i}^{(1)} \sqrt{2} \cos(2\pi t_{ik}) + \varepsilon_{ik}^{(1)},$$

$$X_i^{(2)}(t_{ik}) = \beta_{1i}^{(2)} \sqrt{2} \sin\{2\pi(t_{ik} + \delta)\} + \beta_{2i}^{(2)} \sqrt{2} \cos\{2\pi(t_{ik} + \delta)\} + \varepsilon_{ik}^{(2)},$$

where $\varepsilon_{ik}^{(p)} \sim N(0, 0.25)$, $p = 1, 2$, all other parameters remain the same as in the simulation setup “a.” Using this setup, we recalculate the simulation presented in the second “row” of Table 1, for estimation of the functions $X_i^{(p)}$, $p = 1, 2$, we used the Nadaraya–Watson estimation with Epanechnikov kernel and bandwidth $b = 0.05$. We run the simulations with various bandwidths, the choice of the bandwidth does not have a strong influence on results except by oversmoothing (large bandwidths). The results are printed in Table 2. As we can see, the difference of the simulation results using estimated functions is not significant in comparison to the results printed in the second line of Table 1—directly observed functional values.

The last limitation of this simulation study is the choice of a particular alternative. A more general setup of this simulation study might be based on the following model: $X_i^{(1)}(t) = \beta_{1i}^{(1)} \gamma_1^{(1)}(t) + \beta_{2i}^{(1)} \gamma_2^{(1)}(t)$, $X_i^{(2)}(t) = \beta_{1i}^{(2)} \gamma_1^{(2)}(t) + \beta_{2i}^{(2)} \gamma_2^{(2)}(t)$, where $\gamma_1^{(1)}$, $\gamma_1^{(2)}$, $\gamma_2^{(1)}$ and g are mutually orthogonal functions on $L^2[0, 1]$ and $\gamma_2^{(2)} = (1 + v^2)^{-1/2} \{\gamma_2^{(1)} + vg\}$. Basically we create the alternative by the contamination of one of the “eigenfunctions” (in our case the second one) in the direction g and ensure $\|\gamma_2^{(2)}\| = 1$. The amount of the contamination is controlled by the parameter v . Note that the exact squared integral difference $\|\gamma_2^{(1)} - \gamma_2^{(2)}\|^2$ does not depend on function g . Thus, in the “functional sense” particular “direction of the alternative hypothesis” represented by the function g has no impact on the power of the test. However, since we are using a nonparametric estimation technique, we might expect that rough (highly fluctuating) functions g will yield higher error of estimation and, hence, decrease the precision (and power) of the test. Finally, a higher number of factor functions (L) in simulation may cause less precise approximation of critical values and more bootstrap replications and

larger sample-size may be needed. This can also be expected from Theorem 2 in Section 2.2—the variance of the estimated eigenfunctions depends on all eigenfunctions corresponding to nonzero eigenvalues.

4. Implied volatility analysis. In this section we present an application of the method discussed in previous sections to the implied volatilities of European options on the German stock index (ODAX). Implied volatilities are derived from the Black–Scholes (BS) pricing formula for European options; see Black and Scholes (1973). European call and put options are derivatives written on an underlying asset with price process S_i , which yield the pay-off $\max(S_I - K, 0)$ and $\max(K - S_I, 0)$, respectively. Here i denotes the current day, I the expiration day and K the strike price. Time to maturity is defined as $\tau = I - i$. The BS pricing formula for a Call option is

$$(14) \quad C_i(S_i, K, \tau, r, \sigma) = S_i \Phi(d_1) - K e^{-r\tau} \Phi(d_2),$$

where $d_1 = \frac{\ln(S_i/K) + (r + \sigma^2/2)\tau}{\sigma\sqrt{\tau}}$, $d_2 = d_1 - \sigma\sqrt{\tau}$, r is the risk-free interest rate, σ is the (unknown and constant) volatility parameter, and Φ denotes the c.d.f. of a standard normal distributed random variable. In (14) we assume the zero-dividend case. The Put option price P_i can be obtained from the put–call parity $P_i = C_i - S_i + e^{-r\tau}K$.

The implied volatility $\tilde{\sigma}$ is defined as the volatility σ , for which the BS price C_i in (14) equals the price \tilde{C}_i observed on the market. For a single asset, we obtain at each time point (day i) and for each maturity τ a IV function $\tilde{\sigma}_i^\tau(K)$. Practitioners often rescale the strike dimension by plotting this surface in terms of (futures) moneyness $\kappa = K/F_i(\tau)$, where $F_i(\tau) = S_i e^{r\tau}$.

Clearly, for given parameters S_i, r, K, τ the mapping from prices to IVs is a one-to-one mapping. The IV is often used for quoting the European options in financial practice, since it reflects the “uncertainty” of the financial market better than the option prices. It is also known that if the stock price drops, the IV raises (so-called leverage effect), motivates hedging strategies based on IVs. Consequently, for the purpose of this application, we will regard the BS–IV as an individual financial variable. The practical relevance of such an approach is justified by the volatility based financial products such as VDAX, which are commonly traded on the option markets.

The goal of this analysis is to study the dynamics of the IV functions for different maturities. More specifically, our aim is to construct low dimensional factor models based on the truncated Karhunen–Loève expansions (1) for the log-returns of the IV functions of options with different maturities and compare these factor models using the methodology presented in the previous sections. Analysis of IVs based on a low-dimensional factor model gives directly a descriptive insight into the structure of distribution

of the log-IV-returns—structure of the factors and empirical distribution of the factor loadings may be a good starting point for further pricing models. In practice, such a factor model can also be used in Monte Carlo based pricing methods and for risk-management (hedging) purposes. For comprehensive monographs on IV and IV-factor models, see [Hafner \(2004\)](#) or [Fengler \(2005b\)](#).

The idea of constructing and analyzing the factor models for log-IV-returns for different maturities was originally proposed by [Fengler, Härdle and Villa \(2003\)](#), who studied the dynamics of the IV via PCA on discretized IV functions for different maturity groups and tested the Common Principal Components (CPC) hypotheses (equality of eigenvectors and eigenspaces for different groups). [Fengler, Härdle and Villa \(2003\)](#) proposed a PCA-based factor model for log-IV-returns on (short) maturities 1, 2 and 3 months and grid of moneyness $[0.85, 0.9, 0.95, 1, 1.05, 1.1]$. They showed that the factor functions do not significantly differ and only the factor loadings differ across maturity groups. Their method relies on the CPC methodology introduced by [Flury \(1988\)](#) which is based on maximum likelihood estimation under the assumption of multivariate normality. The log-IV-returns are extracted by the two-dimensional Nadaraya–Watson estimate.

The main aim of this application is to reconsider their results in a functional sense. Doing so, we overcome two basic weaknesses of their approach. First, the factor model proposed by [Fengler, Härdle and Villa \(2003\)](#) is performed only on a sparse design of moneyness. However, in practice (e.g., in Monte Carlo pricing methods), evaluation of the model on a fine grid is needed. Using the functional PCA approach, we may overcome this difficulty and evaluate the factor model on an arbitrary fine grid. The second difficulty of the procedure proposed by [Fengler, Härdle and Villa \(2003\)](#) stems from the data design—on the exchange we cannot observe options with desired maturity on each day and we need to estimate them from the IV-functions with maturities observed on the particular day. Consequently, the two-dimensional Nadaraya–Watson estimator proposed by [Fengler, Härdle and Villa \(2003\)](#) results essentially in the (weighted) average of the IVs (with closest maturities) observed on a particular day, which may affect the test of the common eigenfunction hypothesis. We use the linear interpolation scheme in the *total variance* $\sigma_{\text{TOT},i}^2(\kappa, \tau) \stackrel{\text{def}}{=} (\sigma_i^\tau(\kappa))^2 \tau$, in order to recover the IV functions with fixed maturity (on day i). This interpolation scheme is based on the arbitrage arguments originally proposed by [Kahalé \(2004\)](#) for zero-dividend and zero-interest rate case and generalized for deterministic interest rate by [Fengler \(2005a\)](#). More precisely, having IVs with maturities observed on a particular day i : $\tilde{\sigma}_i^{\tau_{j_i}}(\kappa)$, $j_i = 1, \dots, p_{\tau_i}$, we calculate the corresponding total variance $\tilde{\sigma}_{\text{TOT},i}(\kappa, \tau_{j_i})$. From these total variances

we linearly interpolate the total variance with the desired maturity from the nearest maturities observed on day i . The total variance can be easily transformed to corresponding IV $\tilde{\sigma}_i^\tau(\kappa)$. As the last step, we calculate the log-returns $\Delta \log \tilde{\sigma}_i^\tau(\kappa) \stackrel{\text{def}}{=} \log \tilde{\sigma}_{i+1}^\tau(\kappa) - \log \tilde{\sigma}_i^\tau(\kappa)$. The log-IV-returns are observed for each maturity τ on a discrete grid κ_{ik}^τ . We assume that observed log-IV-return $\Delta \log \tilde{\sigma}_i^\tau(\kappa_{ik}^\tau)$ consists of true log-return of the IV function denoted by $\Delta \log \sigma_i^\tau(\kappa_{ik}^\tau)$ and possibly of some additional error ε_{ik}^τ . By setting $Y_{ik}^\tau := \Delta \log \tilde{\sigma}_i^\tau(\kappa_{ik}^\tau)$, $X_i^\tau(\kappa) := \Delta \log \sigma_i^\tau(\kappa)$, we obtain an analogue of the model (4) with the argument κ :

$$(15) \quad Y_{ik}^\tau = X_i^\tau(\kappa_{ik}) + \varepsilon_{ik}^\tau, \quad i = 1, \dots, n_\tau.$$

In order to simplify the notation and make the connection with the theoretical part clear, we will use the notation of (15).

For our analysis we use a recent data set containing daily data from January 2004 to June 2004 from the German–Swiss exchange (EUREX). Violations of the arbitrage-free assumptions (“obvious” errors in data) were corrected using the procedure proposed by [Fengler \(2005a\)](#). Similarly to [Fengler, Härdle and Villa \(2003\)](#), we excluded options with maturity smaller than 10 days, since these option-prices are known to be very noisy, partially because of a special and arbitrary setup in the pricing systems of the dealers. Using the interpolation scheme described above, we calculate the log-IV-returns for two maturity groups: “1M” group with maturity $\tau = 0.12$ (measured in years) and “3M” group with maturity $\tau = 0.36$. The observed log-IV-returns are denoted by Y_{ik}^{1M} , $k = 1, \dots, K_i^{1M}$, Y_{ik}^{3M} , $k = 1, \dots, K_i^{3M}$. Since we ensured that for no i , the interpolation procedure uses data with the same maturity for both groups, this procedure has no impact on the independence of both samples.

The underlying models based on the truncated version of (3) are as follows:

$$(16) \quad X_i^{1M}(\kappa) = \bar{X}^{1M}(\kappa) + \sum_{r=1}^{L_{1M}} \hat{\beta}_{ri}^{1M} \widehat{\gamma}_r^{1M}(\kappa), \quad i = 1, \dots, n_{1M},$$

$$(17) \quad X_i^{3M}(\kappa) = \bar{X}^{3M}(\kappa) + \sum_{r=1}^{L_{3M}} \hat{\beta}_{ri}^{3M} \widehat{\gamma}_r^{3M}(\kappa), \quad i = 1, \dots, n_{3M}.$$

Models (16) and (17) can serve, for example, in a Monte Carlo pricing tool in the risk management for pricing exotic options where the whole path of implied volatilities is needed to determine the price. Estimating the factor functions in (16) and (17) by eigenfunctions displayed in [Figure 1](#), we only need to fit the (estimated) factor loadings $\hat{\beta}_{ji}^{1M}$ and $\hat{\beta}_{ji}^{3M}$. The pillar of the model is the dimension reduction. Keeping the factor function fixed for a certain time period, we need to analyze (two) multivariate random processes

of the factor loadings. For the purposes of this paper we will focus on the comparison of factors from models (16) and (17) and the technical details of the factor loading analysis will not be discussed here, since in this respect we refer to [Fengler, Härdle and Villa \(2003\)](#), who proposed to fit the factor loadings by centered normal distributions with diagonal variance matrix containing the corresponding eigenvalues. For a deeper discussion of the fitting of factor loadings using a more sophisticated approach, basically based on (possibly multivariate) GARCH models; see [Fengler \(2005b\)](#).

From our data set we obtained 88 functional observations for the 1M group (n_{1M}) and 125 observations for the 3M group (n_{3M}). We will estimate the model on the interval for futures moneyness $\kappa \in [0.8, 1.1]$. In comparison to [Fengler, Härdle and Villa \(2003\)](#), we may estimate models (16) and (17) on an arbitrary fine grid (we used an equidistant grid of 500 points on the interval $[0.8, 1.1]$). For illustration, the Nadaraya–Watson (NW) estimator of resulting log-returns is plotted in [Figure 2](#). The smoothing parameters have been chosen in accordance with the requirements in [Section 2.2](#). As argued in [Section 2.2](#), we should use small smoothing parameters in order to avoid a possible bias in the estimated eigenfunctions. Thus, we use for each i essentially the smallest bandwidth b_i that guarantees that estimator \hat{X}_i is defined on the entire support $[0.8, 1.1]$.

Using the procedures described in [Section 2.1](#), we first estimate the eigenfunctions of both maturity groups. The estimated eigenfunctions are plotted in [Figure 1](#). The structure of the eigenfunctions is in accordance with other empirical studies on IV-surfaces. For a deeper discussion and economical interpretation, see, for example, [Fengler, Härdle and Mammen \(2007\)](#) or [Fengler, Härdle and Villa \(2003\)](#).

Clearly, the ratio of the variance explained by the k th factor function is given by the quantity $\hat{\nu}_k^{1M} = \hat{\lambda}_k^{1M} / \sum_{j=1}^{r_{1M}} \hat{\lambda}_j^{1M}$ for the 1M group and, correspondingly, by $\hat{\nu}_k^{3M}$ for the 3M group. In [Table 3](#) we list the contributions of the factor functions. Looking at [Table 3](#), we can see that 4th factor functions explain less than 1% of the variation. This number was the “threshold” for the choice of L_{1M} and L_{2M} .

We can observe (see [Figure 1](#)) that the factor functions for both groups are similar. Thus, in the next step we use the bootstrap test for testing the

TABLE 3
Variance explained by the eigenfunctions

	Var. explained 1M	Var. explained 3M
$\hat{\nu}_1^T$	89.9%	93.0%
$\hat{\nu}_2^T$	7.7%	4.2%
$\hat{\nu}_3^T$	1.7%	1.0%
$\hat{\nu}_4^T$	0.6%	0.4%

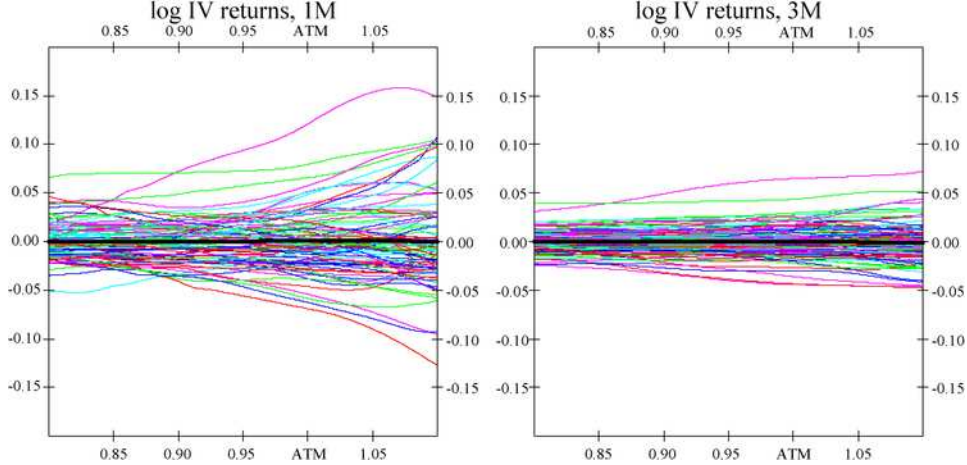


FIG. 2. *Nadaraya–Watson estimate of the log-IV-returns for maturity 1M (left figure) and 3M (right figure). The bold line is the sample mean of the corresponding group.*

equality of the factor functions. We use 2000 bootstrap replications. The test of equality of the eigenfunctions was rejected for the first eigenfunction for the analyzed time period (January 2004–June 2004) at a significance level $\alpha = 0.05$ (P-value 0.01). We may conclude that the (first) factor functions are not identical in the factor model for both maturity groups. However, from a practical point of view, we are more interested in checking the appropriateness of the entire models for a fixed number of factors: $L = 2$ or $L = 3$ in (16) and (17). This requirement translates into the testing of the equality of eigenspaces. Thus, in the next step we use the same setup (2000 bootstrap replications) to test the hypotheses that the first two and first three eigenfunctions span the same eigenspaces \mathcal{E}_L^{1M} and \mathcal{E}_L^{3M} . None of the hypotheses for $L = 2$ and $L = 3$ is rejected at significance level $\alpha = 0.05$ (P-value is 0.61 for $L = 2$ and 0.09 for $L = 3$). Summarizing, even in the functional sense we have no significant reason to reject the hypothesis of common eigenspaces for these two maturity groups. Using this hypothesis, the factors governing the movement of the returns of IV surface are invariant to time to maturity, only their relative importance can vary. This leads to the common factor model: $X_i^\tau(\kappa) = \bar{X}^\tau(\kappa) + \sum_{r=1}^{L_\tau} \hat{\beta}_{ri}^\tau \widehat{\gamma}_r(\kappa)$, $i = 1, \dots, n_\tau$, $\tau = 1M, 3M$, where $\gamma_r := \gamma_r^{1M} = \gamma_r^{3M}$. Beside contributing to the understanding of the structure of the IV function dynamics, the common factor model helps us to reduce the number of functional factors by half compared to models (16) and (17). Furthermore, from the technical point of view, we also obtain an additional dimension reduction and higher estimation precision, since under this hypothesis we may estimate the eigenfunctions from the (individually centered) pooled sample $X_i(\kappa)^{1M}$, $i = 1, \dots, n_{1M}$, $X_i^{3M}(\kappa)$, $i =$

$1, \dots, n_{3M}$. The main improvement compared to the multivariate study by [Fengler, Härdle and Villa \(2003\)](#) is that our test is performed in the functional sense – it does not depend on particular discretization and our factor model can be evaluated on an arbitrary fine grid.

APPENDIX: MATHEMATICAL PROOFS

In the following, $\|v\| = (\int_0^1 v(t)^2 dt)^{1/2}$ will denote the L^2 -norm for any square integrable function v . At the same time, $\|a\| = (\frac{1}{k} \sum_{i=1}^k a_i^2)^{1/2}$ will indicate the Euclidean norm, whenever $a \in \mathbb{R}^k$ is a k -vector for some $k \in \mathbb{N}$.

In the proof of [Theorem 1](#), E_ε and Var_ε denote expectation and variance with respect to ε only (i.e., conditional on t_{ij} and X_i).

PROOF OF THEOREM 1. Recall the definition of the $\chi_i(t)$ and note that $\chi_i(t) = \chi_i^X(t) + \chi_i^\varepsilon(t)$, where

$$\chi_i^\varepsilon(t) = \sum_{j=1}^{T_i} \varepsilon_{i(j)} I\left(t \in \left[\frac{t_{i(j-1)} + t_{i(j)}}{2}, \frac{t_{i(j)} + t_{i(j+1)}}{2}\right)\right),$$

as well as

$$\chi_i^X(t) = \sum_{j=1}^{T_i} X_i(t_{i(j)}) I\left(t \in \left[\frac{t_{i(j-1)} + t_{i(j)}}{2}, \frac{t_{i(j)} + t_{i(j+1)}}{2}\right)\right)$$

for $t \in [0, 1]$, $t_{i(0)} = -t_{i(1)}$ and $t_{i(T_i+1)} = 2 - t_{i(T_i)}$. Similarly, $\chi_i^*(t) = \chi_i^{X^*}(t) + \chi_i^{\varepsilon^*}(t)$.

By [Assumption 2](#), $E(|t_{i(j)} - t_{i(j-1)}|^s) = \mathcal{O}(T^{-s})$ for $s = 1, \dots, 4$, and the convergence is uniform in $j < n$. Our assumptions on the structure of X_i together with some straightforward Taylor expansions then lead to

$$\langle \chi_i, \chi_j \rangle = \langle X_i, X_j \rangle + \mathcal{O}_p(1/T)$$

and

$$\langle \chi_i, \chi_i^* \rangle = \|X_i\|^2 + \mathcal{O}_p(1/T).$$

Moreover,

$$\begin{aligned} E_\varepsilon(\langle \chi_i^\varepsilon, \chi_j^X \rangle) &= 0, & E_\varepsilon(\|\chi_i^\varepsilon\|^2) &= \sigma_i^2, \\ E_\varepsilon(\langle \chi_i^\varepsilon, \chi_i^{\varepsilon^*} \rangle) &= 0, & E_\varepsilon(\langle \chi_i^\varepsilon, \chi_i^{\varepsilon^*} \rangle^2) &= \mathcal{O}_p(1/T), \\ E_\varepsilon(\langle \chi_i^\varepsilon, \chi_j^X \rangle^2) &= \mathcal{O}_p(1/T), & E_\varepsilon(\langle \chi_i^\varepsilon, \chi_j^X \rangle \langle \chi_k^\varepsilon, \chi_l^X \rangle) &= 0 \quad \text{for } i \neq k, \end{aligned}$$

$$E_\varepsilon(\langle \chi_i^\varepsilon, \chi_j^\varepsilon \rangle \langle \chi_i^\varepsilon, \chi_k^\varepsilon \rangle) = 0 \quad \text{for } j \neq k \text{ and } E_\varepsilon(\|\chi_i^\varepsilon\|^4) = \mathcal{O}_p(1)$$

hold (uniformly) for all $i, j = 1, \dots, n$.

Consequently, $E_\varepsilon(\|\bar{\chi}\|^2 - \|\bar{X}\|^2) = \mathcal{O}_p(T^{-1} + n^{-1})$.

When using these relations, it is easily seen that for all $i, j = 1, \dots, n$

$$(18) \quad \begin{aligned} \widehat{M}_{ij} - M_{ij} &= \mathcal{O}_p(T^{-1/2} + n^{-1}) \quad \text{and} \\ \text{tr}\{(\widehat{M} - M)^2\}^{1/2} &= \mathcal{O}_p(1 + nT^{-1/2}). \end{aligned}$$

Since the orthonormal eigenvectors p_q of M satisfy $\|p_q\| = 1$, we furthermore obtain for any $i = 1, \dots, n$ and all $q = 1, 2, \dots$

$$(19) \quad \sum_{j=1}^n p_{jq} \left\{ \widehat{M}_{ij} - M_{ij} - \int_0^1 \chi_i^\varepsilon(t) \chi_j^X(t) dt \right\} = \mathcal{O}_p(T^{-1/2} + n^{-1/2}),$$

as well as

$$(20) \quad \sum_{j=1}^n p_{jq} \int_0^1 \chi_i^\varepsilon(t) \chi_j^X(t) dt = \mathcal{O}_p\left(\frac{n^{1/2}}{T^{1/2}}\right)$$

and

$$(21) \quad \sum_{i=1}^n a_i \sum_{j=1}^n p_{jq} \int_0^1 \chi_i^\varepsilon(t) \chi_j^X(t) dt = \mathcal{O}_p\left(\frac{n^{1/2}}{T^{1/2}}\right)$$

for any further vector a with $\|a\| = 1$.

Recall that the j th largest eigenvalue l_j satisfies $n\hat{\lambda}_j = l_j$. Since by assumption $\inf_{s \neq r} |\lambda_r - \lambda_s| > 0$, the results of [Dauxois, Pousse and Romain \(1982\)](#) imply that $\hat{\lambda}_r$ converges to λ_r as $n \rightarrow \infty$, and $\sup_{s \neq r} \frac{1}{|\hat{\lambda}_r - \hat{\lambda}_s|} = \mathcal{O}_p(1)$, which leads to $\sup_{s \neq r} \frac{1}{|l_r - l_s|} = \mathcal{O}_p(1/n)$. Assertion (a) of Lemma A of [Kneip and Utikal \(2001\)](#) together with (18)–(21) then implies that

$$(22) \quad \begin{aligned} \left| \hat{\lambda}_r - \frac{\hat{l}_r}{n} \right| &= n^{-1} |l_r - \hat{l}_r| = n^{-1} |p_r^\top (\widehat{M} - M) p_r| + \mathcal{O}_p(T^{-1} + n^{-1}) \\ &= \mathcal{O}_p\{(nT)^{-1/2} + T^{-1} + n^{-1}\}. \end{aligned}$$

When analyzing the difference between the estimated and true eigenvectors \hat{p}_r and p_r , assertion (b) of Lemma A of [Kneip and Utikal \(2001\)](#) together with (18) lead to

$$(23) \quad \hat{p}_r - p_r = -\mathcal{S}_r(\widehat{M} - M)p_r + \mathcal{R}_r, \quad \text{with } \|\mathcal{R}_r\| = \mathcal{O}_p(T^{-1} + n^{-1})$$

and $\mathcal{S}_r = \sum_{s \neq r} \frac{1}{l_s - l_r} p_s p_s^\top$. Since $\sup_{\|a\|=1} a^\top \mathcal{S}_r a \leq \sup_{s \neq r} \frac{1}{|l_r - l_s|} = \mathcal{O}_p(1/n)$, we can conclude that

$$(24) \quad \|\hat{p}_r - p_r\| = \mathcal{O}_p(T^{-1/2} + n^{-1}),$$

and our assertion on the sequence $n^{-1} \sum_i (\hat{\beta}_{ri} - \hat{\beta}_{ri;T})^2$ is an immediate consequence.

Let us now consider assertion (ii). The well-known properties of local linear estimators imply that $|\mathbb{E}_\varepsilon\{\hat{X}_i(t) - X_i(t)\}| = \mathcal{O}_p(b^2)$, as well as $\text{Var}_\varepsilon\{\hat{X}_i(t)\} = \mathcal{O}_p\{Tb\}$, and the convergence is uniform for all i, n . Furthermore, due to the independence of the error term ε_{ij} , $\text{Cov}_\varepsilon\{\hat{X}_i(t), \hat{X}_j(t)\} = 0$ for $i \neq j$. Therefore,

$$\left| \hat{\gamma}_r(t) - \frac{1}{\sqrt{l_r}} \sum_{i=1}^n p_{ir} \hat{X}_i(t) \right| = \mathcal{O}_p\left(b^2 + \frac{1}{\sqrt{nTb}}\right).$$

On the other hand, (18)–(24) imply that with $\hat{X}(t) = (\hat{X}_1(t), \dots, \hat{X}_n(t))^\top$

$$\begin{aligned} & \left| \hat{\gamma}_{r;T}(t) - \frac{1}{\sqrt{l_r}} \sum_{i=1}^n p_{ir} \hat{X}_i(t) \right| \\ &= \left| \frac{1}{\sqrt{l_r}} \sum_{i=1}^n (\hat{p}_{ir} - p_{ir}) X_i(t) + \frac{1}{\sqrt{l_r}} \sum_{i=1}^n (\hat{p}_{ir} - p_{ir}) \{\hat{X}_i(t) - X_i(t)\} \right| \\ & \quad + \mathcal{O}_p(T^{-1} + n^{-1}) \\ &= \frac{\|\mathcal{S}_r X(t)\|}{\sqrt{l_r}} \left| p_r^\top (\hat{M} - M) \mathcal{S}_r \frac{X(t)}{\|\mathcal{S}_r X(t)\|} \right| \\ & \quad + \mathcal{O}_p(b^2 T^{-1/2} + T^{-1} b^{-1/2} + n^{-1}) \\ &= \mathcal{O}_p(n^{-1/2} T^{-1/2} + b^2 T^{-1/2} + T^{-1} b^{-1/2} + n^{-1}). \end{aligned}$$

This proves the theorem. \square

PROOF OF THEOREM 2. First consider assertion (i). By definition,

$$\bar{X}(t) - \mu(t) = n^{-1} \sum_{i=1}^n \{X_i(t) - \mu(t)\} = \sum_r \left(n^{-1} \sum_{i=1}^n \beta_{ri} \right) \gamma_r(t).$$

Recall that, by assumption, β_{ri} are independent, zero mean random variables with variance λ_r , and that the above series converges with probability 1. When defining the truncated series

$$V(q) = \sum_{r=1}^q \left(n^{-1} \sum_{i=1}^n \beta_{ri} \right) \gamma_r(t),$$

standard central limit theorems therefore imply that $\sqrt{n}V(q)$ is asymptotically $N(0, \sum_{r=1}^q \lambda_r \gamma_r(t)^2)$ distributed for any possible $q \in \mathbb{N}$.

The assertion of a $N(0, \sum_{r=1}^\infty \lambda_r \gamma_r(t)^2)$ limiting distribution now is a consequence of the fact that for all $\delta_1, \delta_2 > 0$ there exists a q_δ such that $P\{|\sqrt{n}V(q) - \sqrt{n}\sum_r(n^{-1}\sum_{i=1}^n\beta_{ri})\gamma_r(t)| > \delta_1\} < \delta_2$ for all $q \geq q_\delta$ and all n sufficiently large.

In order to prove assertions (i) and (ii), consider some fixed $r \in \{1, 2, \dots\}$ with $\lambda_{r-1} > \lambda_r > \lambda_{r+1}$. Note that Γ as well as $\hat{\Gamma}_n$ are nuclear, self-adjoint and non-negative linear operators with $\Gamma v = \int \sigma(t, s)v(s) ds$ and $\hat{\Gamma}_n v = \int \hat{\sigma}(t, s)v(s) ds$, $v \in L^2[0, 1]$. For $m \in \mathbb{N}$, let Π_m denote the orthogonal projector from $L^2[0, 1]$ into the m -dimensional linear space spanned by $\{\gamma_1, \dots, \gamma_m\}$, that is, $\Pi_m v = \sum_{j=1}^m \langle v, \gamma_j \rangle \gamma_j$, $v \in L^2[0, 1]$. Now consider the operator $\Pi_m \hat{\Gamma}_n \Pi_m$, as well as its eigenvalues and corresponding eigenfunctions denoted by $\hat{\lambda}_{1,m} \geq \hat{\lambda}_{2,m} \geq \dots$ and $\hat{\gamma}_{1,m}, \hat{\gamma}_{2,m}, \dots$, respectively. It follows from well-known results in the Hilbert space theory that $\Pi_m \hat{\Gamma}_n \Pi_m$ converges strongly to $\hat{\Gamma}_n$ as $m \rightarrow \infty$. Furthermore, we obtain (Rayleigh–Ritz theorem)

$$(25) \quad \lim_{m \rightarrow \infty} \hat{\lambda}_{r,m} = \lambda_r \quad \text{and} \quad \lim_{m \rightarrow \infty} \|\hat{\gamma}_r - \hat{\gamma}_{r,m}\| = 0 \quad \text{if } \hat{\lambda}_{r-1} > \hat{\lambda}_r > \hat{\lambda}_{r+1}.$$

Note that under the above condition $\hat{\gamma}_r$ is uniquely determined up to sign, and recall that we always assume that the right “versions” (with respect to sign) are used so that $\langle \hat{\gamma}_r, \hat{\gamma}_{r,m} \rangle \geq 0$. By definition, $\beta_{ji} = \int \gamma_j(t) \{X_i(t) - \mu(t)\} dt$, and therefore, $\int \gamma_j(t) \{X_i(t) - \bar{X}(t)\} dt = \beta_{ji} - \bar{\beta}_j$, as well as $X_i - \bar{X} = \sum_j (\beta_{ji} - \bar{\beta}_j) \gamma_j$, where $\bar{\beta}_j = \frac{1}{n} \sum_{i=1}^n \beta_{ji}$. When analyzing the structure of $\Pi_m \hat{\Gamma}_n \Pi_m$ more deeply, we can verify that $\Pi_m \hat{\Gamma}_n \Pi_m v = \int \hat{\sigma}_m(t, s)v(s) ds$, $v \in L^2[0, 1]$, with

$$\hat{\sigma}_m(t, s) = g_m(t)^\top \hat{\Sigma}_m g_m(s),$$

where $g_m(t) = (\gamma_1(t), \dots, \gamma_m(t))^\top$, and where $\hat{\Sigma}_m$ is the $m \times m$ matrix with elements $\{\frac{1}{n} \sum_{i=1}^n (\beta_{ji} - \bar{\beta}_j)(\beta_{ki} - \bar{\beta}_k)\}_{j,k=1, \dots, m}$. Let $\lambda_1(\hat{\Sigma}_m) \geq \lambda_2(\hat{\Sigma}_m) \geq \dots \geq \lambda_m(\hat{\Sigma}_m)$ and $\hat{\zeta}_{1,m}, \dots, \hat{\zeta}_{m,m}$ denote eigenvalues and corresponding eigenvectors of $\hat{\Sigma}_m$. Some straightforward algebra then shows that

$$(26) \quad \hat{\lambda}_{r,m} = \lambda_r(\hat{\Sigma}_m), \quad \hat{\gamma}_{r,m} = g_m(t)^\top \hat{\zeta}_{r,m}.$$

We will use Σ_m to represent the $m \times m$ diagonal matrix with diagonal entries $\lambda_1 \geq \dots \geq \lambda_m$. Obviously, the corresponding eigenvectors are given by the m -dimensional unit vectors denoted by $e_{1,m}, \dots, e_{m,m}$. Lemma A of [Kneip and Utikal \(2001\)](#) now implies that the differences between eigenvalues and eigenvectors of Σ_m and $\hat{\Sigma}_m$ can be bounded by

$$(27) \quad \hat{\lambda}_{r,m} - \lambda_r = \text{tr}\{e_{r,m} e_{r,m}^\top (\hat{\Sigma}_m - \Sigma_m)\} + \tilde{R}_{r,m},$$

$$\text{with } \tilde{R}_{r,m} \leq \frac{6 \sup_{\|a\|=1} a^\top (\hat{\Sigma}_m - \Sigma_m)^2 a}{\min_s |\lambda_s - \lambda_r|},$$

$$(28) \quad \hat{\zeta}_{r,m} - e_{r,m} = -S_{r,m}(\hat{\Sigma}_m - \Sigma_m)e_{r,m} + R_{r,m}^*,$$

$$\text{with } \|R_{r,m}^*\| \leq \frac{6 \sup_{\|a\|=1} a^\top (\hat{\Sigma}_m - \Sigma_m)^2 a}{\min_s |\lambda_s - \lambda_r|^2},$$

where $S_{r,m} = \sum_{s \neq r} \frac{1}{\lambda_s - \lambda_r} e_{s,m} e_{s,m}^\top$.

Assumption 1 implies $E(\bar{\beta}_r) = 0$, $\text{Var}(\bar{\beta}_r) = \frac{\lambda_r}{n}$, and with $\delta_{ii} = 1$, as well as $\delta_{ij} = 0$ for $i \neq j$, we obtain

$$\begin{aligned}
& E \left\{ \sup_{\|a\|=1} a^\top (\hat{\Sigma}_m - \Sigma_m)^2 a \right\} \\
& \leq E \{ \text{tr}[(\hat{\Sigma}_m - \Sigma_m)^2] \} \\
& = E \left\{ \sum_{j,k=1}^m \left[\frac{1}{n} \sum_{i=1}^n (\beta_{ji} - \bar{\beta}_j)(\beta_{ki} - \bar{\beta}_k) - \delta_{jk} \lambda_j \right]^2 \right\} \\
(29) \quad & \leq E \left\{ \sum_{j,k=1}^\infty \left[\frac{1}{n} \sum_{i=1}^n (\beta_{ji} - \bar{\beta}_j)(\beta_{ki} - \bar{\beta}_k) - \delta_{jk} \lambda_j \right]^2 \right\} \\
& = \frac{1}{n} \left(\sum_j \sum_k E \{ \beta_{ji}^2 \beta_{ki}^2 \} \right) + \mathcal{O}(n^{-1}) = \mathcal{O}(n^{-1}),
\end{aligned}$$

for all m . Since $\text{tr}\{e_{r,m} e_{r,m}^\top (\hat{\Sigma}_m - \Sigma_m)\} = \frac{1}{n} \sum_{i=1}^n (\beta_{ri} - \bar{\beta}_r)^2 - \lambda_r$, (25), (26), (27) and (29) together with standard central limit theorems imply that

$$\begin{aligned}
\sqrt{n}(\hat{\lambda}_r - \lambda_r) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n (\beta_{ri} - \bar{\beta}_r)^2 - \lambda_r + \mathcal{O}_p(n^{-1/2}) \\
(30) \quad &= \frac{1}{\sqrt{n}} \sum_{i=1}^n [(\beta_{ri})^2 - E\{(\beta_{ri})^2\}] + \mathcal{O}_p(n^{-1/2}) \\
&\xrightarrow{\mathcal{L}} N(0, \Lambda_r).
\end{aligned}$$

It remains to prove assertion (iii). Relations (26) and (28) lead to

$$\begin{aligned}
\hat{\gamma}_{r,m}(t) - \gamma_r(t) &= g_m(t)^\top (\hat{\zeta}_{r,m} - e_{r,m}) \\
(31) \quad &= - \sum_{s \neq r}^m \left\{ \frac{1}{n(\lambda_s - \lambda_r)} \sum_{i=1}^n (\beta_{si} - \bar{\beta}_s)(\beta_{ri} - \bar{\beta}_r) \right\} \gamma_s(t) \\
&\quad + g_m(t)^\top R_{r,m}^*,
\end{aligned}$$

where due to (29) the function $g_m(t)^\top R_{r,m}^*$ satisfies

$$\begin{aligned}
E(\|g_m^\top R_{r,m}^*\|) &= E(\|R_{r,m}^*\|) \\
&\leq \frac{6}{n \min_s |\lambda_s - \lambda_r|^2} \left(\sum_j \sum_k E \{ \beta_{ji}^2 \beta_{ki}^2 \} \right) + \mathcal{O}(n^{-1}),
\end{aligned}$$

for all m . By Assumption 1, the series in (31) converge with probability 1 as $m \rightarrow \infty$.

Obviously, the event $\hat{\lambda}_{r-1} > \hat{\lambda}_r > \hat{\lambda}_{r+1}$ occurs with probability 1. Since m is arbitrary, we can therefore conclude from (25) and (31) that

$$\begin{aligned}
& \hat{\gamma}_r(t) - \gamma_r(t) \\
(32) \quad &= - \sum_{s \neq r} \left\{ \frac{1}{n(\lambda_s - \lambda_r)} \sum_{i=1}^n (\beta_{si} - \bar{\beta}_s)(\beta_{ri} - \bar{\beta}_r) \right\} \gamma_s(t) + R_r^*(t) \\
&= - \sum_{s \neq r} \left\{ \frac{1}{n(\lambda_s - \lambda_r)} \sum_{i=1}^n \beta_{si} \beta_{ri} \right\} \gamma_s(t) + R_r(t),
\end{aligned}$$

where $\|R_r^*\| = \mathcal{O}_p(n^{-1})$, as well as $\|R_r\| = \mathcal{O}_p(n^{-1})$. Moreover, $\sqrt{n} \times \sum_{s \neq r} \left\{ \frac{1}{n(\lambda_s - \lambda_r)} \sum_{i=1}^n \beta_{si} \beta_{ri} \right\} \gamma_s(t)$ is a zero mean random variable with variance $\sum_{q \neq r} \sum_{s \neq r} \frac{E[\beta_{ri}^2 \beta_{qi} \beta_{si}]}{(\lambda_q - \lambda_r)(\lambda_s - \lambda_r)} \gamma_q(t) \gamma_s(t) < \infty$. By Assumption 1, it follows from standard central limit arguments that for any $q \in \mathbb{N}$ the truncated series $\sqrt{n} W(q) \stackrel{\text{def}}{=} \sqrt{n} \sum_{s=1, s \neq r}^q \left[\frac{1}{n(\lambda_s - \lambda_r)} \sum_{i=1}^n \beta_{si} \beta_{ri} \right] \gamma_s(t)$ is asymptotically normal distributed. The asserted asymptotic normality of the complete series then follows from an argument similar to the one used in the proof of assertion (i). \square

PROOF OF THEOREM 3. The results of Theorem 2 imply that

$$\begin{aligned}
(33) \quad n\Delta_1 &= \int \left(\sum_r \left(\frac{1}{\sqrt{q_1 n_1}} \sum_{i=1}^{n_1} \beta_{ri}^{(1)} \gamma_r^{(1)}(t) \right. \right. \\
&\quad \left. \left. - \sum_r \frac{1}{\sqrt{q_2 n_2}} \sum_{i=1}^{n_2} \beta_{ri}^{(2)} \gamma_r^{(2)}(t) \right) \right)^2 dt.
\end{aligned}$$

Furthermore, independence of $X_i^{(1)}$ and $X_i^{(2)}$ together with (30) imply that

$$\begin{aligned}
(34) \quad \sqrt{n}[\hat{\lambda}_r^{(1)} - \lambda_r^{(1)} - \{\hat{\lambda}_r^{(2)} - \lambda_r^{(2)}\}] &\xrightarrow{\mathcal{L}} N\left(0, \frac{\Lambda_r^{(1)}}{q_1} + \frac{\Lambda_r^{(2)}}{q_2}\right) \quad \text{and} \\
\frac{n}{\Lambda_r^{(1)}/q_1 + \Lambda_r^{(2)}/q_2} \Delta_{3,r} &\xrightarrow{\mathcal{L}} \chi_1^2.
\end{aligned}$$

Furthermore, (32) leads to

$$\begin{aligned}
(35) \quad n\Delta_{2,r} &= \left\| \sum_{s \neq r} \left\{ \frac{1}{\sqrt{q_1 n_1} (\lambda_s^{(1)} - \lambda_r^{(1)})} \sum_{i=1}^{n_1} \beta_{si}^{(1)} \beta_{ri}^{(1)} \right\} \gamma_s^{(1)} \right. \\
&\quad \left. - \sum_{s \neq r} \left\{ \frac{1}{\sqrt{q_2 n_2} (\lambda_s^{(2)} - \lambda_r^{(2)})} \sum_{i=1}^{n_2} \beta_{si}^{(2)} \beta_{ri}^{(2)} \right\} \gamma_s^{(2)} \right\|^2 + \mathcal{O}_p(n^{-1/2})
\end{aligned}$$

and

$$\begin{aligned}
n\Delta_{4,L} &= n \iint \left[\sum_{r=1}^L \gamma_r^{(1)}(t) \{ \hat{\gamma}_r^{(1)}(u) - \gamma_r^{(1)}(u) \} \right. \\
&\quad + \gamma_r^{(1)}(u) \{ \hat{\gamma}_r^{(1)}(t) - \gamma_r^{(1)}(t) \} \\
&\quad - \sum_{r=1}^L \gamma_r^{(2)}(t) \{ \hat{\gamma}_r^{(2)}(u) - \gamma_r^{(2)}(u) \} \\
&\quad \left. + \gamma_r^{(2)}(u) \{ \hat{\gamma}_r^{(2)}(t) - \gamma_r^{(2)}(t) \} \right]^2 dt du + \mathcal{O}_p(n^{-1/2}) \\
(36) \quad &= \iint \left[\sum_{r=1}^L \sum_{s>L} \left\{ \frac{1}{\sqrt{q_1 n_1} (\lambda_s^{(1)} - \lambda_r^{(1)})} \sum_{i=1}^{n_1} \beta_{si}^{(1)} \beta_{ri}^{(1)} \right\} \right. \\
&\quad \times \{ \gamma_r^{(1)}(t) \gamma_s^{(1)}(u) + \gamma_r^{(1)}(u) \gamma_s^{(1)}(t) \} \\
&\quad - \sum_{r=1}^L \sum_{s>L} \left\{ \frac{1}{\sqrt{q_2 n_2} (\lambda_s^{(2)} - \lambda_r^{(2)})} \sum_{i=1}^{n_2} \beta_{si}^{(2)} \beta_{ri}^{(2)} \right\} \\
&\quad \left. \times \{ \gamma_r^{(2)}(t) \gamma_s^{(2)}(u) + \gamma_r^{(2)}(u) \gamma_s^{(2)}(t) \} \right]^2 dt du \\
&\quad + \mathcal{O}_p(n^{-1/2}).
\end{aligned}$$

In order to verify (36), note that $\sum_{r=1}^L \sum_{s=1, s \neq r}^L \frac{1}{(\lambda_s^{(p)} - \lambda_r^{(p)})} a_r a_s = 0$ for $p = 1, 2$ and all possible sequences a_1, \dots, a_L . It is clear from our assumptions that all sums involved converge with probability 1. Recall that $E(\beta_{ri}^{(p)} \beta_{si}^{(p)}) = 0$, $p = 1, 2$ for $r \neq s$.

It follows that $\tilde{X}_r^{(p)} := \frac{1}{\sqrt{q_p n_p}} \sum_{s \neq r} \sum_{i=1}^{n_p} \frac{\beta_{si}^{(p)} \beta_{ri}^{(p)}}{\lambda_s^{(p)} - \lambda_r^{(p)}} \gamma_s^{(p)}$, $p = 1, 2$, is a continuous, zero mean random function on $L^2[0, 1]$, and, by assumption, $E(\|\tilde{X}_r^{(p)}\|^2) < \infty$. By Hilbert space central limit theorems [see, e.g., [Araujo and Giné \(1980\)](#)], $\tilde{X}_r^{(p)}$ thus converges in distribution to a Gaussian random function $\xi_r^{(p)}$ as $n \rightarrow \infty$. Obviously, $\xi_r^{(1)}$ is independent of $\xi_r^{(2)}$. We can conclude that $n\Delta_{4,L}$ possesses a continuous limit distribution $F_{4,L}$ defined by the distribution of $\iint [\sum_{r=1}^L \{ \xi_r^{(1)}(t) \gamma_r^{(1)}(u) + \xi_r^{(1)}(u) \gamma_r^{(1)}(t) \} - \sum_{r=1}^L \{ \xi_r^{(2)}(t) \gamma_r^{(2)}(u) + \xi_r^{(2)}(u) \gamma_r^{(2)}(t) \}]^2 dt du$. Similar arguments show the existence of continuous limit distributions F_1 and $F_{2,r}$ of $n\Delta_1$ and $n\Delta_{2,r}$.

For given $q \in \mathbb{N}$, define vectors $b_{i1}^{(p)} = (\beta_{1i}^{(p)}, \dots, \beta_{qi}^{(p)})^\top \in \mathbb{R}^q$, $b_{i2}^{(p)} = (\beta_{1i}^{(p)} \beta_{ri}^{(p)}, \dots, \beta_{r-1,i}^{(p)} \beta_{ri}^{(p)}, \beta_{r+1,i}^{(p)} \beta_{ri}^{(p)}, \dots, \beta_{qi}^{(p)} \beta_{ri}^{(p)})^\top \in \mathbb{R}^{q-1}$ and $b_{i3} = (\beta_{1i}^{(p)} \beta_{2i}^{(p)},$

$\dots, \beta_{qi}^{(p)} \beta_{Li}^{(p)} \top \in \mathbb{R}^{(q-1)L}$. When the infinite sums over r in (33), respectively $s \neq r$ in (35) and (36), are restricted to $q \in \mathbb{N}$ components (i.e., \sum_r and $\sum_{s>L}$ are replaced by $\sum_{r \leq q}$ and $\sum_{L < s \leq q}$), then the above relations can generally be presented as limits $n\Delta = \lim_{q \rightarrow \infty} n\Delta(q)$ of quadratic forms

$$(37) \quad \begin{aligned} n\Delta_1(q) &= \begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} b_{i1}^{(1)} \\ \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} b_{i1}^{(2)} \end{pmatrix}^\top Q_1^q \begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} b_{i1}^{(1)} \\ \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} b_{i1}^{(2)} \end{pmatrix}, \\ n\Delta_{2,r}(q) &= \begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} b_{i2}^{(1)} \\ \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} b_{i2}^{(2)} \end{pmatrix}^\top Q_2^q \begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} b_{i2}^{(1)} \\ \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} b_{i2}^{(2)} \end{pmatrix}, \\ n\Delta_{4,L}(q) &= \begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} b_{i3}^{(1)} \\ \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} b_{i3}^{(2)} \end{pmatrix}^\top Q_3^q \begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} b_{i3}^{(1)} \\ \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} b_{i3}^{(2)} \end{pmatrix}, \end{aligned}$$

where the elements of the $2q \times 2q$, $2(q-1) \times 2(q-1)$ and $2L(q-1) \times 2L(q-1)$ matrices Q_1^q , Q_2^q and Q_3^q can be computed from the respective (q -element) version of (33)–(36). Assumption 1 implies that all series converge with probability 1 as $q \rightarrow \infty$, and by (33)–(36), it is easily seen that for all $\epsilon, \delta > 0$ there exist some $q(\epsilon, \delta), n(\epsilon, \delta) \in \mathbb{N}$ such that

$$(38) \quad \begin{aligned} P(|n\Delta_1 - n\Delta_1(q)| > \epsilon) &< \delta, & P(|n\Delta_{2,r} - n\Delta_{2,r}(q)| > \epsilon) &< \delta, \\ P(|n\Delta_{4,L} - n\Delta_{4,L}(q)| > \epsilon) &< \delta \end{aligned}$$

hold for all $q \geq q(\epsilon, \delta)$ and all $n \geq n(\epsilon, \delta)$. For any given q , we have $E(b_{i1}) = E(b_{i2}) = E(b_{i3}) = 0$, and it follows from Assumption 1 that the respective covariance structures can be represented by finite covariance matrices $\Omega_{1,q}$, $\Omega_{2,q}$ and $\Omega_{3,q}$. It therefore follows from our assumptions together with standard multivariate central limit theorems that the vectors $\{\frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} (b_{ik}^{(1)})^\top, \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} (b_{ik}^{(2)})^\top\}^\top$, $k = 1, 2, 3$, are asymptotically normal with zero means and covariance matrices $\Omega_{1,q}$, $\Omega_{2,q}$ and $\Omega_{3,q}$. One can thus conclude that, as $n \rightarrow \infty$,

$$(39) \quad n\Delta_1(q) \xrightarrow{\mathcal{L}} F_{1,q}, \quad n\Delta_{2,r}(q) \xrightarrow{\mathcal{L}} F_{2,r,q}, \quad n\Delta_{4,L}(q) \xrightarrow{\mathcal{L}} F_{4,L,q},$$

where $F_{1,q}, F_{2,r,q}, F_{4,L,q}$ denote the continuous distributions of the quadratic forms $z_1^\top Q_1^q z_1, z_2^\top Q_2^q z_2, z_3^\top Q_3^q z_3$ with $z_1 \sim N(0, \Omega_{1,q}), z_2 \sim N(0, \Omega_{2,q}), z_3 \sim$

$N(0, \Omega_{3,q})$. Since ϵ, δ are arbitrary, (38) implies

$$(40) \quad \lim_{q \rightarrow \infty} F_{1,q} = F_1, \quad \lim_{q \rightarrow \infty} F_{2,r,q} = F_{2,r}, \quad \lim_{q \rightarrow \infty} F_{4,L,q} = F_{4,L}.$$

We now have to consider the asymptotic properties of bootstrapped eigenvalues and eigenfunctions. Let $\bar{X}^{(p)*} = \frac{1}{n_p} \sum_{i=1}^{n_p} X_i^{(p)*}$, $\beta_{ri}^{(p)*} = \int \gamma_r^{(p)}(t) \{X_i^{(p)*}(t) - \mu(t)\}$, $\bar{\beta}_r^{(p)*} = \frac{1}{n_p} \sum_{i=1}^{n_p} \beta_{ri}^{(p)*}$, and note that $\int \gamma_r^{(p)}(t) \{X_i^{(p)*}(t) - \bar{X}^{(p)*}(t)\} = \beta_{ri}^{(p)*} - \bar{\beta}_r^{(p)*}$. When considering unconditional expectations, our assumptions imply that for $p = 1, 2$

$$(41) \quad \begin{aligned} \mathbb{E}[\beta_{ri}^{(p)*}] &= 0, & \mathbb{E}[(\beta_{ri}^{(p)*})^2] &= \lambda_r^{(p)}, \\ \mathbb{E}[(\bar{\beta}_r^{(p)*})^2] &= \frac{\lambda_r^{(p)}}{n_p}, & \mathbb{E}\{[(\beta_{ri}^{(p)*})^2 - \lambda_r^{(p)}]^2\} &= \Lambda_r^{(p)}, \\ \mathbb{E}\left\{ \sum_{l,k=1}^{\infty} \left[\frac{1}{n_p} \sum_{i=1}^{n_p} (\beta_{li}^{(p)*} - \bar{\beta}_l^{(p)*})(\beta_{ki}^{(p)*} - \bar{\beta}_k^{(p)*}) - \delta_{lk} \lambda_l^{(p)} \right]^2 \right\} \\ &= \frac{1}{n_p} \left(\sum_l \Lambda_l^{(p)} + \sum_{l \neq k} \lambda_l^{(p)} \lambda_k^{(p)} \right) + \mathcal{O}(n_p^{-1}). \end{aligned}$$

One can infer from (41) that the arguments used to prove Theorem 1 can be generalized to approximate the difference between the bootstrap eigenvalues and eigenfunctions $\hat{\lambda}_r^{(p)*}$, $\hat{\gamma}_r^{(p)*}$ and the true eigenvalues $\lambda_r^{(p)}$, $\gamma_r^{(p)}$. All infinite sums involved converge with probability 1. Relation (30) then generalizes to

$$(42) \quad \begin{aligned} &\sqrt{n_p}(\hat{\lambda}_r^{(p)*} - \hat{\lambda}_r^{(p)}) \\ &= \sqrt{n_p}(\hat{\lambda}_r^{(p)*} - \lambda_r^{(p)}) - \sqrt{n_p}(\hat{\lambda}_r^{(p)} - \lambda_r^{(p)}) \\ &= \frac{1}{\sqrt{n_p}} \sum_{i=1}^{n_p} (\beta_{ri}^{(p)*} - \bar{\beta}_r^{(p)*})^2 \\ &\quad - \frac{1}{\sqrt{n_p}} \sum_{i=1}^{n_p} (\beta_{ri}^{(p)} - \bar{\beta}_r^{(p)})^2 + \mathcal{O}_p(n_p^{-1/2}) \\ &= \frac{1}{\sqrt{n_p}} \sum_{i=1}^{n_p} \left\{ (\beta_{ri}^{(p)*})^2 - \frac{1}{n_p} \sum_{k=1}^{n_p} (\beta_{rk}^{(p)})^2 \right\} + \mathcal{O}_p(n_p^{-1/2}). \end{aligned}$$

Similarly, (32) becomes

$$(43) \quad \begin{aligned} &\hat{\gamma}_r^{(p)*} - \hat{\gamma}_r^{(p)} \\ &= \hat{\gamma}_r^{(p)*} - \gamma_r^{(p)} - (\hat{\gamma}_r^{(p)} - \gamma_r^{(p)}) \end{aligned}$$

$$\begin{aligned}
&= - \sum_{s \neq r} \left\{ \frac{1}{\lambda_s^{(p)} - \lambda_r^{(p)}} \frac{1}{n_p} \sum_{i=1}^{n_p} (\beta_{si}^{(p)*} - \bar{\beta}_s^{(p)*})(\beta_{ri}^{(p)*} - \bar{\beta}_r^{(p)*}) \right. \\
&\quad \left. - \frac{1}{\lambda_s^{(p)} - \lambda_r^{(p)}} \frac{1}{n_p} \sum_{i=1}^{n_p} (\beta_{si}^{(p)} - \bar{\beta}_s^{(p)})(\beta_{ri}^{(p)} - \bar{\beta}_r^{(p)}) \right\} \gamma_s^{(p)}(t) \\
&\quad + R_r^{(p)*}(t) \\
&= - \sum_{s \neq r} \left\{ \frac{1}{\lambda_s^{(p)} - \lambda_r^{(p)}} \frac{1}{n_p} \sum_{i=1}^{n_p} \left(\beta_{si}^{(p)*} \beta_{ri}^{(p)*} - \frac{1}{n_p} \sum_{k=1}^{n_p} \beta_{sk}^{(p)} \beta_{rk}^{(p)} \right) \right\} \gamma_s^{(p)}(t) \\
&\quad + \tilde{R}_r^{(p)*}(t),
\end{aligned}$$

where due to (28), (29) and (41), the remainder term satisfies $\|R_r^{(p)*}\| = \mathcal{O}_p(n_p^{-1})$.

We are now ready to analyze the bootstrap versions Δ^* of the different Δ . First consider $\Delta_{3,r}^*$ and note that $\{(\beta_{ri}^{(p)*})^2\}$ are i.i.d. bootstrap resamples from $\{(\beta_{ri}^{(p)})^2\}$. It therefore follows from basic bootstrap results that the conditional distribution of $\frac{1}{\sqrt{n_p}} \sum_{i=1}^{n_p} [(\beta_{ri}^{(p)*})^2 - \frac{1}{n_p} \sum_{k=1}^{n_p} (\beta_{rk}^{(p)})^2]$ given \mathcal{X}_p converges to the same $N(0, \Lambda_r^{(p)})$ limit distribution as $\frac{1}{\sqrt{n_p}} \sum_{i=1}^{n_p} [(\beta_{ri}^{(p)})^2 - E\{(\beta_{ri}^{(p)})^2\}]$. Together with the independence of $(\beta_{ri}^{(1)*})^2$ and $(\beta_{ri}^{(2)*})^2$, the assertion of the theorem is an immediate consequence.

Let us turn to Δ_1^* , $\Delta_{2,r}^*$ and $\Delta_{4,L}^*$. Using (41)–(43), it is then easily seen that $n\Delta_1^*$, $n\Delta_{2,r}^*$ and $n\Delta_{4,L}^*$ admit expansions similar to (33), (35) and (36), when replacing there $\frac{1}{\sqrt{n_p}} \sum_{i=1}^{n_p} \beta_{ri}^{(p)}$ by $\frac{1}{\sqrt{n_p}} \sum_{i=1}^{n_p} (\beta_{ri}^{(p)*} - \frac{1}{n_p} \sum_{k=1}^{n_p} \beta_{rk}^{(p)})$, as well as $\frac{1}{\sqrt{n_p}} \sum_{i=1}^{n_p} \beta_{si}^{(p)} \beta_{ri}^{(p)}$ by $\frac{1}{\sqrt{n_p}} \sum_{i=1}^{n_p} (\beta_{si}^{(p)*} \beta_{ri}^{(p)*} - \frac{1}{n_p} \sum_{k=1}^{n_p} \beta_{sk}^{(p)} \beta_{rk}^{(p)})$.

Replacing $\beta_{ri}^{(p)}$, $\beta_{si}^{(p)}$ by $\beta_{ri}^{(p)*}$, $\beta_{si}^{(p)*}$ leads to bootstrap analogs $b_{ik}^{(p)*}$ of the vectors $b_{ik}^{(p)}$, $k = 1, 2, 3$. For any $q \in \mathbb{N}$, define bootstrap versions $n\Delta_1^*(q)$, $n\Delta_{2,r}^*(q)$ and $n\Delta_{4,L}^*(q)$ of $n\Delta_1(q)$, $n\Delta_{2,r}(q)$ and $n\Delta_{4,L}(q)$ by using $(\frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} (b_{ik}^{(1)*} - \frac{1}{n_1} \sum_{k=1}^{n_1} b_{ik}^{(1)})^\top, \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} (b_{ik}^{(2)*} - \frac{1}{n_2} \sum_{k=1}^{n_2} b_{ik}^{(2)})^\top)$ instead of $(\frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} (b_{ik}^{(1)})^\top, \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} (b_{ik}^{(2)})^\top)$, $k = 1, 2, 3$, in (37). Applying again (41)–(43), one can conclude that for any $\epsilon > 0$ there exists some $q(\epsilon)$ such that, as $n \rightarrow \infty$,

$$\begin{aligned}
(44) \quad &P(|n\Delta_1^* - n\Delta_1^*(q)| < \epsilon) \rightarrow 1, \\
&P(|n\Delta_{2,r}^* - n\Delta_{2,r}^*(q)| < \epsilon) \rightarrow 1, \\
&P(|n\Delta_{4,L}^* - n\Delta_{4,L}^*(q)| < \epsilon) \rightarrow 1
\end{aligned}$$

hold for all $q \geq q(\epsilon)$. Of course, (44) generalizes to the conditional probabilities given $\mathcal{X}_1, \mathcal{X}_2$.

In order to prove the theorem, it thus only remains to show that for *any* given q and all δ

$$(45) \quad |\mathbb{P}(n\Delta(q) \geq \delta) - \mathbb{P}(n\Delta^*(q) \geq \delta) | \mathcal{X}_1, \mathcal{X}_2| = o_p(1)$$

hold for either $\Delta(q) = \Delta_1(q)$ and $\Delta^*(q) = \Delta_1^*(q)$, $\Delta(q) = \Delta_{2,r}(q)$ and $\Delta^*(q) = \Delta_{2,r}^*(q)$, or $\Delta(q) = \Delta_{4,L}(q)$ and $\Delta^*(q) = \Delta_{4,L}^*(q)$. But note that for $k = 1, 2, 3$, $\mathbb{E}(b_{ik}) = 0$, $\{b_{ik}^{(j)*}\}$ are i.i.d. bootstrap resamples from $\{b_{ik}^{(p)}\}$, and $\mathbb{E}(b_{ik}^{(p)*} | \mathcal{X}_1, \mathcal{X}_2) = \frac{1}{n_p} \sum_{k=1}^{n_p} b_{ik}^{(p)}$ are the corresponding conditional means. It therefore follows from basic bootstrap results that as $n \rightarrow \infty$ the conditional distribution of $(\frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} (b_{ik}^{(1)*} - \frac{1}{n_1} \sum_{k=1}^{n_1} b_{ik}^{(1)})^\top, \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} (b_{ik}^{(2)*} - \frac{1}{n_2} \sum_{k=1}^{n_2} b_{ik}^{(2)})^\top)$ given $\mathcal{X}_1, \mathcal{X}_2$ converges to the same $N(0, \Omega_{k,q})$ limit distribution as $(\frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} (b_{ik}^{(1)})^\top, \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} (b_{ik}^{(2)})^\top)$. This obviously holds for all $q \in \mathbb{N}$, and (45) is an immediate consequence. The theorem then follows from (38), (39), (40), (44) and (45). \square

REFERENCES

- ARAÚJO, A. and GINÉ, E. (1980). *The Central Limit Theorem for Real and Banach Valued Random Variables*. Wiley, New York. [MR0576407](#)
- BESSE, P. and RAMSAY, J. (1986). Principal components of sampled functions. *Psychometrika* **51** 285–311. [MR0848110](#)
- BLACK, F. and SCHOLES, M. (1973). The pricing of options and corporate liabilities. *J. Political Economy* **81** 637–654.
- DAUXOIS, J., POUSSE, A. and ROMAIN, Y. (1982). Asymptotic theory for the principal component analysis of a vector random function: Some applications to statistical inference. *J. Multivariate Anal.* **12** 136–154. [MR0650934](#)
- FENGLER, M. (2005a). Arbitrage-free smoothing of the implied volatility surface. SFB 649 Discussion Paper No. 2005–019, SFB 649, Humboldt-Universität zu Berlin.
- FENGLER, M. (2005b). *Semiparametric Modeling of Implied Volatility*. Springer, Berlin. [MR2183565](#)
- FENGLER, M., HÄRDLE, W. and VILLA, P. (2003). The dynamics of implied volatilities: A common principle components approach. *Rev. Derivative Research* **6** 179–202.
- FENGLER, M., HÄRDLE, W. and MAMMEN, E. (2007). A dynamic semiparametric factor model for implied volatility string dynamics. *Financial Econometrics* **5** 189–218.
- FERRATY, F. and VIEU, P. (2006). *Nonparametric Functional Data Analysis*. Springer, New York. [MR2229687](#)
- FLURY, B. (1988). *Common Principal Components and Related Models*. Wiley, New York. [MR0986245](#)
- GIHMAN, I. I. and SKOROHOD, A. V. (1973). *The Theory of Stochastic Processes. II*. Springer, New York. [MR0375463](#)
- HALL, P. and HOSSEINI-NASAB, M. (2006). On properties of functional principal components analysis. *J. Roy. Statist. Soc. Ser. B* **68** 109–126. [MR2212577](#)
- HALL, P., MÜLLER, H. G. and WANG, J. L. (2006). Properties of principal components methods for functional and longitudinal data analysis. *Ann. Statist.* **34** 1493–1517. [MR2278365](#)

- HALL, P., KAY, J. W. and TITTERINGTON, D. M. (1990). Asymptotically optimal difference-based estimation of variance in nonparametric regression. *Biometrika* **77** 520–528. [MR1087842](#)
- HAFNER, R. (2004). *Stochastic Implied Volatility*. Springer, Berlin. [MR2090447](#)
- HÄRDLE, W. and SIMAR, L. (2003). *Applied Multivariate Statistical Analysis*. Springer, Berlin. [MR2061627](#)
- KAHALÉ, N. (2004). An arbitrage-free interpolation of volatilities. *Risk* **17** 102–106.
- KNEIP, A. and UTIKAL, K. (2001). Inference for density families using functional principal components analysis. *J. Amer. Statist. Assoc.* **96** 519–531. [MR1946423](#)
- LACANTORE, N., MARRON, J. S., SIMPSON, D. G., TRIPOLI, N., ZHANG, J. T. and COHEN, K. L. (1999). Robust principal component analysis for functional data. *Test* **8** 1–73. [MR1707596](#)
- PEZZULLI, S. D. and SILVERMAN, B. (1993). Some properties of smoothed principal components analysis for functional data. *Comput. Statist.* **8** 1–16. [MR1220336](#)
- RAMSAY, J. O. and DALZELL, C. J. (1991). Some tools for functional data analysis (with discussion). *J. Roy. Statist. Soc. Ser. B* **53** 539–572. [MR1125714](#)
- RAMSAY, J. and SILVERMAN, B. (2002). *Applied Functional Data Analysis*. Springer, New York. [MR1910407](#)
- RAMSAY, J. and SILVERMAN, B. (2005). *Functional Data Analysis*. Springer, New York. [MR2168993](#)
- RAO, C. (1958). Some statistical methods for comparison of growth curves. *Biometrics* **14** 1–17.
- RICE, J. and SILVERMAN, B. (1991). Estimating the mean and covariance structure non-parametrically when the data are curves. *J. Roy. Statist. Soc. Ser. B* **53** 233–243. [MR1094283](#)
- SILVERMAN, B. (1996). Smoothed functional principal components analysis by choice of norm. *Ann. Statist.* **24** 1–24. [MR1389877](#)
- TYLER, D. E. (1981). Asymptotic inference for eigenvectors. *Ann. Statist.* **9** 725–736. [MR0619278](#)
- YAO, F., MÜLLER, H. G. and WANG, J. L. (2005). Functional data analysis for sparse longitudinal data. *J. Amer. Statist. Assoc.* **100** 577–590. [MR2160561](#)

M. BENKO
W. HÄRDLE
CASE—CENTER FOR APPLIED STATISTICS AND ECONOMICS
HUMBOLDT-UNIVERSITÄT ZU BERLIN
SPANDAUERSTR 1
D-10178 BERLIN
GERMANY
E-MAIL: benko@wiwi.hu-berlin.de
haerdle@wiwi.hu-berlin.de
URL: <http://www.case.hu-berlin.de/>

A. KNEIP
STATISTISCHE ABTEILUNG
DEPARTMENT OF ECONOMICS
UNIVERSITÄT BONN
ADENAUERALLEE 24-26
D-53113 BONN
GERMANY
E-MAIL: akneip@uni-bonn.de

GHICA - Risk Analysis with GH Distributions and Independent Components

Ying Chen^{1,2}, Wolfgang Härdle¹ and Vladimir Spokoiny^{1,2}

¹ CASE - Center for Applied Statistics and Economics

Humboldt-Universität zu Berlin

Wirtschaftswissenschaftliche Fakultät

Spandauerstrasse 1, 10178 Berlin, Germany

² Weierstraß - Institute für Angewandte Analysis und Stochastik

Mohrenstrasse 39, 10117 Berlin, Germany

Abstract

Over recent years, study on risk management has been prompted by the Basel committee for regular banking supervisory. There are however limitations of some widely-used risk management methods that either calculate risk measures under the Gaussian distributional assumption or involve numerical difficulty. The primary aim of this paper is to present a realistic and fast method, **GHICA**, which overcomes the limitations in multivariate risk analysis. The idea is to first retrieve independent components (ICs) out of the observed high-dimensional time series and then individually and adaptively fit the resulting ICs in the generalized hyperbolic (GH) distributional framework. For the volatility estimation of each IC, the local exponential smoothing technique is used to achieve the best possible accuracy of estimation. Finally, the fast Fourier transformation technique is used to approximate the density of the portfolio returns.

The proposed GHICA method is applicable to covariance estimation as well. It is compared with the dynamic conditional correlation (DCC) method based on the simulated data with $d = 50$ GH distributed components. We further implement the GHICA method to calculate risk measures given 20-dimensional German DAX portfolios and a dynamic exchange rate portfolio. Several alternative methods are considered as well to compare the accuracy of calculation with the GHICA one.

Keywords: multivariate risk management, independent component analysis, generalized hyperbolic distribution, local exponential estimation, value at risk, expected shortfall

JEL Codes: C14, C16, C32, C61, G20

Acknowledgement: This research was supported by the Deutsche Forschungsgemeinschaft through the SFB 649 "Economic Risk". In the numerical analysis, the Matlab DCC function developed by Kevin Sheppard and the Matlab FastICA function developed by Aapo Hyvärinen are used.

1 Introduction

Over recent years, study on risk management has been prompted by the Basel committee for regular banking supervisory. Given a d -dimensional portfolio, the conditionally heteroscedastic model is widely used to describe the movement of the underlying series:

$$x(t) = \Sigma_x^{1/2}(t)\varepsilon_x(t), \quad (1)$$

where $x(t) \in \mathbb{R}^d$ are risk factors of the portfolio, e.g. (log) returns of the financial instruments. The covariance Σ_x is assumed to be predictable with respect to (w.r.t.) the past information and $\varepsilon_x(t) \in \mathbb{R}^d$ is a sequence of standardized innovations with $\mathbb{E}[\varepsilon_x(t)|\mathcal{F}_{t-1}] = 0$ and $\mathbb{E}[\varepsilon_x^2(t)|\mathcal{F}_{t-1}] = I_d$. There is a sizeable literature on risk management methods. Among others, we refer to Jorion (2001) for a systematic description.

In this paper, we focus on the calculation of two risk measures, value at risk (VaR) and expected shortfall (ES). These two risk measures are inherently related to the joint density of $x(t)$. The VaR is in fact the distributional quantile of loss, i.e. $-x(t)$, at a prescribed level over a target time horizon and the ES measures the size of loss once the loss exceeds the VaR value. Indicated by formula (1), the joint density estimation depends on the covariance estimation and the distributional assumption of the innovations.

The largest challenge of risk management is due to the high-dimensionality of real portfolios. Above all, the covariance estimation is really computationally demanding as high dimensional series, e.g. a dimension $d > 10$, is considered, see Härdle, Herwartz and Spokoiny (2003). For example, the dynamic conditional correlation (DCC) model proposed by Engle (2002), Engle and Sheppard (2001), which is one multivariate GARCH model, is recommended due to the good performance of its univariate version. In the estimation, the covariance matrix is approximated by the product of a diagonal matrix and a correlation matrix, which reduces the number of unknown parameters much relative to the BEKK specification proposed by Engle and Kroner (1995). In spite of the appealing dimensional reduction, the mentioned estimation method is time consuming and numerically difficult to handle given high-dimensional data.

Moreover, many widely-used risk management methods rely on the unrealistic Gaussian distributional assumption, e.g. the RiskMetrics product introduced by JP Morgan in 1994. In the Gaussian framework with an estimate $\hat{\Sigma}_x(t)$ of $\Sigma_x(t)$, the standardized returns $\hat{\varepsilon}_x(t) = \hat{\Sigma}_x^{-1/2}(t)x(t)$ are asymptotically independent and the joint distributional behavior can be easily measured by the marginal distributions. However the Gaussian distributional assumption is merely used for computational and numerical purposes and not for statistical reasons. The conditional Gaussian marginal distributions and the resulting joint Gaussian distribution are at odds with empirical facts, i.e. financial series are heavy tailed distributed.

The heavy tails are typically reduced but not eliminated as the series are standardized by the estimated volatility, see Anderson, Bollerslev, Diebold and Labys (2001).

We illustrate this effect based on two real data sets, the Allianz stock and a DAX portfolio from 1988/01/04 to 1996/12/30. The DAX is the leading index of Frankfurt stock exchange and a 20-dimensional hypothetic portfolio with a static trading strategy $b(t) = (1/20, \dots, 1/20)^\top$ is considered. The portfolio returns $r(t) = b(t)^\top x(t)$ are analyzed in the univariate version of (1). This simplified calculation is used in practice, but it often suffers from low accuracy of calculation. Suppose now that the two return processes have been properly standardized, by using a local volatility estimation technique discussed later. The standardized returns are empirically heavy-tailed distributed, indicated by the sample kurtoses 12.07 for the Allianz and 22.38 for the portfolio respectively.

Figure 1 displays the estimated logarithmic density curves under several distributional assumptions. Among them, the estimate using the nonparametric kernel estimation is considered as benchmark. The comparison w.r.t. the Allianz stock shows that the GH estimate is most close to the benchmark among others. The Gaussian estimate presents lighter tails. To alleviate the limitation, the Student- $t(6)$ distribution with degrees of freedom of 6 has been recommended in practice. However this distribution is found to over-fit the heavy tails, namely the $t(6)$ estimate displays heavier tails relative to the benchmark. The similar result is observed w.r.t. the DAX portfolio. It is rational to surmise that the risk management methods under the Gaussian and $t(6)$ distributional assumptions generate low accurate results.

To overcome these limitations, Chen, Härdle and Spokoiny (2006) present a simple VaR calculation approach that achieves much better accuracy than the alternative RiskMetrics method. In their study, univariate approaches that involve more realistic but complex procedures can be easily extended for multivariate risk measurement. To be more specific, financial risk factors are first converted to independent components (ICs) using a linear filtering and the univariate method is applied to identify the distributional behavior of each IC. We name here two univariate approaches which measure the risk exposure in the realistic distributional framework. One is the univariate VaR calculation proposed by Chen, Härdle and Jeong (2005), which implements local constant model to estimate volatility and fit the standardized returns under the GH distributional assumption. The other is proposed by Chen and Spokoiny (2006), who apply the local exponential smoothing method to estimate volatility and calculate the risk measure in the GH distributional framework. The standardization of the Allianz and DAX returns in Figure 1 is in fact based on the local exponential smoothing technique.

The primary aim of this paper is to present an realistic and fast multivariate risk management method, **GHICA**, by implementing the IC analysis (ICA) to the high dimensional series and adaptively fitting the ICs in the GH distributional framework. The GHICA

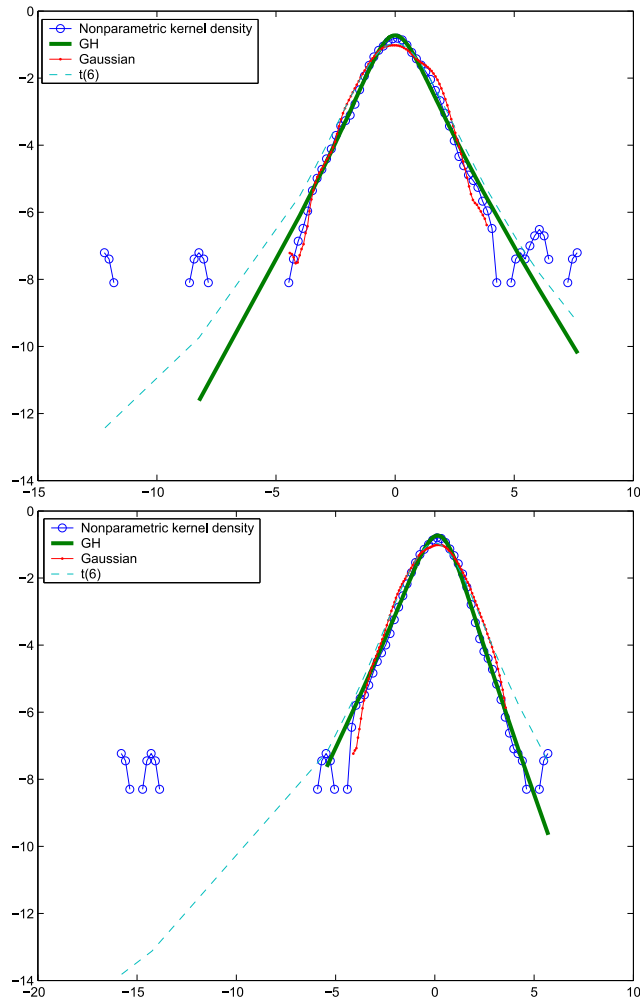


Fig. 1: Density comparisons of the standardized returns in log scale based on the Allianz stock (top) and the DAX portfolio (bottom) with static weights $b(t) = \text{unit}(1/20)$. Time interval: 1988/01/04 - 1996/12/30. The nonparametric kernel density is considered as benchmark. The GH distributional parameters are respectively $\text{GH}(-0.5, 1.01, 0.05, 1.11, -0.03)$ for the Allianz and $\text{GH}(-0.5, 1.21, -0.21, 1.21, 0.24)$ for the DAX portfolio. Data source: FEDC (<http://sfb649.wiwi.hu-berlin.de>).

method improves the work of Chen et al. (2006) from two aspects. The volatility estimation is driven by the local exponential smoothing technique to achieve the best possible accuracy of estimation. The fast Fourier transformation (FFT) technique is used to approximate the density of the portfolio returns. Compared to the Monte Carlo simulation technique used in the former study, it significantly speeds up the calculation.

In addition, the proposed GHICA method is easily applicable for covariance estimation. Relative to the widely used DCC setup, the GHICA method is fast and delivers sensitive estimates. We demonstrate the comparison based on simulated data. Furthermore, the

GHICA method is implemented to risk management on the base of DAX stocks and foreign exchange rates. Several hypothetic portfolios are constructed by assigning static and dynamic trading strategies to the data sets. The results are compared with those calculated using alternative methods, i.e. the RiskMetrics method, the method using the exponential smoothing to estimate volatility and assuming the Student- $t(6)$ distribution, and the method using the DCC to estimate covariance in the Gaussian distributional framework. All the results are analyzed from the viewpoints of regulatory, investors and internal supervisory. The GHICA method, in general, produces better results than the others.

The paper is organized as follows. The GHICA method is described in Section 2, by which the ICA method, the local exponential smoothing technique and the FFT technique are detailed. Section 3 compares the covariance estimation using the GHICA and DCC methods based on the simulated data with $d = 50$ GH components. The real data analysis in Section 4 demonstrates the implementation of the GHICA method in risk management based on the 20-dimensional German DAX portfolios and a dynamic exchange rate portfolio. Several alternative methods are considered as well to compare the accuracy of calculation with the GHICA one.

2 GHICA Methodology

Given multidimensional time series, for example prices of financial assets, $s(t) \in \mathbb{R}^d$, the (log) returns are calculated as $x(t) = \log\{s(t)/s(t-1)\}$. Without loss of generality, the drift of the returns is set to be 0. Given the time homogeneous model, $x(t) = \Sigma_x^{1/2} \varepsilon_x(t)$ with standardized innovations $\varepsilon_x(t)$, the maximum Gaussian likelihood estimate of the time independent covariance Σ_x is the sample covariance based on the whole past information. Since the covariance is in fact time dependent, one considers the conditional heteroscedastic model:

$$x(t) = \Sigma_x^{1/2}(t) \varepsilon_x(t).$$

Many techniques have been used to approximate the local covariance by specifying a “local homogeneous” interval (e.g. one year or 250 trading days). Inside the homogeneous interval, the unknown covariance should be time-invariant and can be identified using the ML estimation. Among many others, the multivariate GARCH setup such as the DCC is successful in characterizing the clustering feature of covariance under the Gaussian distributional assumption. As the dimension d increases, it however needs to estimate many parameters and becomes numerically difficult. Moreover, the standardized returns $\hat{\varepsilon}_x(t) = \hat{\Sigma}_x^{-1/2}(t)x(t)$ are empirically not Gaussian distributed. Under a realistic distributional assumption, on the other hand, by which the distributional behaviors such as asymmetry and heavy tails are well matched, it is hard to identify the unknown distributional parameters due to complex density form.

The GHICA method proposes a solution to balance the numerical tractability and the realistic distributional assumption on the risk factors. It first converts the return series using a linear transformation and filters out ICs: $y(t) = Wx(t)$. The transformation matrix W is assumed to be time constant and nonsingular and $y(t)$ is the independent vector. The heteroscedastic model is now reformulated as:

$$x(t) = W^{-1}y(t) = W^{-1}\Sigma_y^{1/2}(t)\varepsilon_y(t) = W^{-1}D_y^{1/2}(t)\varepsilon_y(t).$$

Due to the statistical property of independence, the covariance of the ICs $\Sigma_y(t)$ is a diagonal matrix and is denoted as $D_y(t)$ to emphasize this feature. Its diagonal elements are the time varying variances of the ICs. The stochastic innovations $\varepsilon_y(t) = \{\varepsilon_{y_1}(t), \dots, \varepsilon_{y_d}(t)\}^\top$ are cross independent and can be individually identified in the realistic and univariate distributional framework. By doing so, the GHICA method converts the high dimensional analysis to univariate study and significantly speeds up the calculation.

In this section, the building blocks of the GHICA method are detailed: The FastICA procedure is used to estimate the transformation matrix W ; The resulting ICs are individually analyzed, by which the univariate volatility process is estimated using the local exponential smoothing technique and the innovations are assumed to be GH distributed; The quantile of the portfolio return is approximated using the FFT technique.

The GHICA algorithm is summarized as follows:

1. Do ICA to the given risk factors to get ICs.
2. Implement local exponential smoothing to estimate the variance of each IC
3. Identify the distribution of every IC's innovation in the GH distributional framework
4. Estimate the density of the portfolio return using the FFT technique
5. Calculate risk measures

In addition, the GHICA method can be used to estimate the covariance matrix $\Sigma_x(t)$. Given the matrix estimate \hat{W} in the ICA and the variance estimates of the ICs, the covariance of the observed time series are: $\hat{\Sigma}_x(t) = \hat{W}^{-1}\hat{D}_y(t)\hat{W}^{-1\top}$. An alternative covariance estimation approach, the DCC, is briefly described as well. We will compare the GHICA-based covariance estimation with the DCC estimation in the later simulation study.

2.1 Independent component analysis (ICA) and FastICA approach

The aim of ICA is to retrieve, out of high dimensional time series, stochastically ICs through a linear transformation: $y(t) = Wx(t)$, where the transformation matrix $W = (w_1, \dots, w_d)^\top$

is nonsingular. It is essential to use high order moments in the ICA. In the Gaussian framework, high order moments are however fixed such as skewness with value of 0 and kurtosis with value of 3. Therefore the ICs are assumed to be nongaussian distributed. Furthermore, the ICA transformation has scale identification problem, i.e. the equation holds true by simultaneously multiplying the same constants to the unknown terms $y(t)$ and W : $\{cy(t)\} = \{cW\}x(t)$. To avoid this problem, it is natural to standardize the dependent series and assume that every IC has unit variance $\mathbf{E}(y_j) = 1$ with $j = 1, \dots, d$. The Mahalanobis transformation $\tilde{x}(t) = \tilde{\Sigma}_x^{-1/2}x(t)$ helps to standardize the return series and the resulting series are considered:

$$y(t) = \tilde{W}\tilde{x}(t),$$

where $\tilde{\Sigma}_x$ is the sample covariance based on the available data. It is easy to show that after the standardization the transformation matrix \tilde{W} turns to be an orthogonal matrix with unit norm. The corresponding matrix w.r.t. the return series is $W = \tilde{W}\tilde{\Sigma}_x^{-1/2}$. For notational simplification, we eliminate the mark $\tilde{\cdot}$ in the following text in this section.

Various ideas have been proposed to estimate the transformation matrix W . Among others, one intuitive ICA estimation is motivated by the definition of mutual information. The mutual information is a natural measure of independence. It is defined as the difference of the sum of marginal entropy and the mutual entropy:

$$I(y) = \sum_{j=1}^d H(y_j) - H(y) \quad (2)$$

where $H(y_j) = - \int f_{y_j}(u) \log f_{y_j}(u) du$

The mutual information is nonnegative and goes to 0 if the vector y is cross independent, see Cover and Thomas (1991). Hence for a candidate transformation W , one can minimize the mutual information to achieve independence. Based on the linear transformation of the ICA, the mutual information in (2) can be reformulated as:

$$I(W, y) = \sum_{j=1}^d H(y_j) - H(x) - \log |\det(W)|.$$

Notice that the entropy of the return series $H(x)$ is a fixed value and does not depend on the ICs, and the last term in the equation is 0 due to the orthogonality of the transformation matrix W . The optimization problem is: $\min_W \sum_{j=1}^d H(y_j)$ and can be further simplified to d optimization problems according to the inequality:

$$\min_W \sum_{j=1}^d H(y_j) \geq \sum_{j=1}^d \min_{w_j} H(y_j)$$

This simplification leads to some loss in the W estimation but it extensively speeds up the estimation procedure by merely considering d elements of W every time. Equivalently, one can formulate the optimization problem concerning negentropy $J(y_j) = H(y_0) - H(y_j)$ since the entropy and the negentropy are in one-to-one correspondence, where $y_0 \sim \mathcal{N}(0, 1)$ is a standard Gaussian vector and $H(y_0)$ is merely a constant. The negentropy is always nonnegative since the Gaussian random variable has the largest entropy given the same variance, see Hyvärinen (1998).

$$\hat{w}_j = \operatorname{argmin} H(y_j) = \operatorname{argmax} J(w_j, y_j).$$

In the estimation, the approximation of negentropy is used to construct the optimization object function w.r.t. the j -th row of the transformation matrix W :

$$\begin{aligned} \hat{w}_j &= \operatorname{argmin} H(y_j) = \operatorname{argmax} J(y_j) \\ J(y_j) &\approx \operatorname{const.} \{ \mathbb{E}[G(y)] - \mathbb{E}[G(y_0)] \}^2 \\ &= \operatorname{const.} \{ \mathbb{E}[G(w_j^\top x)] - \mathbb{E}[G(y_0)] \}^2 \\ G(y_j) &= \log \cosh(y_j) \end{aligned} \tag{3}$$

This optimization problem is solved by using the symmetric FastICA algorithm, see Hyvärinen, Karhunen and Oja (2001):

1. Initialization: Choose initial vectors $\hat{w}_j^{(1)}$ for $W = \{w_1, \dots, w_d\}^\top$ with $j = 1, \dots, d$, each has a unit norm.

2. Loop:

- At step n , Calculate $\hat{w}_j^{(n)} = \mathbb{E} \left[x^\top(t) g \left\{ \hat{w}_j^{(n-1)\top} x(t) \right\} \right] - \mathbb{E} \left[g' \left\{ \hat{w}_j^{(n-1)\top} x(t) \right\} \right] \hat{w}_j^{(n-1)}$, where g is the first derivative of $G(y)$ in form (3) and g' is the second derivative. The expectation $\mathbb{E}[\cdot]$ is approximated by the sample mean.
- Do a symmetric orthogonalization of the estimated transformation matrix $\hat{W}^{(n)}$:

$$\hat{W}^{(n)} = \{ \hat{W}^{(n)} \hat{W}^{(n)\top} \}^{-1/2} \hat{W}^{(n)}$$

- If not converged, i.e. $\det\{\hat{W}^{(n)} - \hat{W}^{(n-1)}\} \neq 0$, go back to 2. Otherwise, the algorithm stops.

3. Final result: the last (converged) estimate is the final estimate \hat{W} .

2.2 Local exponential smoothing and dynamically conditional correlation

Suppose that the ICs and the transformation matrix W are given. The covariance matrices of the ICs and the original return series are respectively:

$$\begin{aligned} D_y(t) &= \text{diag}\{\sigma_{y_1}^2(t), \dots, \sigma_{y_d}^2(t)\} \\ \Sigma_x(t) &= W^{-1}D_y(t)W^{-1\top} \end{aligned} \quad (4)$$

where $\sigma_{y_j}(t)$ is the heteroscedastic volatility of the j -th IC with $j = 1, \dots, d$. Recall that (4) has a similar decomposition structure as the often-used principal component analysis (PCA), by which the covariance is decomposed as: $\Sigma_x = \Gamma\Lambda\Gamma^\top$ with the eigenvector matrix Γ and the diagonal eigenvalue matrix Λ , see Flury (1998). Among other distinctions, the PCA method orders the resulting PCs whereas the ICs have equal importance. In the estimation of the unknown variance, the local exponential smoothing method is used.

Local exponential smoothing: Given the univariate conditional heteroscedastic model: $y_j(t) = \sigma_{y_j}(t)\varepsilon_{y_j}(t)$ with $\mathbb{E}[\varepsilon_{y_j}(t)|\mathcal{F}_{t-1}] = 0$ and $\mathbb{E}[\varepsilon_{y_j}^2(t)|\mathcal{F}_{t-1}] = 1$, we now focus on the adaptive estimation of the volatility σ_{y_j} for $j = 1, \dots, d$. For notational simplification, the subscripts y_j in σ_{y_j} and j in y_j are eliminated here.

Suppose that a finite set $\{\eta_k, k = 1, \dots, K\}$ of values of smoothing parameter is given. Every value η_k leads to a localizing weighting scheme $\{\eta_k^{t-s}\}$ for $s \leq t$ to the local Gaussian MLE $\tilde{\sigma}^{(k)}(t)$

$$\tilde{\sigma}^{(k)}(t) = \left[\frac{\sum_{m=0}^{\infty} \eta_k^m y^2(t-m-1)}{\sum_{m=0}^{\infty} \eta_k^m} \right]^{1/2}$$

In practice, one truncates the smoothing window at M_k such that $\eta_k^{M_k+1} \leq c \rightarrow 0$:

$$\tilde{\sigma}^{(k)}(t) = \left[\frac{\sum_{m=0}^{M_k} \eta_k^m y^2(t-m-1)}{\sum_{m=0}^{M_k} \eta_k^m} \right]^{1/2}$$

where the Gaussian log-likelihood function given η_k is:

$$\begin{aligned} L(\eta_k, \tilde{\sigma}^{(k)}(t)) &= -\frac{N_k}{2} \log(2\pi\{\sigma^{(k)}(t)\}^2) - \frac{1}{2\{\sigma^{(k)}(t)\}^2} \sum_{m=0}^{M_k} \eta_k^m y^2(t-m-1) \\ \text{where } N_k &= \sum_{m=0}^{M_k} \eta_k^m \end{aligned} \quad (5)$$

The fitted log-likelihood ratio $L(\eta_k, \tilde{\sigma}^{(k)}(t), \sigma(t))$ reads as:

$$L(\eta_k, \tilde{\sigma}^{(k)}(t), \sigma(t)) = L(\eta_k, \tilde{\sigma}^{(k)}(t)) - L(\eta_k, \sigma(t))$$

The idea of local exponential smoothing is to aggregate all the local likelihood estimate to achieve the best possible accuracy of estimation. In this sense, the local MLEs $\tilde{\sigma}^{(k)}(t)$ are referred as “weak” estimates.

In our study, we concern the heavy-tailedness of financial time series and assume the normal inverse Gaussian (NIG) distribution, one subclass of the GH distribution, see Section 2.3 for more details. Since the NIG distributional parameters of the innovations are unknown at this stage, we use the quasi ML estimation instead of estimating the variance based on the NIG density form. The quasi ML estimation is applicable if the exponential moment of the squared innovations $\mathbb{E}[\exp\{\rho\varepsilon^2(t)\}]$ exists. A power transformation guarantees that:

$$\begin{aligned} y_p(t) &= \text{sign}\{y(t)\}|y(t)|^p \\ \theta(t) &= \text{Var}\{y_p(t)|\mathcal{F}_{t-1}\} = \mathbb{E}\{y_p^2(t)|\mathcal{F}_{t-1}\} = \mathbb{E}\{|y(t)|^{2p}|\mathcal{F}_{t-1}\} \\ &= \sigma^{2p}(t) \mathbb{E}|\varepsilon(t)|^{2p} = \sigma^{2p}(t)C_p \end{aligned} \quad (6)$$

where $C_p = \mathbb{E}(|\varepsilon(t)|^{2p}|\mathcal{F}_{t-1})$ is a constant and only relies on $0 \leq p < 1/2$. Notice that the power transformed variable $\theta(t)$ is one-to-one correspondence to the variance $\sigma^2(t)$ and can be estimated on the base of the transformed observations $|y(t)|^{2p}$:

$$\tilde{\theta}^{(k)}(t) = \left\{ \sum_{m=0}^{M_k} \eta_k^m |y(t-m-1)|^{2p} \right\} / N_k$$

Here the smoothing parameter η_k is designed to run over a wide range from values close to zero to one, so that the variability of the unknown process $\theta(t)$ reduces and at least one of the resulting MLEs is good in the sense of small estimation bias. Polzehl and Spokoiny (2006) show that the inverse of N_k in (5) is positively related to the variation of the MLEs. This result is used to construct the sequence of the smoothing parameter $\{\eta_k\}$:

$$\frac{N_{k+1}}{N_k} \approx \frac{1 - \eta_k}{1 - \eta_{k+1}} = a > 1, \quad (7)$$

where the coefficient a controls the decreasing speed of the variations.

The procedure is sequential and starts with the estimate $\tilde{\theta}^{(1)}(t)$ that has the largest variability but small bias, i.e. we set $\hat{\theta}^{(1)}(t) = \tilde{\theta}^{(1)}(t)$. At every step $k \geq 2$, the new estimate $\hat{\theta}^{(k)}(t)$ is constructed by aggregating the next “weak” estimate $\tilde{\theta}^{(k)}(t)$ and the previously constructed estimate $\hat{\theta}^{(k-1)}(t)$. Following to Belomestny and Spokoiny (2006), the aggregation is done in terms of the parameter $v = -1/(2\theta)$ so that the variable $y(t)$

belongs to the exponential distributional family with a density form: $p(y, v) = p(y) \exp\{yv - d(v)\}$:

$$\begin{aligned} \hat{v}^{(k)}(t) &= \gamma_k \tilde{v}^{(k)}(t) + (1 - \gamma_k) \hat{v}^{(k-1)}(t) \\ \text{or equivalently, } \hat{\theta}^{(k)}(t) &= \left(\frac{\gamma_k}{\tilde{\theta}^{(k)}(t)} + \frac{1 - \gamma_k}{\hat{\theta}^{(k-1)}(t)} \right)^{-1} \end{aligned}$$

The mixing weights $\{\gamma_k\}$ are computed on the base of the fitted log-likelihood ratio by checking that the previously accepted estimate $\hat{\theta}^{(k-1)}(t)$ is in agreement with the next “weak” estimate $\tilde{\theta}^{(k)}(t)$, i.e. the difference between these two estimates is bounded by critical values δ_k :

$$\gamma_k = K_{ag} \left\{ L \left(\eta_k, \tilde{\theta}^{(k)}(t), \hat{\theta}^{(k-1)}(t) \right) / \delta_k \right\}$$

The aggregation kernel K_{ag} guarantees that the mixing coefficient γ_k is one if there is no essential difference between $\tilde{\theta}^{(k)}(t)$ and $\hat{\theta}^{(k-1)}(t)$, and zero if the difference is significant. The significance level is measured by the critical value ζ_k . In the intermediate case, the mixing coefficient γ_k is between zero and one. The procedure terminates after step k if $\gamma_k = 0$ and we define in this case $\hat{\theta}^{(m)}(t) = \hat{\theta}^{(k-1)}(t)$ for all $m \geq k$.

The critical values $\{\zeta_k\}$ are calculated by using Monte Carlo simulation. We briefly summarize the procedure here. Since the NIG distributional parameters of the innovations are unknown and the transformed variable is close to Gaussian variable, we start from the Gaussian assumption. To be more specific, we generate $y(t) = \sigma^* \varepsilon(t)$ with $\varepsilon(t) \sim N(0, 1)$ and $\sigma^* \stackrel{\text{def}}{=} 1$. The “weak” estimates are calculated given the sequence of $\{\eta_k\}$. For $k = 2, \dots, K$ with $\zeta_1, \infty, \dots, \infty$, the value ζ_1 is selected as the minimal one to fulfill

$$\mathbb{E}_{\theta^*} |L \left(\eta_k, \tilde{\theta}^{(k)}(t), \hat{\theta}_{\zeta_1}^{(k)}(t) \right)|^r \leq \frac{\alpha \tau_r}{K - 1}, \quad (8)$$

where $\tau_r = 2r \int_{\zeta \geq 0} \zeta^{r-1} e^{-\zeta} d\zeta = 2r\Gamma(r)$, and $r = 0.5$ and $\alpha = 1$ have been suggested in Chen and Spokoiny (2006). Consequently for $l = k + 1, \dots, K$ with the parameters $\zeta_1, \dots, \zeta_k, \infty, \dots, \infty$, we select ζ_k as the minimal value which fulfills

$$\mathbb{E}_{\theta^*} |L \left(\eta_l, \tilde{\theta}^{(l)}(t), \hat{\theta}_{\zeta_1, \dots, \zeta_k}^{(l)}(t) \right)|^r \leq \frac{k\alpha\tau_r}{K - 1}. \quad (9)$$

As said before, the transformed variable is close to Gaussian variable, we use the generated critical values under the Gaussian assumption to estimate the volatility. The constant C_p is calculated based on the estimates $\hat{\theta}(t)$ such that the innovation is standardized, i.e. $\text{Var}\{\hat{\varepsilon}(t)\} = \text{Var} \left[y(t) \{ \hat{C}_p / \hat{\theta}(t) \}^{\frac{1}{2p}} \right] = 1$. One then estimates the NIG distributional parameters of $\hat{\varepsilon}(t) = y(t) / \hat{\sigma}(t)$ where $\hat{\sigma}(t) = \{ \hat{\theta}(t) / \hat{C}_p \}^{\frac{1}{2p}}$. To get more accurate results, one

generates NIG innovations with the estimated distributional parameters and recalculates the critical values as in the Gaussian case.

The local exponential smoothing algorithm is described as follows:

1. Initialization: $\hat{\theta}^{(1)}(t) = \tilde{\theta}^{(1)}(t)$.

2. Loop: for $k \geq 2$,

$$\hat{\theta}^{(k)}(t) = \left(\frac{\gamma_k}{\hat{\theta}^{(k)}(t)} + \frac{1 - \gamma_k}{\hat{\theta}^{(k-1)}(t)} \right)^{-1}$$

where the aggregating parameter γ_k is computed as:

$$\gamma_k = \text{Kag}(L(\eta_k, \tilde{\theta}^{(k)}(t), \hat{\theta}^{(k-1)}(t)) / \zeta_{k-1}) \quad (10)$$

If $\gamma_k = 0$ then terminate by letting $\hat{\theta}^{(k)}(t) = \dots = \hat{\theta}^{(K)}(t) = \hat{\theta}^{(k-1)}(t)$.

3. Aggregation estimate: $\hat{\theta}(t) = \hat{\theta}^{(K)}(t)$.

4. Final estimate: $\hat{\sigma}(t) = \{\hat{\theta}(t)/C_p\}^{\frac{1}{2p}}$, where the constant C_p is computed such that the residuals $\hat{\varepsilon}(t) = y(t)/\hat{\sigma}(t)$ have a unit variance as assumed in the heteroscedastic model.

Consequently, the covariance matrices $D_y(t)$ and $\Sigma_x(t)$ are calculated.

Dynamic conditional correlation (DCC) model: Alternatively, the covariance of the return series can be estimated by the DCC model:

$$\Sigma_x(t) = D_x(t)R_x(t)D_x(t)^\top.$$

This technique first identifies the elements of the diagonal matrix $D_x(t)$ in the GARCH(1,1) setup and adaptively specifies the correlation matrix as:

$$R_x(t) = \tilde{R}_x(1 - \theta_1 - \theta_2) + \theta_1\{\varepsilon_x(t-1)\varepsilon_x(t-1)^\top\} + \theta_2R_x(t-1),$$

where \tilde{R}_x is the sample correlation of the risk factors, $\varepsilon_x \in \mathbb{R}^d$ are the standardized returns, i.e. risk factors divided by the univariate GARCH(1,1) volatilities, or equivalently by the squared diagonal elements in $D_x(t)$. The standardized returns are assumed to be Gaussian distributed. The parameters θ_1 and θ_2 are identified by the ML estimation.

2.3 Normal inverse Gaussian (NIG) distribution and fast Fourier transformation (FFT)

The estimated ICs are assumed to be NIG distributed. The NIG is a subclass of the GH distribution with a fixed value of $\lambda = -1/2$, see Eberlein and Prause (2002). With 4

distributional parameters, the NIG distribution is flexible to well match the behavior of real data. Compared to many other subclasses of GH distribution, the NIG distribution has a desirable property, saying that the scaled NIG variable belongs to the NIG distribution as well. The density of NIG random variable has a form of:

$$f_{\text{NIG}}(y; \alpha, \beta, \delta, \mu) = \frac{\alpha\delta}{\pi} \frac{K_1 \left\{ \alpha \sqrt{\delta^2 + (y - \mu)^2} \right\}}{\sqrt{\delta^2 + (y - \mu)^2}} \exp\{\delta \sqrt{\alpha^2 - \beta^2} + \beta(y - \mu)\},$$

where the distributional parameters fulfill $\mu \in \mathbb{R}$, $\delta > 0$ and $|\beta| \leq \alpha$. The modified Bessel function of the third kind $K_\lambda(\cdot)$ with an index $\lambda = 1$ has a form of:

$$K_\lambda(y) = \frac{1}{2} \int_0^\infty y^{\lambda-1} \exp\left\{-\frac{y}{2}(y + y^{-1})\right\} dy$$

The characteristic function of the NIG variable is:

$$\varphi_y(z) = \exp \left[\mathbf{i}z\mu + \delta \left\{ \sqrt{\alpha^2 - \beta^2} - \sqrt{\alpha^2 - (\beta + \mathbf{i}z)^2} \right\} \right]$$

Proof: The characteristic function of the GH random variable has a form of:

$$\varphi_y(z) = \exp(\mathbf{i}z\mu) \left\{ \frac{\alpha^2 - \beta^2}{\alpha^2 - (\beta + \mathbf{i}z)^2} \right\}^{\lambda/2} \frac{K_\lambda\{\delta \sqrt{\alpha^2 - (\beta + \mathbf{i}z)^2}\}}{K_\lambda(\delta \sqrt{\alpha^2 - \beta^2})}$$

Using the representation of the modified Bessel function with a fixed index $\lambda = -1/2$ derived in Barndorff-Nielsen and Blæsild (1981):

$$K_\lambda(y) = \sqrt{\frac{2}{\pi}} y^{-1/2} e^{-y},$$

it is straightforwardly to show that the assertion holds. \square

One desirable feature of the NIG distribution is its explicit scaling transformation. Multiplying the random variable by c , the resulting variable $y' = cy$ belongs to the NIG distribution as well:

$$f_{\text{NIG}}(y'; \alpha', \beta', \delta', \mu') = f_{\text{NIG}}(cy; \alpha/|c|, \beta/c, |c|\delta, c\mu). \quad (11)$$

Proof: It is easy to show the result by using the Jacobian transformation, see Härdle and Simar (2003). Given the density of y and let $\alpha' = \alpha/|c|$, $\beta' = \beta/c$, $\delta' = |c|\delta$ and $\mu' = c\mu$, the density of $y' = cy$ has a form of:

$$\begin{aligned} f(y') &= \frac{1}{|c|} f_y\left(\frac{y'}{c}\right) = \frac{\alpha'\delta'}{\pi} \frac{K_1 \left\{ \alpha' \sqrt{\delta'^2 + (y' - \mu')^2} \right\}}{\sqrt{\delta'^2 + (y' - \mu')^2}} \exp\{\delta' \sqrt{\alpha'^2 - \beta'^2} + \beta'(y' - \mu')\} \\ &= f_{\text{NIG}}(y'; \alpha', \beta', \delta', \mu'). \end{aligned}$$

□

To calculate risk measures, it requires the identification of the portfolio returns' density. Based on the GHICA model, the portfolio returns are calculated as:

$$r(t) = b(t)^\top W^{-1} D_y(t)^{1/2} \varepsilon_y(t)$$

where $b(t)$ is the trading strategy. Notice that the linear transformation of the NIG variable is not necessarily NIG distributed. In other words, the density of the return is unknown although the marginal densities are clear. On the meanwhile its characteristic function is explicitly writable. This is the same case as approximating the α -stable distribution in Menn and Rachev (2004), by which the Fourier transformation is used to approximate the density of the variable based on its characteristic function. This motivates us to use the technique to approximate the density of the return in the GHICA procedure.

Set $a = (a_1, \dots, a_d) = b(t)^\top W^{-1} D_y(t)^{1/2}$, the variable $\zeta_j = a_j \varepsilon_j$ is NIG distributed with $j = 1, \dots, d$, according to (11):

$$\zeta_j \sim \text{NIG}(\zeta_j, \check{\alpha}_j, \check{\beta}_j, \check{\delta}_j, \check{\mu}_j) = \text{NIG}(\zeta_j, \alpha_j/|a_j|, \beta_j/a_j, |a_j|\delta_j, a_j\mu_j).$$

The characteristic function of the return $r = \sum_{j=1}^d \zeta_j$ at time t is:

$$\varphi_r(z) = \prod_{j=1}^d \varphi_{\zeta_j}(z) = \exp \left[\mathbf{i}z \sum_{j=1}^d \check{\mu}_j + \sum_{j=1}^d \check{\delta}_j \left\{ \sqrt{\check{\alpha}_j^2 - \check{\beta}_j^2} - \sqrt{\check{\alpha}_j^2 - (\check{\beta}_j + \mathbf{i}z)^2} \right\} \right]$$

The density function is approximated by the Fourier transformation:

$$f(r) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} \exp(-\mathbf{i}tr) \psi(z) dt \approx \frac{1}{2\pi} \int_{-s}^s \exp(-\mathbf{i}tr) \psi(z) dt$$

The procedure of quantile estimation is summarized as follows:

- Implement the discrete fast Fourier transformation (DFT) to approximate the density of r at every time point t :
 1. Let $N = 2^m$ with $m \in \mathbb{N}$ and define an equidistance grid over the integral interval $[-s, s]$ by setting $h = \frac{2s}{N}$ and the grid points $z_j = -s + j * h$ with $j = 0, \dots, N$.
 2. Calculate the input of the DFT: $y_j = (-1)^j \psi(z_j^*)$ with $z_j^* = 0.5(z_j + z_{j+1})$ are the middle points. Notice that the characteristic function is time dependent.
 3. The density $f(r) = \frac{1}{2\pi} C_k \text{DFT}(y)_k$ with $C_k = \frac{2s}{N} (-1)^k \exp(-\frac{\mathbf{i}k\pi}{N}) \mathbf{i}$ with $k = 0, \dots, N - 1$. We refer to Borak, Detlefsen and Härdle (2005) and Menn and Rachev (2004) for more details. The corresponding values of $r = -\frac{N\pi}{2a} + \frac{\pi k}{a}$.
- The cumulative density function and the quantile are then approximated based on

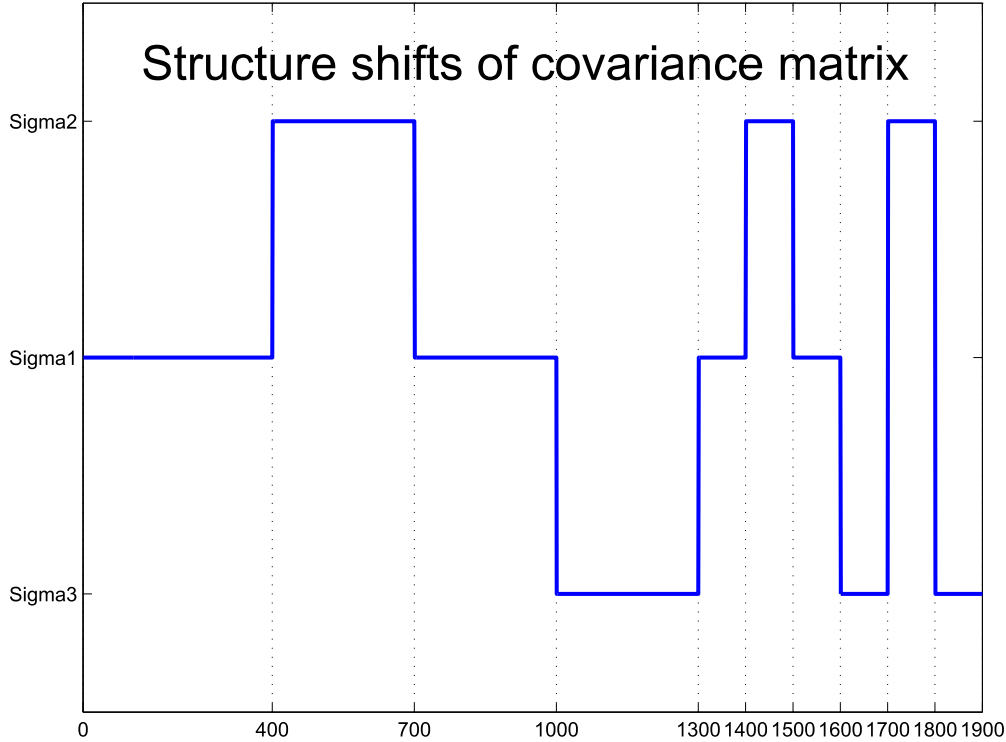


Fig. 2: Structure shifts of the generated covariance through time. Notice that there are shifts among matrices not up-and-down movements.

the resulting density.

3 Covariance estimation with simulated data

In this section, the GHICA versus the DCC, are implemented to estimate covariance of simulated data. The dimension is set to be $d = 50$. The simulation study is designed to include structure shifts of covariance. To be more specific, the designed covariance changes among three matrices over time, one is an identity matrix denoted as Σ_1 , meaning uncorrelatedness, and two symmetric and semi-positive defined matrices Σ_2 and Σ_3 . (Here we first generate $d * d$ matrix U_1 whose elements are uniform random variables for Σ_2 and standard Gaussian variables for Σ_3 , then calculate a new matrix $U_2 = U_1 * U_1'$ to guarantee the semi-positiveness. The elements $\Sigma(i, j)$ of the target matrix are calculated as $\Sigma(i, j) = U_2(i, j) / \sqrt{U_2(i, i)U_2(j, j)}$.) The eigenvalues of these two matrices are distributed in $[5.92e-004, 3.779]$ (Σ_2) and $[0.002, 3.573]$ (Σ_3) respectively. The off-diagonal values span over $[-0.433, 0.468]$ in the first self-correlated matrix (Σ_2) and $[-0.447, 0.464]$ in the second one (Σ_3). Temporal stationarity is assumed to be long for 400 time units and short for 100 units. The structure shifts of the generated covariance are illustrated in Figure 2. The level of the shifts is either small with a shift from one self-correlated matrix (Σ_2 or Σ_3) to the identity matrix or contrariwise, e.g. at the point 700, or large with a shift between the two

self-correlated matrices, e.g. at the point 1800.

Furthermore, two distributional parameters μ and β of the standardized NIG innovations $\varepsilon_x(t)$ are set to be 0, meaning that the innovations are centered around 0 and symmetric distributed, see Barndorff-Nielsen and Blæsild (1981). By doing so, the mean and variance of the NIG innovations only depend on α and δ :

$$\begin{aligned} \mathbf{E}(\varepsilon_x) &= \mu + \frac{\beta\delta}{\sqrt{\alpha^2 - \beta^2}} = 0 \\ \text{Var}(\varepsilon_x) &= \frac{\delta}{\sqrt{\alpha^2 - \beta^2}} + \frac{\beta^2}{\delta^3\sqrt{\alpha^2 - \beta^2}} = \frac{\delta}{\alpha} = 1 \end{aligned}$$

This result is used to generate the standardized innovations, by which $\alpha \sim U[1, 2]$ is suggested by our experience on real data analysis and $\delta = \alpha$.

In the Monte Carlo simulation, we generate $d = 50$ NIG variables with the designed covariance and distributional parameters:

$$x(t) = \Sigma_x^{1/2}(t)\varepsilon_x(t).$$

The sample size is $T = 1900$ and the scenarios are repeated $N = 100$ times. The covariance matrix is estimated using the GHICA procedure and the DCC method respectively.

The GHICA method first converts the underlying series to ICs by a linear transformation:

$$x(t) = W^{-1}y(t) = W^{-1}D_y^{1/2}(t)\varepsilon_y(t),$$

by which the elements of $D_y(t)$ on the diagonal are estimated using the local exponential smoothing method. In the local exponential smoothing estimation, we set the involved parameters $c = 0.01$, $a = 1.25$ and $p = 0.25$. The sequence of the smoothing parameters $\{\eta_k\}$ are $0.600, \dots, 0.982$ with $K = 15$, based on the condition $(1 - \eta_k)/(1 - \eta_{k+1}) = a$ in (7). The first 300 observations are reserved as training set for the very beginning estimations, since the largest smoothing parameter used in this study corresponds to a window with 259 observations.

The covariance of $x(t)$ is calculated by the basic statistical property:

$$\Sigma_x(t) = W^{-1}D_y(t)W^{-1\top}$$

The DCC method assumes that the underlying series are Gaussian distributed. It decomposes the covariance matrix to a product of diagonal variance matrix and correlation matrix:

$$\Sigma_x(t) = D_x(t)R_x(t)D_x(t)^\top.$$

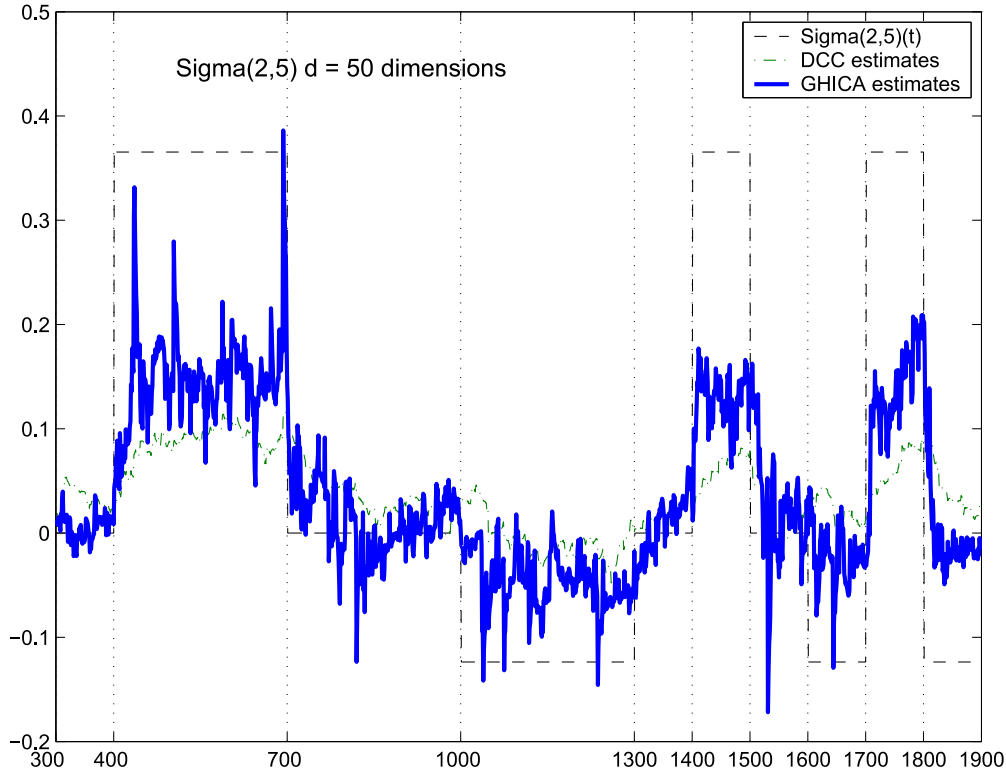


Fig. 3: Realized estimates of $\Sigma(2, 5)$ based on the GHICA and DCC methods. The generated data consists of 50 NIG distributed components.

where $D_x(t)$ consists of the variances of $x(t)$ on the diagonal that are estimated in the GARCH(1,1) setup.

Figure 3 displays one realization of $\Sigma(2, 5)$, i.e. the covariance of the second and fifth risk factors $x_2(t)$ and $x_5(t)$, based on one simulation data. The true values are 0.365 in Σ_2 and -0.124 in Σ_3 . As expected, the GHICA estimates are sensitive to structure shifts through time. The DCC estimates, on the contrary, are over-smooth and slowly follow the shifts. Given more often shifts around the last hundreds of time points, the DCC estimates deliver less information on the movements. Recall that 100 points correspond to 4 months observations of daily returns. It is rational to surmise that structure shifts happen so often in the active financial markets, see Merton (1973). The similar estimation results are observed in the other elements of the covariance, which are eliminated here.

To measure the accuracy of estimation, ratio of absolute estimation error (RAE) of the estimates w.r.t. the true covariance are calculated pointwise.

$$\text{RAE}(i, j) = \frac{\sum_{t=301}^T |\hat{\Sigma}_{(i,j)}^{\text{GHICA}}(t) - \Sigma_{(i,j)}(t)|}{\sum_{t=301}^T |\hat{\Sigma}_{(i,j)}^{\text{DCC}}(t) - \Sigma_{(i,j)}(t)|}$$

If $\text{RAE}(i, j) \leq 1$, it means that the GHICA method reaches higher accuracy in the estima-

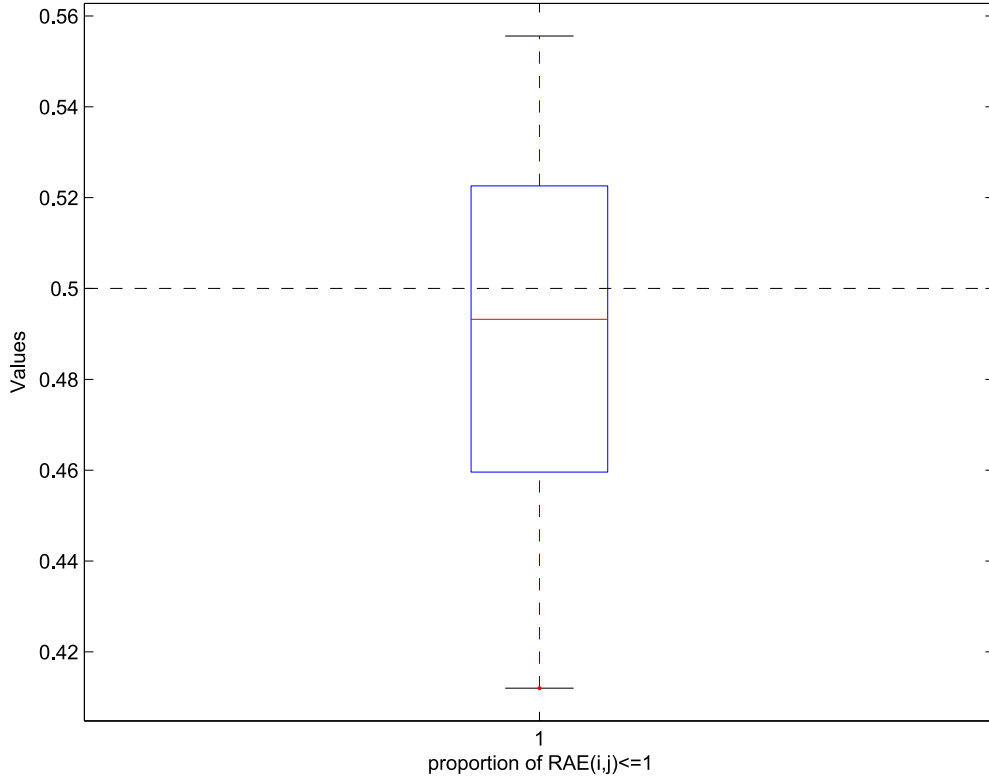


Fig. 4: Boxplot of the proportion $\frac{\sum_i \sum_j \mathbf{1}(\text{RAE}(i,j) \leq 1)}{d \times d}$ for $i, j = 1, \dots, d$. Here $d = 50$ and the proportions on the base of 100 simulations are considered.

tion of $\Sigma(i, j)$ than the DCC. To compare the general performance of these two methods in covariance estimation, we check the proportion of the RAEs among the 2500 ($d \times d$) elements that are smaller or equal to one, i.e. $\frac{\sum_i \sum_j \mathbf{1}(\text{RAE}(i,j) \leq 1)}{d \times d}$ for $i, j = 1, \dots, d$. Notice that the proportion with value of 0.5 indicates that half elements are better estimated by using the GHICA and the other half are better done by the DCC. In other words, the considered methods have a comparable accuracy of estimation. Figure 4 displays the boxplot of the 100 proportions. The mean of the proportion is 0.4904 among the 100 simulations. It states that the DCC method performs a little bit better than the GHICA in the sense of accuracy. On the meanwhile, the GHICA method is much fast and sensitive to structure shifts.

4 Risk management with real data

In this section, we implement the proposed GHICA method to calculate risk measures using real data sets: 20-dimensional German DAX portfolio and 7-dimensional exchange rate portfolio. The results are compared with those based on alternative risk management models. The data sets have been kindly provided by the financial and economic data center (FEDC) of the Collaborative Research Center 649 on Economic Risk of the Humboldt-

Universität zu Berlin (<http://sfb649.wiwi.hu-berlin.de>). Before giving detailed description of the data sets, we analyze the risk measures from the viewpoints of regulatory, investors and internal supervisory.

Regulatory requirement: Financial institutions generally face market risk that arises from the uncertainty due to changes in market prices and rates such as share prices, foreign exchange rates and interest rates, the correlations among them and their levels of volatility, see Jorion (2001). The market risk is the main risk source and has a great negative influence on the development of economic. The famous example is the stock crashes in the autumn 1929 and 1987 which caused a violent depression in the United States and some other countries, with the collapse of financial markets and the contraction of production and employment. To alleviate the down influence of market risks, regulation on banking and other financial institutions has been strengthened since the mid-1990s. The goals of the regulation are to restrict the happening of extremely large losses and require banks to reserve adequate capital. In 1998 the Basel accord officially allowed financial institutions to use their internal models to measure market risks. Among others, Value at Risk (VaR) has been considered as industry standard risk measure:

$$\text{VaR}_{t,\text{pr}} = -\text{quantile}_{\text{pr}}\{r(t)\}.$$

where pr is the $h = 1$ -day or $h = 5$ -day forecasted probability of the portfolio returns. Internal models for risk management are verified in accordance with the “traffic light” rule that counts the number of exceptions over VaR at 1% probability spanning the last 250 days and identifies the multiplicative factor M_f in the market risk charge calculation, see Franke, Härdle and Hafner (2004):

$$\text{Risk charge}_t = \max \left(M_f \frac{1}{60} \sum_{i=1}^{60} \text{VaR}_{t-i,1\%}, \text{VaR}_{t,1\%} \right)$$

The multiplicative factor M_f has a floor value 3. It increases corresponding to the number of exceptions, see Table 1. For example, if an internal model generates 7 exceptions at 1% probability over the last 250 days, the model is in the yellow zone and its multiplicative factor is $M_f = 3.65$. Financial institutions whose internal model is located in the yellow or red zone, with a very high probability, are required to reserve more risk capital than their internal-model-based VaRs. Notice that the increase of risk charge will reduce the ratio of profit since the reserved capital can not be invested. On the meanwhile, an internal model is automatically accepted if the number of exceptions does not exceed 4. This regulatory rule in fact suggests banks to control VaR at 1.6% (i.e. 4/250) instead of 1% probability. It is clear that 1.6%-VaR is smaller than 1%-VaR. Therefore an internal model is particularly desirable by financial institutions if its empirical probability is smaller or equal to 1.6%, and simultaneously requires risk charge as small as possible. Here a simplified calculation

No. exceptions	Increase of M_f	Zone
0 bis 4	0	green
5	0.4	yellow
6	0.5	yellow
7	0.65	yellow
8	0.75	yellow
9	0.85	yellow
More than 9	1	red

Tab. 1: Traffic light as a factor of the exceeding amount, cited from Franke, Härdle and Hafner (2004).

on the average value of VaRs is used as risk charge for comparison:

$$\text{Risk charge (RC)} = \text{mean}(\text{VaR}_{t,\text{pr}})$$

Investor: It is known that VaR is inappropriate for the measurement of capital adequacy, since it controls only the probability of default, i.e. the frequency of losses, but not the size of losses in the case of default. For this reason, investors concern expected shortfall (ES) more than VaR to measure and control their risks.

$$\text{ES} = \text{E}\{-r(t) \mid -r(t) > \text{VaR}_{t,\text{pr}}\}$$

Investors suffer loss once bankruptcy happens. Even in the “best” situation, their loss equals to the difference between the total loss and the reserved risk capital, i.e. the value of ES. Generally risk-averse investors care the amount of loss and thus prefer an internal model with small value of ES. Risk-seeking investors, on the other hand, care profit and hence the small value of risk charge favors their requirement.

Internal supervisory: It is important for internal supervisory to exactly measure the market risk exposures before risk controlling. For this reason, internal supervisory prefers the model delivering accurate probability prediction, i.e. the empirical probability $\hat{\text{pr}}$ is as close to the expected values as possible:

$$\hat{\text{pr}} = \frac{\text{No. exceptions}}{\text{No. total observations}}$$

Given two models with the same empirical probability, the model has a smaller value of ES is considered better than the other. Here two extreme probabilities are considered, i.e. $\text{pr} = 1\%$ for regulatory reason and $\text{pr} = 0.5\%$ used by financial institutions with AAA rating.

4.1 Data analysis 1: DAX portfolio

The primary target of the real data analysis is to compare the forecasting ability of the GHICA method with two alternatives, the RiskMetrics method under the Gaussian distributional assumption and a modification with the Student- $t(6)$ distributional assumption (abbreviated as $t(6)$ method) in the market. The comparison is demonstrated based on 20 DAX stocks over a long time period, starting on 1974/01/02 and ending on 1996/12/30 (5748 observations). The return series are all centered around 0 and have heavy tails (kurtosis > 3), the smallest correlation coefficient is 0.3654. Hypothetical German DAX portfolios are constructed with two static trading strategies $b(t) = b^{(1)} = (1/d, \dots, 1/d)^\top$ and $b(t) = b^{(2)} \sim U[0, 1]^d$. Such a simple portfolio construction eliminates the influence of strategy adjustments on the calculation. The portfolio returns are analyzed using the RiskMetrics or the $t(6)$ method. Here the unknown volatility process of the portfolio is estimated using the exponential smoothing method with $\eta = 0.94$:

$$\begin{aligned} r(t) &= b^\top x(t) = \sigma_r(t) \varepsilon_r(t) \\ \sigma_r^2(t) &= \left\{ \sum_{m=0}^M \eta^m r^2(t-m-1) \right\} / \left(\sum_{m=0}^M \eta^m \right) \end{aligned}$$

where the truncated value M fulfills the condition $\eta^{(M+1)} \leq 0.01$. Notice that given a dynamic trading strategy, this simplification needs to repeatedly estimate the density of the time varying hypothetical portfolio returns, and it often suffers from a low accuracy of estimation.

Figure 5 depicts the one day log-returns of the DAX portfolio with the static trading strategy $b(t) = b^{(1)}$. The VaRs from 1975/03/17 to 1996/12/30 at $\text{pr} = 0.5\%$ are displayed w.r.t. three methods, the GHICA, the RiskMetrics and the $t(6)$. The most volatile time period over $t \in [3300, 4300]$ is detailed in the bottom diagram. Recall that on the Monday, 19 October 1987, the worldwide downward jump of stocks happened. Dow Jones Industrial Average for example dropped by over 500 points. At this market quiver around $t = 3446$, the GHICA method exactly achieves the locations of extreme losses whereas the RiskMetrics and $t(6)$ methods over-react to them. Such over reactions induce large risk charges unnecessarily. On the other hand, it is observed that these two alternative methods give close forecasts to some extreme losses, e.g. around time points 4000 and 4500. As a result, the associating values of ES are small and satisfy the requirement of risk-averse investors.

Table 2 reports the risk measures based on the three methods. In general, the RiskMetrics is successful in fulfilling the minimal requirement of regulatory. The $t(6)$ method is preferred by investors who consider risk happened with 1% probability. The GHICA method performs better than the other two for internal supervisory and requirement of

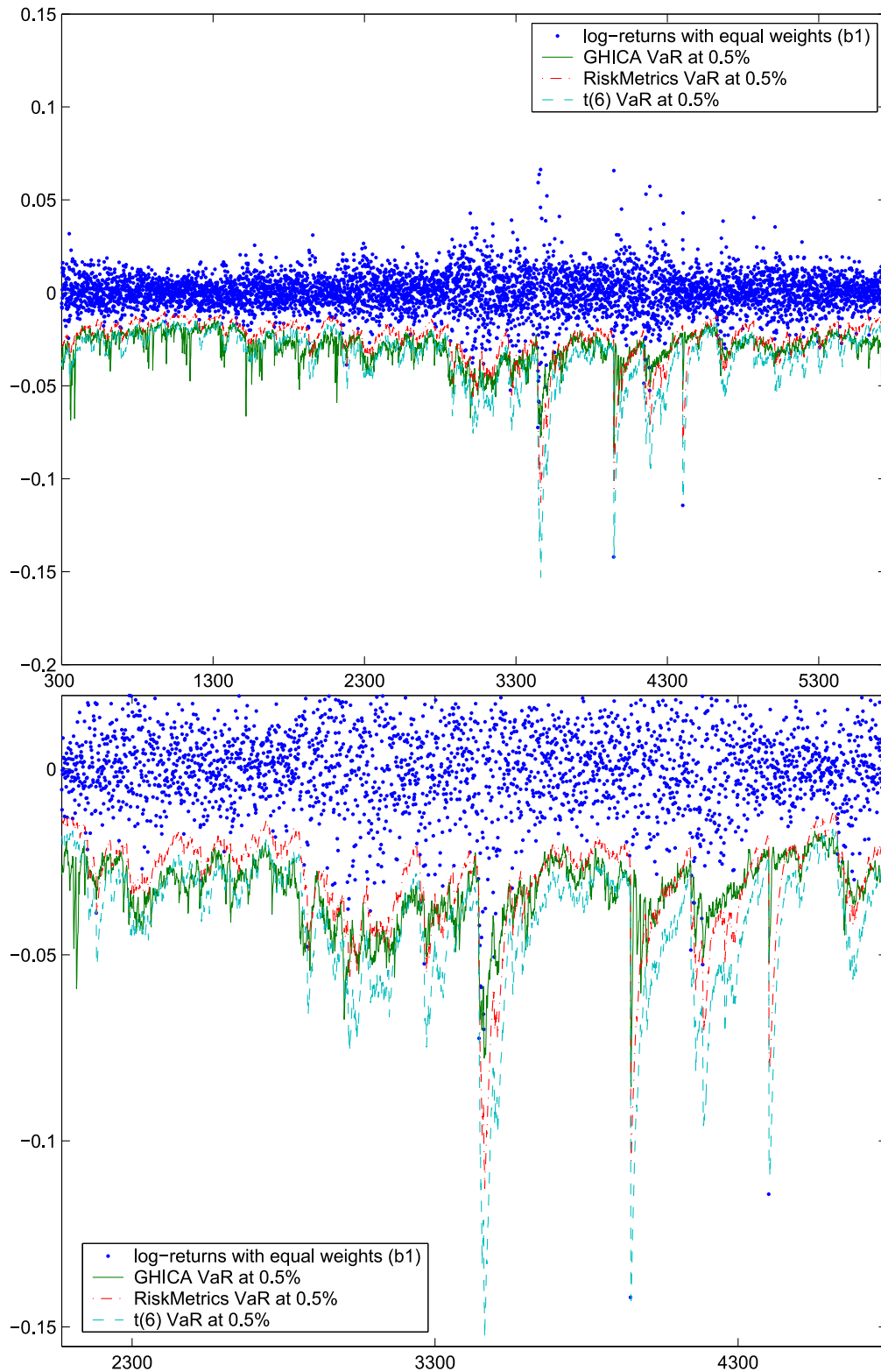


Fig. 5: One day log-returns of the DAX portfolio with the static trading strategy $b(t) = b^{(1)}$. The VaRs are from 1975/03/17 to 1996/12/30 at $pr = 0.5\%$ w.r.t. three methods, the GHICA, the RiskMetrics and the $t(6)$. Part of the VaR time plot is enlarged and displayed on the bottom.

h	$b(t)$	pr	GHICA			RiskMetrics $N(\mu, \sigma^2)$			Exponential smoothing $t(6)$		
			$\hat{p}r$	RC	ES	$\hat{p}r$	RC	ES	$\hat{p}r$	RC	ES
1	$b^{(1)}$	1%	0.55%	0.0264	0.0456	1.18% ^s	0.0229 ^r	0.0279	0.40%	0.0292	0.0269 ⁱ
	$b^{(1)}$	0.5%	0.44% ^s	0.0297	0.0472 ⁱ	0.75%	0.0254	0.0317	0.23%	0.0345	0.0506
	$b^{(2)}$	1%	0.59%	0.0265	0.0448	1.03% ^s	0.0231 ^r	0.0288	0.38%	0.0294	0.0406 ⁱ
	$b^{(2)}$	0.5%	0.42% ^s	0.0298	0.0476 ⁱ	0.71%	0.0256	0.0315	0.21%	0.0347	0.0514
5	$b^{(1)}$	1%	0.83%	0.0550	0.0841	1.15% ^s	0.0481 ^r	0.0602	0.19%	0.0665	0.0833 ⁱ
	$b^{(1)}$	0.5%	0.51% ^s	0.0612	0.0939 ⁱ	0.64%	0.0536	0.0683	0.09%	0.0784	0.1067
	$b^{(2)}$	1%	0.83% ^s	0.0554	0.0828 ⁱ	1.18%	0.0488 ^r	0.0613	0.16%	0.0673	0.0852
	$b^{(2)}$	0.5%	0.50% ^s	0.0617	0.0943 ⁱ	0.63%	0.0543	0.0676	0.07%	0.0794	0.1218

Tab. 2: Risk analysis of the DAX portfolios with two static trading strategies. The concerned forecasting interval is $h = 1$ or $h = 5$ days. The best results to fulfill the regulatory requirement are marked by ^r. The method preferred by investor is marked by ⁱ. For the internal supervisory, the method marked by ^s is recommended.

risk-averse investors who care the extreme risk happened with 0.5% probability.

4.2 Data analysis 2: Foreign exchange rate portfolio

In financial markets, traders adjust trading strategy according to information obtained. The GHICA is easily applicable to dynamic portfolios. We consider here 7 actively traded exchange rates, Euro (EUR), the US dollar (USD), the British pounds (GBP), the Japanese yen (JPY) and the Singapore dollar (SGD) from 1997/01/02 to 2006/01/05 (2332 observations). The foreign exchange rate (FX) market is the most active and liquid financial market in the world. It is realistic to analyze a dynamic portfolio with daily time varying trading strategy $b^{(3)}(t)$. The strategy at time point t relies on the realized returns at $t - 1$, the proportions of which w.r.t the sum of returns:

$$b^{(3)}(t) = \frac{x(t-1)}{\sum_{j=1}^d x_j(t-1)}$$

where $x(t) = \{x_1(t), \dots, x_d(t)\}^\top$. Among these data sets, the returns of the EUR/SGD and USD/JPY rates are least correlated with the correlation coefficient 0.0071 whereas the returns of the EUR/USD and EUR/SGD rates are most correlated with the coefficient 0.6745. The resulting portfolio returns span over $[-0.7962, 0.7074]$.

The GHICA method is compared with an alternative method, abbreviated as DCCN, that applies the DCC covariance estimation under the Gaussian distributional assumption.

$$r(t) = b(t)^\top x(t) = b(t)^\top \Sigma_x^{(1/2)}(t) \varepsilon_x(t)$$

where $\varepsilon_x \sim N(\mu, \Sigma_\varepsilon)$ with the diagonal covariance matrix Σ_ε . Notice that the quantile

			GHICA			DCCN		
h	$b(t)$	pr	$\hat{p}r$	RC	ES	$\hat{p}r$	RC	ES
1	$b^{(3)}(t)$	1%	1.28% ^s	0.0453^r	0.0778	1.59%	0.0494	0.0254ⁱ
	$b^{(3)}(t)$	0.5%	0.59% ^s	0.0493	0.1944 ⁱ	0.94%	0.0547	0.0289
5	$b^{(3)}(t)$	1%	1.53% ^s	0.0806^r	0.2630 ⁱ	4.17%	0.0993	0.1735
	$b^{(3)}(t)$	0.5%	0.79% ^s	0.1092	0.2801 ⁱ	3.44%	0.1100	0.1389

Tab. 3: Risk analysis of the dynamic exchange rate portfolio. The best results to fulfill the regulatory requirement are marked by ^r. The recommended method to the investor is marked by ⁱ. For the internal supervisory, we recommend the method marked by ^s.

vector with pr-quantiles of individual innovations does not necessarily correspond to the pr-quantile of the portfolio return. Under the Gaussian distributional assumption, the standardized DCCN returns are theoretically cross independent and the Gaussian quantiles of the portfolio can be easily calculated. The dynamic mean, variance of the portfolio's returns have values of:

$$\begin{aligned} \mathbf{E}\{r(t)\} &= b(t)^\top \Sigma_x^{(1/2)}(t) \mathbf{E}\{\varepsilon_x(t)\} \\ \mathit{Var}\{r(t)\} &= b(t)^\top \Sigma_x^{(1/2)}(t) \mathit{Var}\{\varepsilon_x(t)\} \Sigma_x^{(1/2)\top}(t) b(t) \end{aligned}$$

The GHICA method in general presents better results than the DCCN. Except the value of ES at 1% level, the GHICA fulfills the requirements of regulatory, internal supervisory and investors, see Table 3. For $h = 1$ day forecasts, the DCCN gives although a closer VaR value to 1.6%, i.e. the ideal probability for regulatory, its risk charge with a value of 0.0494 is larger than that based on the GHICA, 0.0453. Therefore the GHICA is more favored in fulfilling the minimal regulatory requirement.

The two real data studies show that the GHICA method fulfills the minimal regulatory requirement by controlling the risk inside 1.6% level and requiring small risk charge, in particular satisfies the internal supervisory requirement by precisely measuring risk level as expected and favors the investors' requirement by delivering small size of loss. In summary, the GHICA method is not only a realistic and fast procedure given either static or dynamic portfolios but also produces better results than several alternative risk management methods.

References

- Anderson, T., Bollerslev, T., Diebold, F. and Labys, P. (2001). The distribution of realized exchange rate volatility, *Journal of the American Statistical Association* pp. 42–55.
- Barndorff-Nielsen, O. and Blæsild, P. (1981). Hyperbolic distribution and ramifications: Contributions to theory and applications, in C. Taillie, P. Patil and A. Baldessari (eds), *Statistical Distributions in Scientific Work*, Vol. 4, D. Reidel, pp. 19–44.
- Belomestny, D. and Spokoiny, V. (2006). Spatial aggregation of local likelihood estimates with applications to classification, *WIAS Preprint*.
- Borak, S., Detlefsen, K. and Härdle, W. (2005). FFT-based option pricing, in P. Cizek, W. Härdle and R. Weron (eds), *Statistical Tools for Finance and Insurance*, Springer Verlag.
- Chen, Y. and Spokoiny, V. (2006). Local exponential smoothing with applications to volatility estimation and risk management, *working paper*.
- Chen, Y., Härdle, W. and Jeong, S. (2005). Nonparametric risk management with generalized hyperbolic distributions, *SFB 649, Discussion paper 2005-001*, <http://sfb649.wiwi.hu-berlin.de>.
- Chen, Y., Härdle, W. and Spokoiny, V. (2006). Portfolio value at risk based on independent components analysis, *Journal of Computational and Applied Mathematics, forthcoming*.
- Cover, T. and Thomas, J. (1991). *Elements of information theory*, Wiley.
- Eberlein, E. and Prause, K. (2002). The generalized hyperbolic model: financial derivatives and risk measures, in H. Geman, D. Madan, S. Pliska and T. Vorst (eds), *Mathematical Finance-Bachelier Congress 2000*, Springer Verlag.
- Engle, R. (2002). Dynamic conditional correlation - a simple class of multivariate garch models, *Journal of Business and Economic Statistics*, 20(3) pp. 339–350.
- Engle, R. and Kroner, F. (1995). Multivariate simultaneous generalized arch, *Econometric Theory* 11 pp. 122–150.
- Engle, R. and Sheppard, K. (2001). Theoretical and empirical properties of dynamic conditional correlation multivariate garch, *NBER Working Paper 8554*.
- Flury, B. (1998). *Common Principal Components and Related Multivariate Models*, John Wiley & Sons, Inc.
- Franke, J., Härdle, W. and Hafner, C. (2004). *Statistics of Financial Markets*, Springer-Verlag Berlin Heidelberg New York.

- Härdle, W. and Simar, L. (2003). *Applied Multivariate Statistical Analysis*, Springer-Verlag Berlin Heidelberg New York.
- Härdle, W., Herwartz, H. and Spokoiny, V. (2003). Time inhomogeneous multiple volatility modelling, *Journal of Financial Econometrics* **1**: 55–95.
- Hyvärinen, A. (1998). *New Approximations of Differential Entropy for Independent Component Analysis and Projection Pursuit*, MIT Press, pp. 273–279.
- Hyvärinen, A., Karhunen, J. and Oja, E. (2001). *Independent Component Analysis*, John Wiley & Sons, Inc.
- Jorion, P. (2001). *Value at Risk*, McGraw-Hill.
- Menn, C. and Rachev, S. (2004). Calibrated FFT-based density approximations for α -stable distributions.
- Merton, R. (1973). Theory of rational option pricing, *The Bell Journal of Economics and Management Science* **4**: 141–183.
- Polzehl, J. and Spokoiny, V. (2006). Propagation-separation approach for local likelihood estimation, *Probability Theory and Related Fields* pp. 335–362.

Empirical Pricing Kernels and Investor Preferences

K. Detlefsen¹, W. K. Härdle², R. A. Moro³,

¹CASE – Center for Applied Statistics and Economics, Humboldt-Universität zu Berlin, Spandauer Straße 1, 10178 Berlin, Germany; e-mail: detlefsen@wiwi.huberlin.de; phone: +49(0)30 2093-5807

²CASE – Center for Applied Statistics and Economics, Humboldt-Universität zu Berlin, Spandauer Straße 1, 10178 Berlin, Germany; e-mail: haerdle@wiwi.huberlin.de; phone: +49(0)30 2093-5630

³German Institute for Economic Research, Königin-Luise-Straße 5, 14195 Berlin, Germany; e-mail: rmoro@diw.de; phone: +49(0)30 8978-9262 and CASE – Center for Applied Statistics and Economics, Humboldt-Universität zu Berlin, Spandauer Straße 1, 10178 Berlin

Abstract

This paper analyzes empirical market utility functions and pricing kernels derived from the DAX and DAX option data for three market regimes. A consistent parametric framework of stochastic volatility is used. All empirical market utility functions show a region of risk proclivity that is reproduced by adopting the hypothesis of heterogeneous individual investors whose utility functions have a switching point between bullish and bearish attitudes. The inverse problem of finding the distribution of individual switching points is formulated in the space of stock returns by discretization as a quadratic optimization problem. The resulting distributions vary over time and correspond to different market regimes.

JEL classification: G12, G13, C50

Keywords: Utility function, pricing kernel, behavioral finance, risk aversion, risk proclivity, Heston model

1 Introduction

Numerous attempts have been undertaken to describe basic principles on which the behaviour of individuals are based. Expected utility theory was originally proposed by J. Bernoulli in 1738. In his work J. Bernoulli used such terms as risk aversion and risk premium and proposed a concave (logarithmic) utility function, see Bernoulli (1956). The utilitarianism theory that emerged in the 18th century considered utility maximization as a principle for the organisation of society. Later the expected utility idea was applied to game theory and formalized by von Neumann and Morgenstern (1944). A utility function relates some observable variable, in most cases consumption, and an unobservable utility level that this consumption delivers. It was suggested that individuals' preferences are based on this unobservable utility: such bundles of goods are preferred that are associated with higher utility levels. It was claimed that three types of utility functions – concave, convex and linear – correspond to three types of individuals – risk averse, risk neutral and risk seeking. A typical economic agent was considered to be risk averse and this was quantified by coefficients of relative or absolute risk aversion. Another important step in the development of utility theory was the prospect theory of Kahneman and Tversky (1979). By behavioural experiments they found that people act risk averse above a certain reference point and risk seeking below it. This implies a concave form of the utility function above the reference point and a convex form below it.

Besides these individual utility functions, market utility functions have recently been analyzed in empirical studies by Jackwerth (2000), Rosenberg and Engle (2002) and others. Across different markets, the authors observed a common pattern in market utility functions: There is a reference point near the initial wealth and in a region around this reference point the market utility functions are convex. But for big losses or gains they show a concave form – risk aversion. Such utility functions disagree with the classical utility functions of von Neumann and Morgenstern (1944) and also with the findings of Kahneman and Tversky (1979). They are however in concordance with the utility function form proposed by Friedman and Savage (1948).

In this paper, we analyze how these market utility functions can be explained by aggregating individual investors' attitudes. To this end, we first determine empirical pricing kernels from DAX data. Our estimation procedure is based on historical and risk neutral densities and these distributions are derived with stochastic volatility models that are widely used in industry. From these pricing kernels we construct the corresponding market utility functions. Then we describe our method of aggregating individual utility functions to a market utility function. This leads to an inverse problem for

the density function that describes how many investors have the utility function of each type. We solve this problem by discrete approximation. In this way, we derive utility functions and their distribution among investors that allow to recover the market utility function. Hence, we explain how (and what) individual utility functions can be used to form the behaviour of the whole market.

The paper is organized as follows: In section 2, we describe the theoretical connection between utility functions and pricing kernels. In section 3, we present a consistent stochastic volatility framework for the estimation of both the historical and the risk neutral density. Moreover, we discuss the empirical pricing kernel implied by the DAX in 2000, 2002 and 2004. In section 4, we explain the utility aggregation method that relates the market utility function and the utility functions of individual investors. This aggregation mechanism leads to an inverse problem that is analyzed and solved in this section. In section 5, we conclude and discuss related approaches.

2 Pricing kernels and utility functions

In this section, we derive the fundamental relationship between utility functions and pricing kernels. It describes how a representative utility function can be derived from historical and risk-neutral distributions of assets. In the following sections, we estimate the empirical pricing kernel and observe in this way the market utility function.

First, we derive the price of a security in an equilibrium model: we consider an investor with a utility function U who has as initial endowment one share of stock. He can invest into the stock and a bond up to a final time when he can consume. His problem is to choose a strategy that maximizes the expected utility of his initial and terminal wealth. In continuous time, this leads to a well known optimization problem introduced by Merton (1973) for stock prices modelled by diffusions. In discrete time, it is a basic optimization problem, see Cochrane (2001).

From this result, we can derive the asset pricing equation

$$P_0 = E^P [\psi(S_T)M_T]$$

for a security on the stock (S_t) with payoff function ψ at maturity T . Here, P_0 denotes the price of the security at time 0 and E^P is the expectation with respect to the real/historical measure P . The stochastic discount factor M_T is given by

$$M_T = \beta U'(S_T)/U'(S_0) \tag{1}$$

where β is a fixed discount factor. This stochastic discount factor is actually the projection of the general stochastic discount factor on the traded asset (S_t). The stochastic discount factor can depend on more variables in general. But as discussed in Cochrane (2001) this projection has the same interpretation for pricing as the general stochastic discount factor.

Besides this equilibrium based approach, Black and Scholes (1973) derived the price of a security relative to the underlying by constructing a perfect hedge. The resulting continuous delta hedging strategy is equivalent to pricing under a risk neutral measure Q under which the discounted price process of the underlying becomes a martingale. Hence, the price of a security is given by an expected value with respect to a risk neutral measure Q :

$$P_0 = E^Q [\exp(-rT)\psi(S_T)]$$

If p denotes the historical density of S_T (i.e. $P(S_T \leq s) = \int_{-\infty}^s p(x) dx$) and q the risk neutral density of S_T (i.e. $Q(S_T \leq s) = \int_{-\infty}^s q(x) dx$) then we get

$$\begin{aligned} P_0 &= \exp(-rT) \int \psi(x)q(x)dx \\ &= \exp(-rT) \int \psi(x) \frac{q(x)}{p(x)} p(x)dx \\ &= E^P \left[\exp(-rT)\psi(S_T) \frac{q(S_T)}{p(S_T)} \right] \end{aligned} \tag{2}$$

Combining equations (1) and (2) we see

$$\beta \frac{U'(s)}{U'(S_0)} = \exp(-rT) \frac{q(s)}{p(s)}.$$

Defining the pricing kernel by $K = q/p$ we conclude that the form of the market utility function can be derived from the empirical pricing kernel by integration:

$$\begin{aligned} U(s) &= U(S_0) + \int_{S_0}^s U'(S_0) \frac{\exp(-rT)}{\beta} \frac{q(x)}{p(x)} dx \\ &= U(S_0) + \int_{S_0}^s U'(S_0) \frac{\exp(-rT)}{\beta} K(x) dx \end{aligned}$$

because S_0 is known.

As an example, we consider the model of Black and Scholes (1973) where the stock follows a geometric Brownian motion

$$dS_t/S_t = \mu dt + \sigma dW_t \quad (3)$$

Here the historical density p of S_t is log-normal, i.e.

$$p(x) = \frac{1}{x} \frac{1}{\sqrt{2\pi\tilde{\sigma}^2}} \exp\left\{-\frac{1}{2}\left(\frac{\log x - \tilde{\mu}}{\tilde{\sigma}}\right)^2\right\}, \quad x > 0$$

where $\tilde{\mu} = (\mu - \sigma^2/2)t + \log S_0$ and $\tilde{\sigma} = \sigma\sqrt{t}$. Under the risk neutral measure Q the drift μ is replaced by the riskless interest rate r , see e.g. Harrison and Pliska (1981). Thus, also the risk neutral density q is log-normal. In this way, we can derive the pricing kernel

$$K(x) = \left(\frac{x}{S_0}\right)^{-\frac{\mu-r}{\sigma^2}} \exp\{(\mu - r)(\mu + r - \sigma^2)T/(2\sigma^2)\}.$$

This pricing kernel has the form of a derivative of a power utility

$$K(x) = \lambda \left(\frac{x}{S_0}\right)^{-\gamma}$$

where the constants are given by $\lambda = e^{\frac{(\mu-r)(\mu+r-\sigma^2)T}{2\sigma^2}}$ and $\gamma = \frac{\mu-r}{\sigma^2}$. This gives a utility function corresponding to the underlying (3)

$$U(S_T) = \left(1 - \frac{\mu - r}{\sigma^2}\right)^{-1} S_T^{(1 - \frac{\mu-r}{\sigma^2})}$$

where we ignored additive and multiplicative constants. In this power utility function the risk aversion is not given by the market price of risk $(\mu - r)/\sigma$. Instead investors take the volatility more into account. The expected return $\mu - r$ that is adjusted by the riskfree return is related to the variance. This results in a higher relative risk aversion than the market price of risk.

A utility function corresponding to the Black-Scholes model is shown in the upper panel of figure 1 as a function of returns. In order to make different market situations comparable we consider utility functions as functions of (half year) returns $R = S_{0.5}/S_0$. We chose the time horizon of half a year ahead for our analysis. Shorter time horizons are interesting economically and moreover the historical density converges to the Dirac measure so that results become trivial (in the end). Longer time horizons are economically

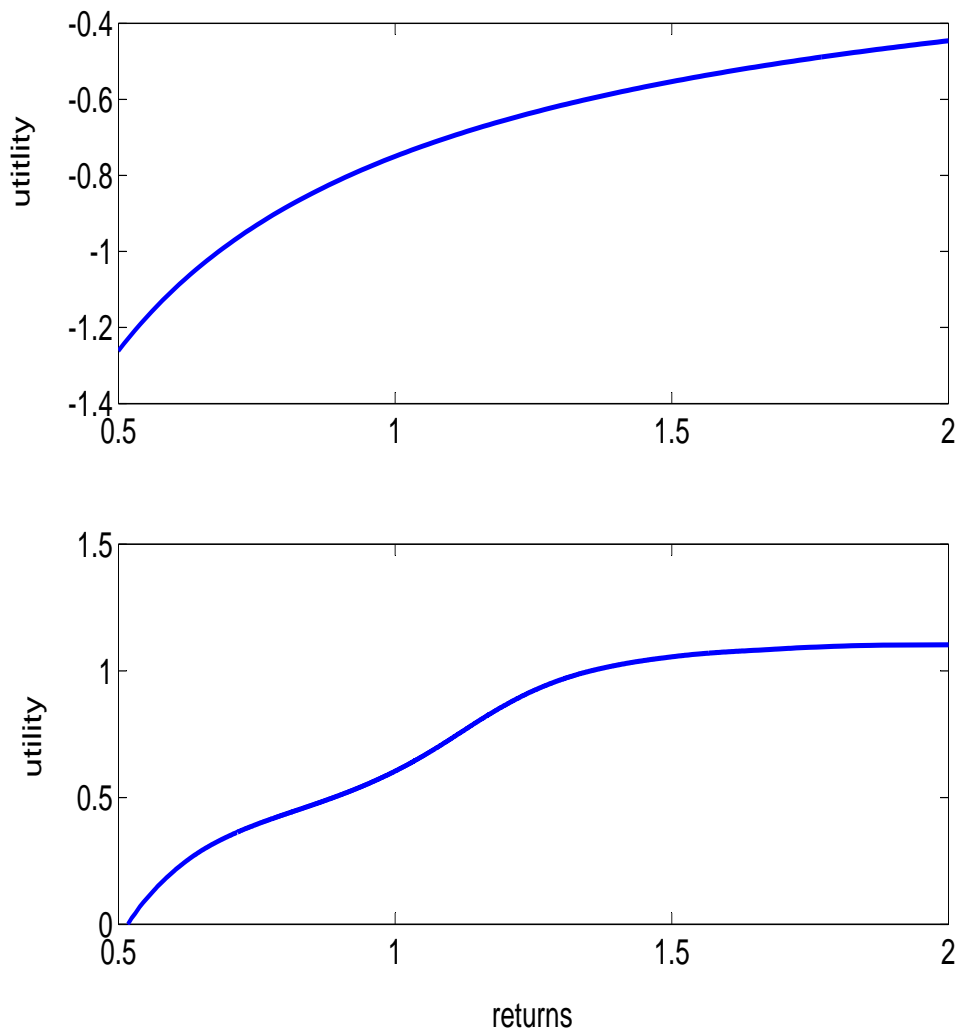


Figure 1: up: Utility function in the Black Scholes model for $T = 0.5$ years ahead and drift $\mu = 0.1$, volatility $\sigma = 0.2$ and interest rate $r = 0.03$. down: Market utility function on 06/30/2000 for $T = 0.5$ years ahead.

more interesting but it is hardly possible to estimate the historical density for a long time ahead. It neither seems realistic to assume that investors have clear ideas where the DAX will be in e.g. 10 years. For these reasons we use half a year as future horizon. Utility functions \tilde{U} of returns are defined by:

$$\tilde{U}(R) := U(RS_0), \quad R > 0$$

where S_0 denotes the value of the DAX on the day of estimation. Because of $U' = cK$ for a constant c we have $\tilde{U}'(R) = cK(RS_0)S_0$ and we see that also utility functions of returns are given as integrals of the pricing kernel. The change to returns allows us to compare different market regimes independently of the initial wealth. In the following we denote the utility functions of returns by the original notation U . Hence, we suppress in the notation the dependence of the utility function U on the day of estimation t .

The utility function corresponding to the model of Black and Scholes (1973) is a power utility, monotonically increasing and concave. But such classical utility functions are not observed on the market. Parametric and nonparametric models that replicate the option prices all lead to utility functions with a hump around the initial wealth level. This is described in detail later but is shown already in figure 1. The upper panel presents the utility function corresponding to Black-Scholes model with a volatility of 20% and an expected return of 10%. The function is concave and implies a constant relative risk aversion. The utility function estimated on the bullish market in summer 2000 is presented in the lower panel. Here, the hump around the money is clearly visible. The function is no more concave but has a region where investors are risk seeking. This risk proclivity around the money is reflected in a negative relative risk aversion.

3 Estimation

In this section, we start by reviewing some recent approaches for estimating the pricing kernel. Then we describe our method that is based on estimates of the risk neutral and the historical density. The risk neutral density is derived from option prices that are given by an implied volatility surface and the historical density is estimated from the independent data set of historical returns. Finally, we present the empirical pricing kernels and the inferred utility and relative risk aversion functions.

3.1 Estimation approaches for the pricing kernel

There exist several ways and methods to estimate the pricing kernel. Some of these methods assume parametric models while others use nonparametric techniques. Moreover, some methods estimate first the risk neutral and subjective density to infer the pricing kernel. Other approaches estimate directly the pricing kernel.

Ait-Sahalia and Lo (1998) derive a nonparametric estimator of the risk neutral density based on option prices. In Ait-Sahalia and Lo (2000), they consider the empirical pricing kernel and the corresponding risk aversion using this estimator. Moreover, they derive asymptotic properties of the estimator that allow e.g. the construction of confidence bands. The estimation procedure consists of two steps: First, the option price function is determined by nonparametric kernel regression and then the risk neutral density is computed by the formula of Breeden and Litzenberger (1978). Advantages of this approach are the known asymptotic properties of the estimator and the few assumptions necessary.

Jackwerth (2000) analyses risk aversion by computing the risk neutral density from option prices and the subjective density from historical data of the underlying. For the risk neutral distribution, he applies a variation of the estimation procedure described in Jackwerth and Rubinstein (1996): A smooth volatility function derived from observed option prices gives the risk neutral density by differentiating it twice. The subjective density is approximated by a kernel density computed from historical data. In this method bandwidths have to be chosen as in the method of Ait-Sahalia and Lo (1998).

Rosenberg and Engle (2002) use a different approach and estimate the subjective density and directly (the projection of) the pricing kernel. This gives the same information as the estimation of the two densities because the risk neutral density is the product of the pricing kernel and the subjective density. For the pricing kernel, they consider two parametric specifications as power functions and as exponentials of polynomials. The evolution of the underlying is modelled by GARCH processes. As the parametric pricing kernels lead to different results according to the parametric form used this parametric approach appears a bit problematic.

Chernov (2003) also estimates the pricing kernel without computing the risk neutral and subjective density explicitly. Instead of assuming directly a parametric form of the kernel he starts with a (multi dimensional) modified model of Heston (1993) and derives an analytic expression for the pricing kernel by the Girsanov theorem, see Chernov (2000) for details. The ker-

nel is estimated by a simulated method of moments technique from equity, fixed income and commodities data and by reprojection. An advantage of this approach is that the pricing kernel is estimated without assuming an equity index to approximate the whole market portfolio. But the estimation procedure is rather complex and model dependent.

In a recent paper, Barone-Adesi et al. (2004) price options in a GARCH framework allowing the volatility to differ between historical and risk neutral distribution. This approach leads to acceptable calibration errors between the observed option prices and the model prices. They estimate the historical density as a GARCH process and consider the pricing kernel only on one day. This kernel is decreasing which coincides with standard economic theory. But the general approach of changing explicitly the volatility between the historical and risk neutral distribution is not supported by the standard economic theory.

We estimate the pricing kernel in this paper by estimating the risk neutral and the subjective density and then deriving the pricing kernel. This approach does not impose a strict structure on the kernel. Moreover, we use accepted parametric models because nonparametric techniques for the estimation of second derivatives depend a lot on the bandwidth selection although they yield the same pricing kernel behaviour over a wide range of bandwidths. For the risk neutral density we use a stochastic volatility model that is popular both in academia and in industry. The historical density is more difficult to estimate because the drift is not fixed. Hence, the estimation depends more on the model and the length of the historical time series. In order to get robust results we consider different (discrete) models and different lengths. In particular, we use a GARCH model that is the discrete version of the continuous model for the risk neutral density. In the following, we describe these models, their estimation and the empirical results.

3.2 Estimation of the risk neutral density

Stochastic volatility models are popular in industry because they replicate the observed smile in the implied volatility surfaces (IVS) rather well and moreover imply rather realistic dynamics of the surfaces. Nonparametric approaches like the local volatility model of Dupire (1994) allow a perfect fit to observed price surfaces but their dynamics are in general contrary to the market. As Bergomi (2005) points out the dynamics are more important for modern products than a perfect fit. Hence, stochastic volatility models are popular.

We consider the model of Heston (1993) for the risk neutral density be-

cause it can be interpreted as the limit of GARCH models. The Heston model has been refined further in order to improve the fit, e.g. by jumps in the stock price or by a time varying mean variance level. We use the original Heston model in order to maintain a direct connection to GARCH processes. Although it is possible to estimate the historical density also with the Heston model e.g. by Kalman filter methods we prefer more direct approaches in order to reduce the dependence of the results on the model and the estimation technique.

The stochastic volatility model of Heston (1993) is given by the two stochastic differential equations:

$$\frac{dS_t}{S_t} = rdt + \sqrt{V_t}dW_t^1$$

where the variance process is modelled by a square-root process:

$$dV_t = \xi(\eta - V_t)dt + \theta\sqrt{V_t}dW_t^2$$

and W^1 and W^2 are Wiener processes with correlation ρ and r is the risk free interest rate. The first equation models the stock returns by normal innovations with stochastic variance. The second equation models the stochastic variance process as a square-root diffusion.

The parameters of the model all have economic interpretations: η is called the long variance because the process always returns to this level. If the variance V_t is e.g. below the long variance then $\eta - V_t$ is positive and the drift drives the variance in the direction of the long variance. ξ controls the speed at which the variance is driven to the long variance. In calibrations, this parameter changes a lot and makes also the other parameters instable. To avoid this problem, the reversion speed is kept fixed in general. We follow this approach and choose $\xi = 2$ as Bergomi (2005) does. The volatility of variance θ controls mainly the kurtosis of the distribution of the variance. Moreover, there are the initial variance V_0 of the variance process and the correlation ρ between the Brownian motions. This correlation models the leverage effect: When the stock goes down then the variance goes up and vice versa. The parameters also control different aspects of the implied volatility surface. The short (long) variance determines the level of implied volatility for short (long) maturities. The correlation creates the skew effect and the volatility of variance controls the smile.

The variance process remains positive if the volatility of variance θ is small enough with respect to the product of the mean reversion speed ξ and

the long variance level η (i.e. $2\xi\eta > \theta^2$). As this constraint leads often to significantly worse fits to implied volatility surfaces it is in general not taken into account and we follow this approach.

The popularity of this model can probably be attributed to the semiclosed form of the prices of plain vanilla options. Carr and Madan (1999) showed that the price $C(K, T)$ of a European call option with strike K and maturity T is given by

$$C(K, T) = \frac{\exp\{-\alpha \ln(K)\}}{\pi} \int_0^{+\infty} \exp\{-\mathbf{i}v \ln(K)\} \psi_T(v) dv$$

for a (suitable) damping factor $\alpha > 0$. The function ψ_T is given by

$$\psi_T(v) = \frac{\exp(-rT) \phi_T\{v - (\alpha + 1)\mathbf{i}\}}{\alpha^2 + \alpha - v^2 + \mathbf{i}(2\alpha + 1)v}$$

where ϕ_T is the characteristic function of $\log(S_T)$. This characteristic function is given by

$$\begin{aligned} \phi_T(z) &= \exp\left\{\frac{-(z^2 + \mathbf{i}z)V_0}{\gamma(z) \coth \frac{\gamma(z)T}{2} + \xi - \mathbf{i}\rho\theta z}\right\} \\ &\times \frac{\exp\left\{\frac{\xi\eta T(\xi - \mathbf{i}\rho\theta z)}{\theta^2} + \mathbf{i}zTr + \mathbf{i}z \log(S_0)\right\}}{\left(\cosh \frac{\gamma(z)T}{2} + \frac{\xi - \mathbf{i}\rho\theta z}{\gamma(z)} \sinh \frac{\gamma(z)T}{2}\right)^{\frac{2\xi\eta}{\theta^2}}} \end{aligned} \quad (4)$$

where $\gamma(z) \stackrel{\text{def}}{=} \sqrt{\theta^2(z^2 + \mathbf{i}z) + (\xi - \mathbf{i}\rho\theta z)^2}$, see e.g. Cizek et al. (2005).

For the calibration we minimize the absolute error of implied volatilities based on the root mean square error:

$$\text{ASE}_t \stackrel{\text{def}}{=} \sqrt{\sum_{i=1}^n n^{-1} \{IV_i^{\text{mod}}(t) - IV_i^{\text{mar}}(t)\}^2}$$

where *mod* refers to a model quantity, *mar* to a quantity observed on the market and $IV(t)$ to an implied volatility on day t . The index i runs over all n observations of the surface on day t .

It is essential for the error functional ASE_t which observed prices are used for the calibration. As we investigate the pricing kernel for half a year to maturity we use only the prices of options that expire in less than 1.5 years. In order to exclude liquidity problems occurring at expiry we consider for the

calibration only options with more than 1 month time to maturity. In the moneyness direction we restrict ourselves to strikes 50% above or below the spot for liquidity reasons.

The risk neutral density is derived by estimation of the model parameters by a least squares approach. This amounts to the minimization of the error functional ASE_t . Cont and Tankov (2004) provided evidence that such error functionals may have local minima. In order to circumvent this problem we apply a stochastic optimization routine that does not get trapped in a local minimum. To this end, we use the method of differential evolution developed by Storn and Price (1997).

Having estimated the model parameters we know the distribution of $X_T = \log S_T$ in form of the characteristic function ϕ_T , see (4). Then the corresponding density f of X_T can be recovered by Fourier inversion:

$$f(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{itx} \phi_T(t) dt,$$

see e.g. Billingsley (1995). This integral can be computed numerically.

Finally, the risk neutral density q of $S_T = \exp(X_T)$ is given as a transformed density:

$$q(x) = \frac{1}{x} f\{\log(x)\}.$$

This density q is risk neutral because it is derived from option prices and options are priced under the risk neutral measure. This measure is applied because banks replicate the payoff of options so that no arbitrage conditions determine the option price, see e.g. Rubinstein (1994). An estimated risk neutral density is presented in figure 2. It is estimated from the implied volatility shown in figure 3 for the day 24/03/2000. The distribution is right skewed and its mean is fixed by the martingale property. This implies that the density is low for high profits and high for high losses. Moreover, the distribution is not symmetrical around the neutral point where there are neither profits nor losses. For this and all the following estimations we approximate the risk free interest rates by the EURIBOR. On each trading day we use the yields corresponding to the maturities of the implied volatility surface. As the DAX is a performance index it is adjusted to dividend payments. Thus, we do not have to consider dividend payments explicitly.

3.3 Estimation of the historical density

While the risk neutral density is derived from option prices observed on the day of estimation we derive the subjective density from the historical time

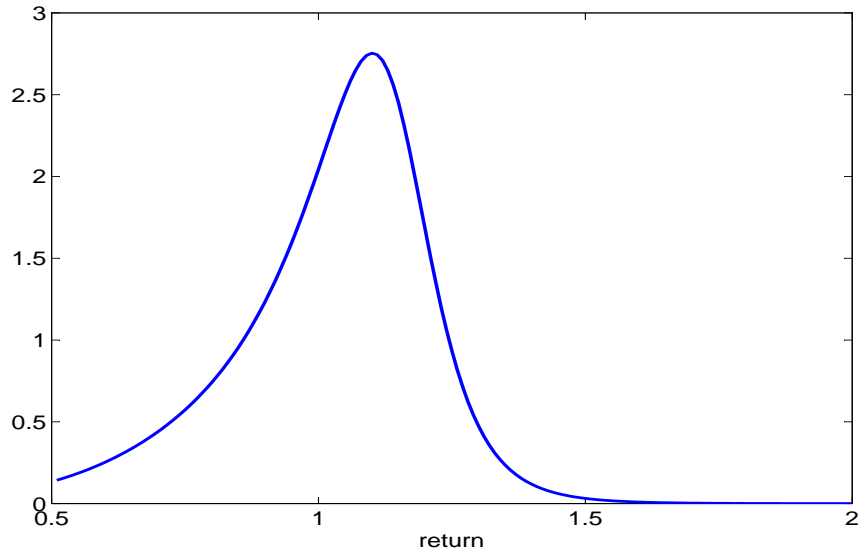


Figure 2: Risk neutral density on 24/03/2000 half a year ahead.

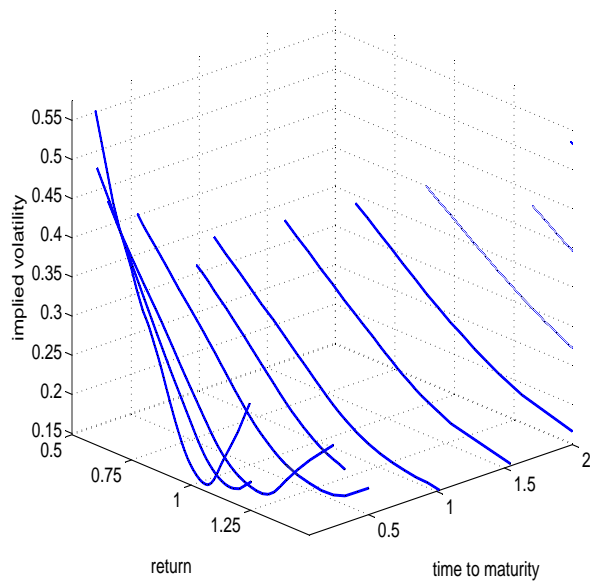


Figure 3: Implied volatility surface on 24/03/00.

model	time period
GARCH in mean	2.0y
discrete Heston	2.0y
observed returns	1.0y

Table 1: Models and the time periods used for their estimation.

series of the index. Hence, the two data sets are independent in the sense that the option prices reflect the future movements and the historical time series the past.

The estimation of the historical density seems more difficult than the estimation of the risk neutral density because the drift is not fixed and it depends in general on the length of the time series. Because of these difficulties we use different models and time horizons for the historical density: First, we estimate a GARCH in mean model for the returns. Returns are generally assumed to be stationary and we confirmed this at least in the time intervals we consider. The mean component in the GARCH model is important to reflect different market regimes. We estimate the GARCH model from the time series of the returns of the last two year because GARCH models require quite long time series for the estimation in order to make the standard error reasonably small. We do not choose longer time period for the estimation because we want to consider special market regimes. Besides this popular model choice we apply a GARCH model that converges in the limit to the Heston model that we used for the risk neutral density. As this model is also hard to estimate we use again the returns of the last 2 years for this model. Moreover, we consider directly the observed returns of the last year. The models and their time period for the estimation are presented in table 1. All these models give by simulation and smoothing the historical density for half a year ahead.

The GARCH estimations are based on the daily log-returns

$$R_i = \log(S_{t_i}) - \log(S_{t_{i-1}})$$

where (S_t) denotes the price process of the underlying and t_i , $i = 1, 2, \dots$ denote the settlement times of the trading days. Returns of financial assets have been analyzed in numerous studies, see e.g. Cont (2001). A model that has often been successfully applied to financial returns and their stylized facts

is the GARCH(1,1) model. This model with a mean is given by

$$\begin{aligned} R_i &= \mu + \sigma_i Z_i \\ \sigma_i^2 &= \omega + \alpha R_{i-1}^2 + \beta \sigma_{i-1}^2 \end{aligned}$$

where (Z_i) are independent identically distributed innovations with a standard normal distribution, see e.g. Franke et al. (2004). On day t_j the model parameters μ, ω, α and β are estimated by quasi maximum likelihood from the observations of the last two years, i.e. R_{j-504}, \dots, R_j assuming 252 trading days per year.

After the model parameters have been estimated on day t_j from historical data the process of logarithmic returns (R_i) is simulated half a year ahead, i.e. until time $t_j + 0.5$. In such a simulation μ, ω, α and β are given and the time series (σ_i) and (R_i) are unknown. The values of the DAX corresponding to the simulated returns are then given by inverting the definition of the log returns:

$$S_{t_i} = S_{t_{i-1}} \exp(R_i)$$

where we start with the observed DAX value on day t_j . Repeating the simulation N times we obtain N samples of the distribution of $S_{t_j+0.5}$. We use $N = 2000$ simulations because tests have shown that the results become robust around this number of simulations.

From these samples we estimate the probability density function of $S_{t_j+0.5}$ (given $(S_{t_j-126}, \dots, S_{t_j})$) by kernel density estimation. We apply the Gaussian kernel and choose the bandwidth by Silverman's rule of thumb, see e.g. Silverman (1986). This rule provides a trade-off between oversmoothing – resulting in a high bias – and undersmoothing – leading to big variations of the density. We have moreover checked the robustness of the estimate relative to this bandwidth choice. The estimation results of a historical density are presented in figure 4 for the day 24/03/2000. This density that represents a bullish market is has most of its weight in the profit region and its tail for the losses is relatively light.

As we use the Heston model for the estimation of the risk neutral density we consider in addition to the described GARCH model a GARCH model that is a discrete version of the Heston model. Heston and Nandi (2000) show that the discrete version of the square-root process is given by

$$V_i = \omega + \beta V_{i-1} + \alpha(Z_{i-1} - \gamma\sqrt{V_{i-1}})$$

and the returns are modelled by

$$R_i = \mu - \frac{1}{2}V_i + \sqrt{V_i}Z_i$$

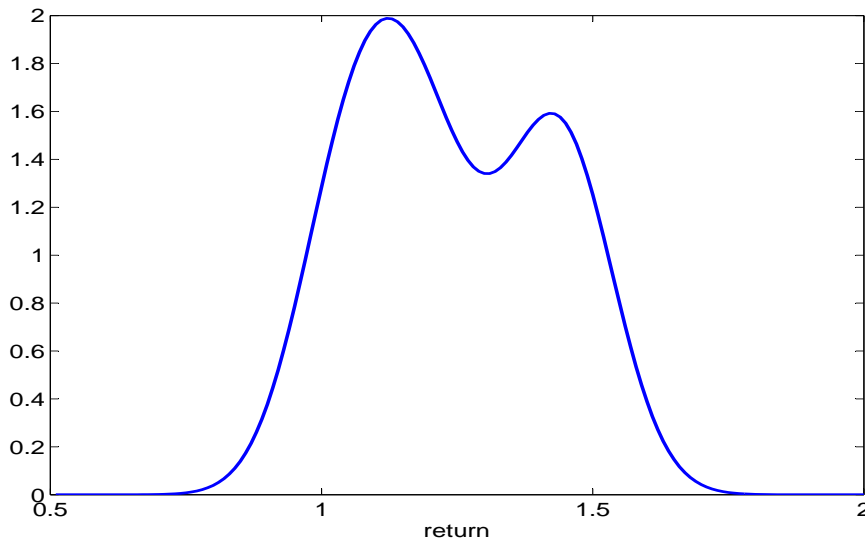


Figure 4: Historical density on 24/03/2000 half a year ahead.

where (Z_i) are independent identically distributed innovations with a standard normal distribution. Having estimated this model by maximum likelihood on day t_j we simulate it half a year ahead and then smooth the samples of $S_{t_j+0.5}$ in the same way as in the other GARCH model.

In addition to these parametric models, we consider directly the observed returns over half a year

$$\tilde{R}_i = S_{t_i}/S_{t_i-126}.$$

In this way, we interpret these half year returns as samples from the distribution of the returns for half a year ahead. Smoothing these historical samples of returns gives an estimate of the density of returns and in this way also an estimate of the historical density of $S_{t_j+0.5}$.

3.4 Empirical pricing kernels

In contrast to many other studies that concentrate on the S&P500 index we analyze the German economy by focusing on the DAX, the German stock index. This broad index serves as an approximation to the German economy. We use two data sets: A daily time series of the DAX for the estimation of the subjective density and prices of European options on the DAX for the estimation of the risk neutral density.

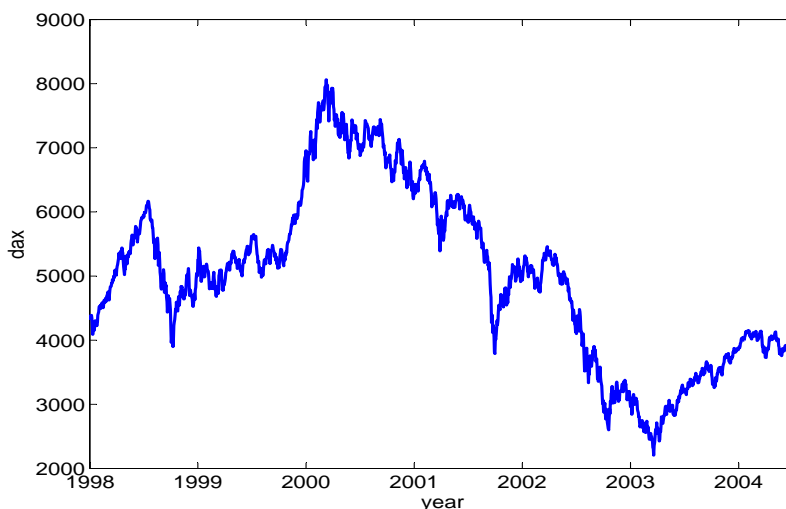


Figure 5: DAX, 1998 - 2004.

	1.0y	2.0y
03/2000	1.63	1.57
07/2002	0.66	0.54
06/2004	1.11	0.98

Table 2: Market regimes in 2000, 2002 and 2004 described by the return $S_0/S_{0-\Delta}$ for periods $\Delta = 1.0y, 2.0y$.

In figure 5, we present the DAX in the years 1998 to 2004. This figure shows that the index reached its peak in 2000 when all the internet firms were making huge profits. But in the same year this bubble burst and the index fell afterwards for a long time. The historical density is estimated from the returns of this time series. We analyze the market utility functions in March 2000, July 2002 and June 2004 in order to consider different market regimes. We interpret 2000 as a bullish, 2002 as a bearish and 2004 as a unsettled market. These interpretations are based on table 2 that describes the changes of the DAX over the preceding 1 or 2 years. (In June 2004 the market went up by 11% in the last 10 months.)

A utility function derived from the market data is a market utility function. It is estimated as an aggregate for all investors as if the representative investor existed. A representative investor is however just a convenient con-

struction because the existence of the market itself implies that the asset is bought and sold, i.e. at least two counterparties are required for each transaction.

In section 2 we identified the market utility function (up to linear transformations) as

$$U(R) = \int_{R_0}^R K(x) dx$$

where K is the pricing kernel for returns. It is defined by

$$K(x) = q(x)/p(x)$$

in terms of the historical and risk neutral densities p and q of returns. Any utility function (both cardinal and ordinal) can be defined up to a linear transformation, therefore we have identified the utility functions sufficiently. In section 3.3 we proposed different models for estimating the historical density. In figure 6 we show the pricing kernels resulting from the different estimation approaches for the historical density. The figure shows that all three kernels are quite similar: They have the same form, the same characteristic features like e.g. the hump and differ in absolute terms only a little. This demonstrates the economic equivalence of the three estimation methods on this day and this equivalence holds also for the other days. In the following we work with historical densities that are estimated by the observed returns.

Besides the pricing kernel and the utility function we consider also the risk attitudes in the markets. Such risk attitudes are often described in terms of relative risk aversion that is defined by

$$RRA(R) = -R \frac{U''(R)}{U'(R)}.$$

Because of $U' = cK = cq/p$ for a constant c the relative risk aversion is also given by

$$RRA(R) = -R \frac{q'(R)p(R) - q(R)p'(R)}{p^2(R)} \bigg/ \frac{q(R)}{p(R)} = R \left(\frac{p'(R)}{p(R)} - \frac{q'(R)}{q(R)} \right).$$

Hence, we can estimate the relative risk aversion from the estimated historical and risk neutral densities.

In figure 7 we present the empirical pricing kernels in March 2000, July 2002 and June 2004. The dates represent a bullish, a bearish and an unsettled markets, see table 2. All pricing kernels have a proclaimed hump located

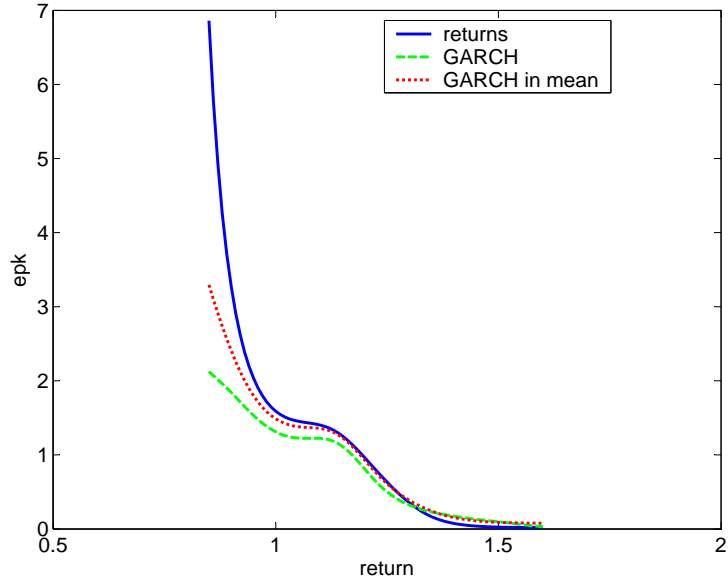


Figure 6: Empirical pricing kernel on 24/03/2000 (bullish market).

at small profits. Hence, the market utility functions do not correspond to standard specification of utility functions. We present the pricing kernels only in regions around the initial DAX (corresponding to a return of 1) value because the kernels explode outside these regions. This explosive behaviour reflects the typical pricing kernel form for losses. The explosion of the kernel for large profits is due to numerical problems in the estimation of the very low densities in this region. But we can see that in the unsettled market the kernel is concentrated on a small region while the bullish and bearish markets have wider pricing kernels. The hump of the unsettled market is also narrower than in the other two regimes. The bullish and bearish regimes have kernels of similar width but the bearish kernel is shifted to the loss region and the bullish kernel is located mainly in the profit area. Moreover, the figures show that the kernel is steeper in the unsettled markets than in the other markets. But this steepness cannot be interpreted clearly because pricing kernels are only defined up to a multiplicative constant.

The pricing kernels are the link between the relative risk aversion and the utility functions that are presented in figure 8. These utility functions are only defined up to linear transformations, see section 2. All the utility functions are increasing but only the utility function of the bullish market is concave. This concavity can be seen from the monotonicity of the kernel, see figure 7. Actually, this non convexity can be attributed to the quite special

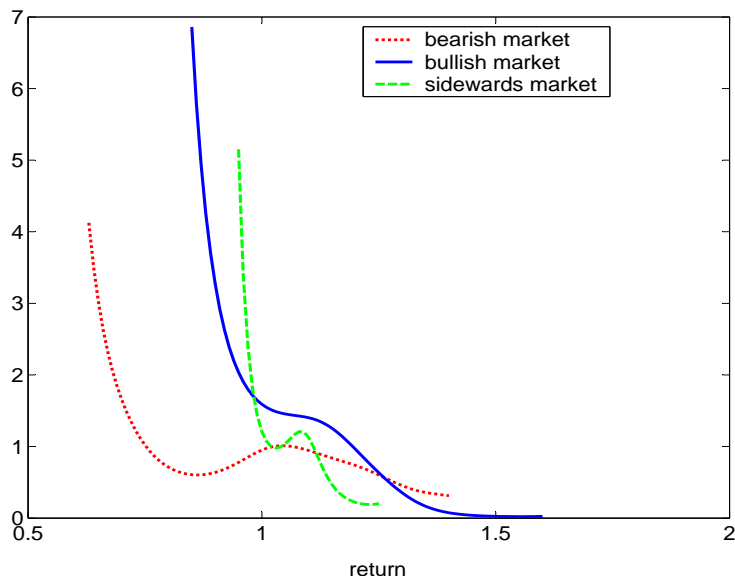


Figure 7: Empirical pricing kernel on 24/03/2000 (bullish), 30/07/2002 (bearish) and 30/06/2004 (unsettled or sideways market).

form of the historical density which has two modes on this date, see figure 4. Hence, we presume that also this utility function has in general a region of convexity. The other two utility functions are convex in a region of small profits where the bullish utility is almost convex. The derivatives of the utility functions cannot be compared directly because utility functions are identified only up to multiplicative constants. But we can compare the ratio of the derivatives in the loss and profit regions for the three dates because the constants cancel in these ratios. We see that the derivatives in the loss region are highest in the bullish and lowest in the bearish market and vice versa in the profit region. Economically these observations can be interpreted in such a way that in the bullish market a loss (of 1 unit) reduces the utility stronger than in the bearish market. On the other hand, a gain (of 1 unit) increases the utility less than in the bearish market. The unsettled market shows a behaviour between these extreme markets. Hence, investors fear in a good market situation losses more than in a bad situation and they appreciate profits in a good situation less than in a bad situation.

Finally, we consider the relative risk aversions in the three market regimes. These risk aversions are presented in figure 9, they do not depend on any constants but are completely identified. We see that the risk aversion is smallest in all markets for a small profit that roughly corresponds to the

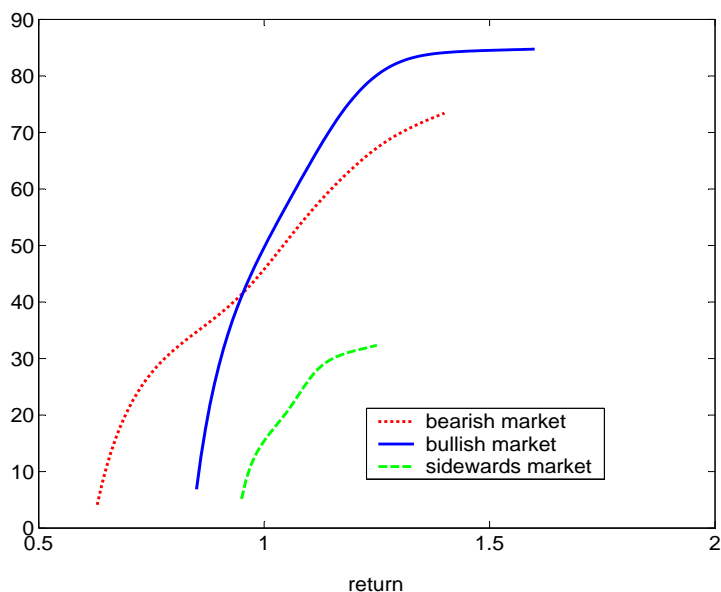


Figure 8: Market utility functions on 24/03/2000 (bullish), 30/07/2002 (bearish) and 30/06/2004 (unsettled or sideways market).

initial value plus a riskless interest on it. In the unsettled regime the market is risk seeking in a small region around this minimal risk aversion. But then the risk aversion increases quite fast. Hence, the representative agent in this market is willing to take small risks but is sensitive to large losses or profits. In the bullish and bearish regimes the representative agent is less sensitive to large losses or profits than in the unsettled market. In the bearish situation the representative agent is willing to take more risks than in the bullish regime. In the bearish regime the investors are risk seeking in a wider region than in the unsettled regime. In this sense they are more risk seeking in the bearish market. In the bullish market – on the other hand – the investors are never risk seeking so that they are less risk seeking than in the unsettled market.

The estimated utility functions most closely follow the specification proposed by Friedman & Savage (1948). The utility function proposed by Kahneman & Tversky (1979) consists of one concave and one convex segment and is less suitable for describing the observed behaviour, see figure 10. Both utility functions were proposed to account for two opposite types of behaviour with respect to risk attitudes: buying insurance and gambling. Any utility function that is strictly concave fails to describe both risk attitudes. Most notable examples are the quadratic utility function with the linear pricing

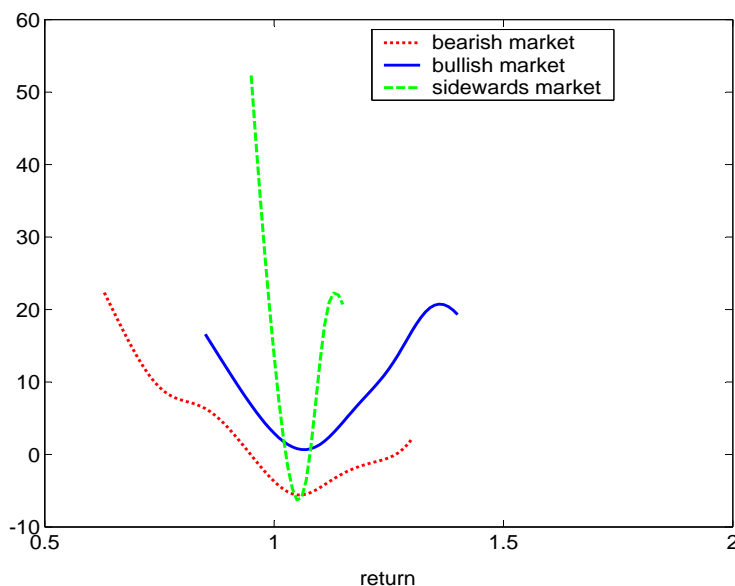


Figure 9: Relative risk aversions on 24/03/2000 (bullish), 30/07/2002 (bearish) and 30/06/2004 (unsettled or sideways market).

kernel as in the CAPM model and the CRRA utility function. These functions are presented in figure 10. Comparing this theoretical figure with the empirical results in figure 7 we see clearly the shortcoming of the standard specifications of utility functions to capture the characteristic hump of the pricing kernels.

4 Individual investors and their utility functions

In this section, we introduce a type of utility function that has two regions of different risk aversion. Then we describe how individual investors can be aggregated to a representative agent that has the market utility function. Finally, we solve the resulting estimation problem by discretization and estimate the distribution of individual investors.

4.1 Individual Utility Function

We learn from figures 10 and 7 that the market utility differs significantly from the standard specification of utility functions. Moreover, we can observe

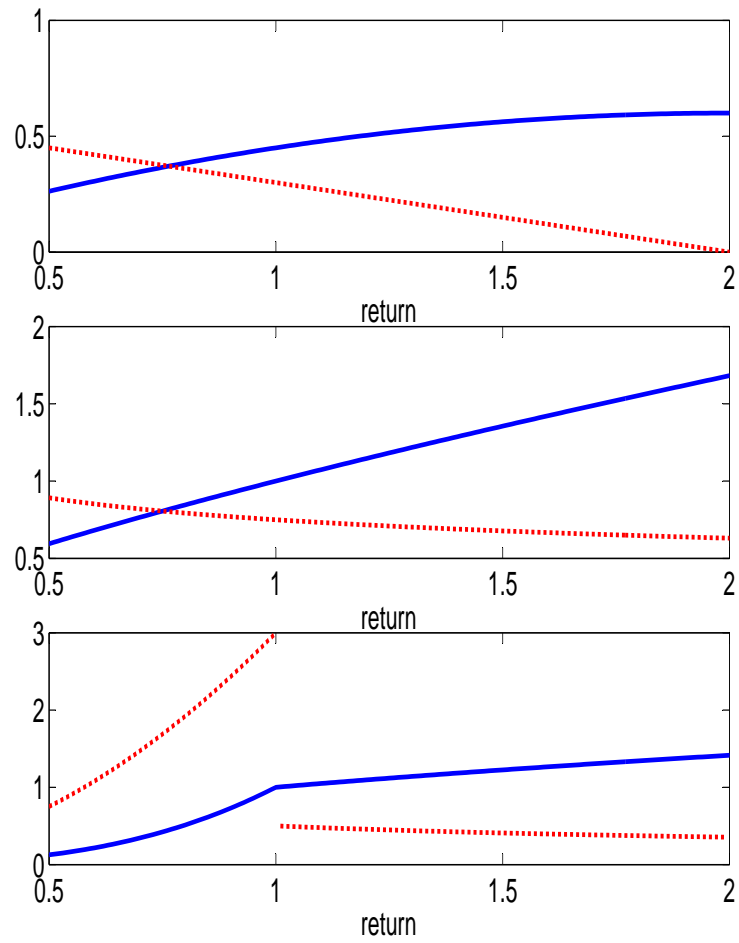


Figure 10: Common utility functions (solid) and their pricing kernels (dotted) (upper: quadratic, middle: power, lower panel: Kahneman and Tversky utility function).

from the estimated utility functions 8 that the loss part and the profit part of the utility functions can be quite well approximated with hyperbolic absolute risk aversion (HARA) functions, $k = 1, 2$:

$$U^{(k)}(R) = a_k(R - c_k)^{\gamma_k} + b_k,$$

where the shift parameter is c_k . These power utility functions become infinitely negative for $R = c_k$ and can be extended by $U^{(k)}(R) = -\infty$ for $R \leq c_k$, i.e. investors will avoid by all means the situation when $R \leq c_k$. The CRRA utility function has $c_k = 0$.

We try to reconstruct the market utility of the representative investor by individual utility functions and hence assume that there are many investors on the market. Investor i will be attributed with a utility function that consists of two HARA functions:

$$U_i(R) = \begin{cases} \max \{U(R, \theta_1, c_1); U(R, \theta_2, c_{2,i})\}, & \text{if } R > c_1 \\ -\infty, & \text{if } R \leq c_1 \end{cases}$$

where $U(R, \theta, c) = a(R - c)^\gamma + b$, $\theta = (a, b, \gamma)^\top$, $c_{2,i} > c_1$. If $a_1 = a_2 = 1$, $b_1 = b_2 = 0$ and $c_1 = c_2 = 0$, we get the standard CRRA utility function.

The parameters θ_1 and θ_2 and c_1 are the same for all investors who differ only with the shift parameter c_2 . θ_1 and c_1 are estimated from the lower part of the utility market function, where all investors probably agree that the market is “bad”. θ_2 is estimated from the upper part of the utility function where all investors agree that the state of the world is “good”. The distribution of c_2 uniquely defines the distribution of switching points and is computed in section 4.3. In this way a bear part $U_{bear}(R) = U(R, \theta_1, c_1)$ and a bull part $U_{bull}(R) = U(R, \theta_1, c_2)$ can be estimated by least squares.

The individual utility function can then be denoted conveniently as:

$$U_i(R) = \begin{cases} \max \{U_{bear}(R); U_{bull}(R, c_i)\}, & \text{if } R > c_1; \\ -\infty, & \text{if } R \leq c_1. \end{cases} \quad (5)$$

Switching between U_{bear} and U_{bull} happens at the *switching point* z , whereas $U_{bear}(z) = U_{bull}(z, c_i)$. The switching point is uniquely determined by $c_i \equiv c_{2,i}$. The notations *bear* and *bull* have been chosen because U_{bear} is activated when returns are low and U_{bull} when returns are high.

Each investor is characterised by a switching point z . The smoothness of the market utility function is the result of the aggregation of different attitudes. U_{bear} characterizes more cautious attitudes when returns are low and U_{bull} describes the attitudes when the market is booming. Both U_{bear}

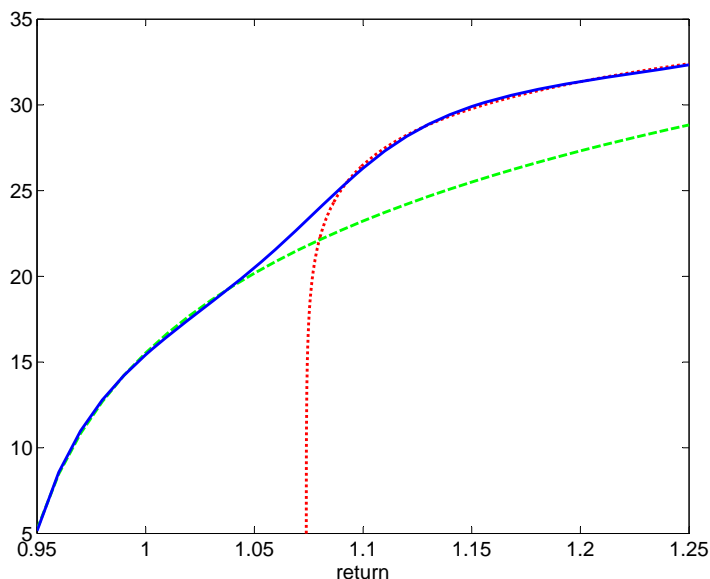


Figure 11: Market utility function (solid) with bearish (dashed) and bullish (dotted) part of an individual utility function 5 estimated in the unsettled market of 30/06/2004.

and U_{bull} are concave. However, due to switching the total utility function can be locally convex.

These utility functions are illustrated in figure 11 that shows the results for the unsettled market. We observe/estimate the market utility function that does not correspond to standard utility approaches because of the convex region. We propose to reconstruct this phenomenon by individual utility functions that consist of a bearish part and a bullish part. While the bearish part is fixed for all investors the bullish part starts at the switching point that characterizes an individual investor. By aggregating investors with different switching points we reconstruct the market utility function. We describe the aggregation in section 4.2 and estimate the distribution of switching points in section 4.3. In this way we explain the special form of the observed market utility functions.

4.2 Market Aggregation Mechanism

We consider the problem of aggregating individual utility functions to a representative market utility function. A simple approach to this problem is to identify the market utility function with an average of the individual utility functions. To this end one needs to specify the *observable* states of the world

in the future by returns R and then find a weighted average of the utility functions for each state. If the importance of the investors is the same, then the weights are equal:

$$U(R) = \frac{1}{N} \sum_{i=1}^N U_i(R),$$

where N is the number of investors. The problem that arises in this case is that utility functions of different investors can not be summed up since they are incomparable.

Therefore, we propose an alternative aggregation technique. First we specify the *subjective* states of the world given by utility levels u and then aggregate the outlooks concerning the returns in the future R for each perceived state. For a *subjective* state described with the utility level U , such that

$$u = U_1(R_1) = U_2(R_2) = \dots = U_N(R_N)$$

the aggregate estimate of the resulting returns is

$$R_A(u) = \frac{1}{N} \sum_{i=1}^N U_i^{-1}(u) \tag{6}$$

if all investors have the same market power. The market utility function U_M resulting from this aggregation is given by the inverse R_A^{-1} .

In contrast to the naive approach described at the beginning of this section, this aggregation mechanism is consistent under transformations: if all individual utility functions are changed by the same transformation then the resulting market utility is also given by the transformation of the original aggregated utility. We consider the individual utility functions U_i and the resulting aggregate U_M . In addition, we consider the transformed individual utility functions $U_i^\phi(x) = \phi\{U_i(x)\}$ and the corresponding aggregate U_M^ϕ where ϕ is a transformation. Then the aggregation is consistent in the sense that $U_M^\phi = \phi(U_M)$. This property can be seen from

$$\begin{aligned} (U_M^\phi)^{-1}(u) &= \frac{1}{N} \sum_{i=1}^N (U_i^\phi)^{-1}(u) \\ &= \frac{1}{N} \sum_{i=1}^N U_i^{-1}\{\phi^{-1}(u)\} \\ &= U_M^{-1}\{\phi^{-1}(u)\} \end{aligned}$$

The naive aggregation is not consistent in the above sense as the following example shows: We consider the two individual utility functions $U_1(x) = \sqrt{x}$

and $U_2(x) = \sqrt{x}/2$ under the logarithmic transformation $\phi = \log$. Then the naively aggregated utility is given by $U_M(x) = 3\sqrt{x}/4$. Hence, the transformed aggregated utility is $\phi\{U_M(x)\} = \log(3/4) + \log(x)/2$. But the aggregate of the transformed individual utility functions is

$$\begin{aligned} U_M^\phi(x) &= \frac{1}{2} \{ \log(\sqrt{x}) + \log(\sqrt{x}/2) \} \\ &= \frac{1}{2} \log\left(\frac{1}{2}\right) + \log(x)/2. \end{aligned}$$

This implies that $U_M^\phi \neq \phi(U_M)$ in general.

This described aggregation approach can be generalized in two ways: If the individual investors have different market power then we use the corresponding weights w_i in the aggregation (6) instead of the uniform weights. As the number of market participants is in general big and unknown it is better to use a continuous density f instead of the discrete distributions given by the weights w_i . These generalizations lead to the following aggregation

$$R_A(u) = \int U^{-1}(\cdot, z)(u) f(z) dz$$

where $U(\cdot, z)$ is the utility function of investor z . We assume in the following that the investors have utility function of the form described in section 4.1. In the next section we estimate the distribution of the investors who are parametrized by z .

4.3 The Estimation of the Distribution of Switching Points

Using the described aggregation procedure, we consider now the problem of replicating the market utility by aggregating individual utility functions. To this end, we choose the parametric utility functions $U(\cdot, z)$ described in 4.1 and try to recover with them the market utility U_M . We do not consider directly the utility functions but minimize instead the distance between the inverse functions:

$$\min_f \left\| \int U^{-1}(\cdot, z) f(z) dz - U_M^{-1} \right\|_{L^2(\tilde{P})} \quad (7)$$

where \tilde{P} is image measure of the historical measure P on the returns under the transformation U_M . As the historical measure has the density p the

transformation theorem for densities implies that \tilde{P} has the density

$$\tilde{p}(u) = p\{U_M^{-1}(u)\}/U'_M\{U_M^{-1}(u)\}.$$

With this density the functional to be minimized in problem (7) can be stated as

$$\begin{aligned} & \int \left(\int U^{-1}(u, z) f(z) dz - U_M^{-1}(u) \right)^2 \tilde{p}(u) du \\ &= \int \left(\int U^{-1}(u, z) f(z) dz - U_M^{-1}(u) \right)^2 p\{U_M^{-1}(u)\}/U'_M\{U_M^{-1}(u)\} du \\ &= \int \left(\int U^{-1}(u, z) f(z) dz - U_M^{-1}(u) \right)^2 p\{U_M^{-1}(u)\}(U_M^{-1})'(u) du \end{aligned}$$

because the derivative of the inverse is given by $(g^{-1})'(y) = 1/g'\{g^{-1}(y)\}$. Moreover, we can apply integration by substitution to simplify this expression further

$$\begin{aligned} & \int \left(\int U^{-1}(u, z) f(z) dz - U_M^{-1}(u) \right)^2 p\{U_M^{-1}(u)\}(U_M^{-1})'(u) du \\ &= \int \left(\int U^{-1}\{U_M(x), z\} f(z) dz - x \right)^2 p(x) dx. \end{aligned}$$

For replicating the market utility by minimizing (7) we observe first that we have samples of the historical distribution with density p . Hence, we can replace the outer integral by the empirical expectation and the minimization problem can be restated as

$$\min_f \frac{1}{n} \sum_{i=1}^n \left(\int g\{U_M(x_i), z\} f(z) dz - x_i \right)^2$$

where x_1, \dots, x_n are the samples from the historical distribution and $g = U^{-1}$.

Replacing the density f by a histogram $f(z) = \sum_{j=1}^J \theta_j I_{B_j}(z)$ with bins B_j , $h_j = |B_j|$, the problem is transformed into

$$\min_{\theta_j} \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{j=1}^J \tilde{g}(i, j) \theta_j - x_i \right\}^2$$

where $\tilde{g}(i, j) = \int_{B_j} g\{U_M(x_i), z\} dz$.

Hence, the distribution of switching points can be estimated by solving the quadratic optimization problem

$$\begin{aligned} \min_{\theta_j} \frac{1}{n} \sum_{i=1}^n \left\{ \sum_{j=1}^J \tilde{g}(i, j) \theta_j - x_i \right\}^2, \\ \text{s.t.} \quad \theta_j \geq 0, \\ \sum_{j=1}^J \theta_j h_j = 1. \end{aligned}$$

Such quadratic optimization problems are well known and their solutions can be obtained using standard techniques, see e.g. Mehrotra (1992) or Wright (1998).

We present in figures 12–14 the estimated distribution of switching points in the bullish (24/03/2000), bearish (30/07/2002) and unsettled (30/06/2004) markets. The distribution density f was computed for 100 bins but we checked the broad range of binwidths. The width of the distribution varies greatly depending on the regularisation scheme, for example as represented by the number of bins. The location of the distribution maximum, however, remains constant and independent from the computational method.

The maximum and the median of the distribution, i.e. the returns at which half of investors have bearish and bullish attitudes, depend on the year. For example, in the bullish market (Figure 12) the peak of the switching point distribution is located in the area of high returns around $R = 1.07$ for half a year. On the contrary, in the bearish market (Figure 13) the peak of switching points is around $R = 0.93$. This means that when the market is booming, such as in year 1999–2000 prior to the dot-com crash, investors get used to high returns and switch to the bullish attitude only for comparatively high R 's. An overall high level of returns serves in this respect as a reference level and investors form their judgements about the market relative to it. Since different investors have different initial wealth, personal habits, attitudes and other factors that our model does not take into account, we have a distribution of switching points. In the bearish market the average level of returns is low and investors switch to bullish attitudes already at much lower R 's.

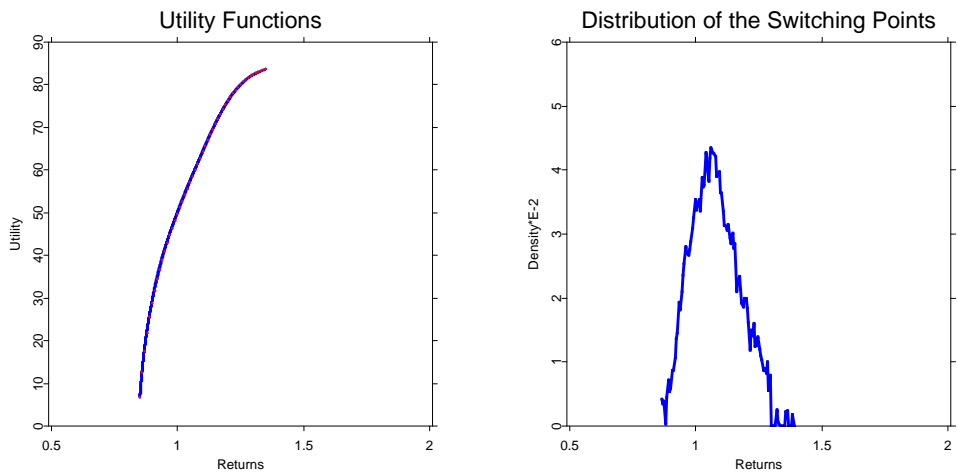


Figure 12: Left panel: the market utility function (red) and the fitted utility function (blue). Right panel: the distribution of the reference points. 24 March 2000, a bullish market.

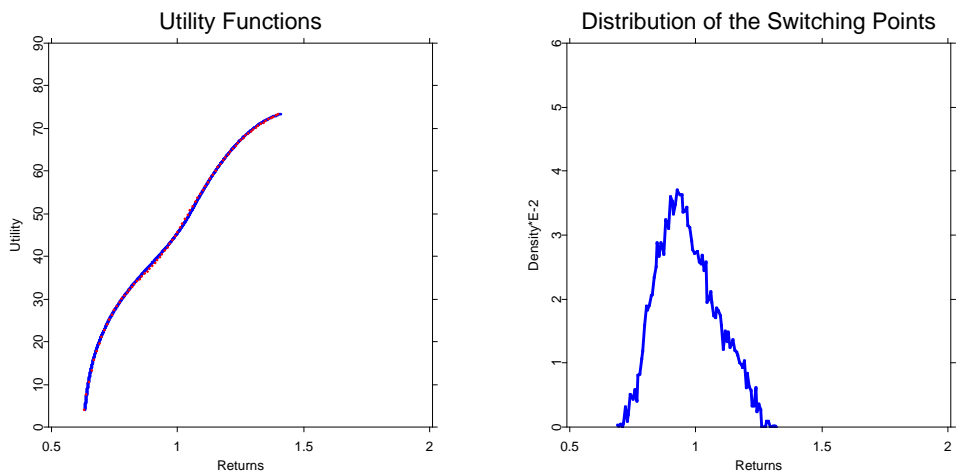


Figure 13: Left panel: the market utility function (red) and the fitted utility function (blue). Right panel: the distribution of the reference points. 30 July 2002, a bearish market.

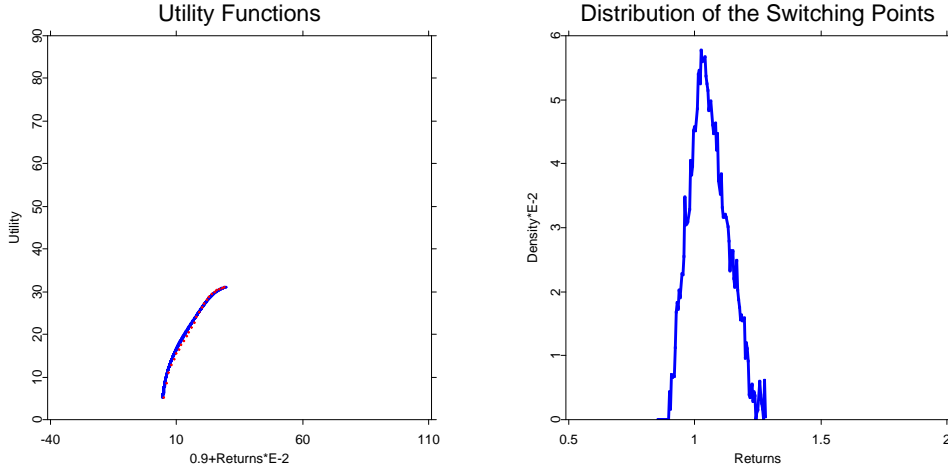


Figure 14: Left panel: the market utility function (red) and the fitted utility function (blue). Right panel: the distribution of the reference points. 30 June 2004, an unsettled market.

5 Conclusion

We have analyzed in this paper empirical pricing kernels in three market regimes using data on the German stock index and options on this index. In the bullish, bearish and unsettled market regime we estimate the pricing kernel and derive the corresponding utility functions and relative risk aversions.

In the unsettled market of June 2004, the market investor is risk seeking in a small region around the riskless return but risk aversion increases fast for high absolute returns. In the bullish market of March 2000, the investor is on the other hand never risk seeking while he becomes more risk seeking in the bearish market of July 2002. Before the stock market crash in 1987 European options did not show the smile and the Black-Scholes model captured the data quite well. Hence, utility functions could be estimated at that times by power utility functions with a constant positive risk aversion. Our analysis shows that this simple structure does not hold anymore and discusses different structures corresponding to different market regimes.

The empirical pricing kernels of all market regimes demonstrate that the corresponding utility functions do not correspond to standard specifications of utility functions including Kahneman and Tversky (1979). The observed utility functions are closest to the general utility functions of Friedman and Savage (1948). We propose a parametric specification of these functions,

estimate it and explain the observed market utility function by aggregating individual utility functions. In this way, we can estimate a distribution of individual investors.

The proposed aggregation mechanism is based on homogeneous investors in the sense that they differ only with switching points. Future research can reveal how nonlinear aggregation procedures could be applied to heterogeneous investors.

6 Acknowledgements

The research work of R. A. Moro was supported by the German Academic Exchange Service (DAAD). K. Detlefsen was supported by Bankhaus Sal. Oppenheim. This research was supported by Deutsche Forschungsgemeinschaft through the SFB 649 “Economic Risk”.

References

- Ait-Sahalia, Y. and A. Lo, 1998: Nonparametric estimation of state-price densities implicit in financial asset prices. *Journal of Finance*, **53**(2).
- Ait-Sahalia, Y. and A. Lo, 2000: Nonparametric risk-management and implied risk aversion. *Journal of Econometrics*, **94**(9).
- Barone-Adesi, G., R. Engle, and L. Mancini, 2004: Garch options in incomplete markets. working paper, University of Lugano.
- Bergomi, L., 2005: Smile dynamics 2. *Risk*, **18**(10).
- Bernoulli, D., 1956: Exposition of a new theory on the measurement of risk. *Econometrica*, **22**, 23–36.
- Billingsley, P., 1995: *Probability and Measure*. Wiley-Interscience.
- Black, F. and M. Scholes, 1973: The pricing of options and corporate liabilities. *Journal of Political Economy*, **81**, 637–659.
- Breedon, D. and R. Litzenberger, 1978: Prices of state-contingent claims implicit in option prices. *Journal of business*, **51**, 621–651.
- Carr, P. and D. Madan, 1999: Option valuation using the fast fourier transform. *Journal of Computational Finance*, **2**, 61–73.

- Chernov, M., 2000: Essays in financial econometrics. Phd thesis, Pennsylvania State University.
- Chernov, M., 2003: Empirical reverse engineering of the pricing kernel. *Journal of Econometrics*, **116**, 329–364.
- Cizek, P., W. Härdle, and R. Weron, 2005: *Statistical Tools in Finance and Insurance*. Springer, Berlin.
- Cochrane, J., 2001: *Asset Pricing*. Princeton University Press.
- Cont, R., 2001: Empirical properties of asset returns: stylized facts and statistical issues. 223–349.
- Cont, R. and P. Tankov, 2004: Nonparametric calibration of jump-diffusion option pricing models. *Journal of Computational Finance*, **7**(3), 1–49.
- Dupire, B., 1994: Pricing with a smile. *Risk*, **7**, 327–343.
- Franke, J., W. Härdle, and C. Hafner, 2004: *Statistics of Financial Markets*. Springer Verlag, Berlin.
- Friedman, M. and L. P. Savage, 1948: The utility analysis of choices involving risk. *Journal of Political Economy*, **56**, 279–304.
- Harrison, M. and S. Pliska, 1981: Martingales and stochastic integrals in the theory of continuous trading. *Stochastic Processes and their Applications*, **11**, 215–260.
- Heston, S., 1993: A closed-form solution for options with stochastic volatility with applications to bond and currency options. *Review of Financial Studies*, **6**(2), 327–343.
- Heston, S. and S. Nandi, 2000: A closed form garch option pricing model. *Review of Financial Studies*, **13**, 585–625.
- Jackwerth, J., 2000: Recovering risk aversion from option prices and realized returns. *Review of Financial Studies*, **13**(2), 433–451.
- Jackwerth, J. and M. Rubinstein, 1996: Recovering probability distributions from option prices. *Journal of Finance*, **51**(5), 1611–1631.
- Kahneman, D. and A. Tversky, 1979: Prospect theory: An analysis of decision under risk. *Econometrica*, **47**, 263–291.

- Mehrotra, S., 1992: On the implementation of a primal-dual interior point method. *SIAM Journal on Optimization*, **2**(4), 575–601.
- Merton, R. C., 1973: An intertemporal capital asset pricing model. *Econometrica*, **41**(5), 867–887.
- Rosenberg, J. and R. Engle, 2002: Empirical pricing kernels. *Journal of Financial Economics*, **64**(7), 341–372.
- Rubinstein, M., 1994: Implied binomial trees. *Journal of Finance*, **69**, 771–818.
- Silverman, B., 1986: *Density Estimation*. Chapman and Hall, London.
- Storn, R. and K. Price, 1997: Differential evolution - a simple and efficient heuristic for global optimization over continuous space. *Journal of Global Optimization*, **11**, 341–359.
- von Neumann, J. and O. Morgenstern, 1944: *The Theory of Games and Economic Behavior*. Princeton University Press.
- Wright, S., 1998: Primal-dual interior-point methods. *Mathematics of Computation*, **67**(222), 867–870.

De copulis non est disputandum*

Copulae: An Overview

Wolfgang Karl Härdle[†], Ostap Okhrin[‡]

May 27, 2009

Abstract: Normal distribution of the residuals is the traditional assumption in the classical multivariate time series models. Nevertheless it is not very often consistent with the real data. Copulae allows for an extension of the classical time series models to nonelliptically distributed residuals. In this paper we apply different copulae to the calculation of the static and dynamic Value-at-Risk of portfolio returns and Profit-and-Loss function. In our findings copula based multivariate model provide better results than those based on the normal distribution.

Keywords: copula; multivariate distribution; value-at-risk; multivariate dependence.

JEL Classification: C13, C14, C50.

1 Introduction

Understanding the joint distribution of high dimensional data is fundamental in applied statistics. The conventional procedure to model joint distributions is to approximate them with *multivariate normal distributions*.

That implies, however, that the dependence structures is reduced to a fixed type. Pre-termining a multivariate normal distribution means that the tails of the distribution are not too heavy, the distribution is symmetric and that the dependence between variables is linear.

Empirical evidence for these assumptions are barely verified and an alternative model is needed, with more flexible dependence structure and arbitrary marginal distributions. These are exactly the characteristics of *copulae*.

Copulae are very useful for modelling and estimating multivariate distributions. The flexibility of copulae basically follows from *Sklar's Theorem*, which says that each joint

*The financial support from the Deutsche Forschungsgemeinschaft via SFB 649 "Ökonomisches Risiko", Humboldt-Universität zu Berlin is gratefully acknowledged.

[†]CASE - Center for Applied Statistics and Economics, Institute for Statistics and Econometrics of Humboldt-Universität zu Berlin, Spandauer Straße 1, 10178 Berlin, Germany. Email: haerdle@wiwi.hu-berlin.de.

[‡]CASE - Center for Applied Statistics and Economics, Institute for Statistics and Econometrics of Humboldt-Universität zu Berlin, Spandauer Straße 1, 10178 Berlin, Germany. Email: ostap.okhrin@wiwi.hu-berlin.de

distribution can be “decomposed” into its marginal distributions and a copula C “responsible” for the dependence structure:

$$F(x_1, \dots, x_d) = C\{F_1(x_1), \dots, F_d(x_d)\}.$$

Two important factors for practical applications rely on this theorem:

1. The construction of multivariate distributions may be done in two independent steps: the specification of marginal distributions - not necessarily identical - and the specification of a dependence structure. Copulae “couple together” the marginal distributions into a multivariate distribution with the desired dependence structure.
2. Joint distributions can be separately estimated from a sample of observations: the marginal distributions are estimated first, the dependence structure later.

The copula approach gives us more freedom than the normality assumptions, marginal distributions with asymmetric heavy tails (typical for financial returns) can be combined with different dependence structures, resulting in multivariate distributions (far different from the multivariate normal) that better describe the empirical characteristics of financial returns distribution.

Moreover, copulae allow for dynamical modelling and adaption to portfolios, different copulae with distinct properties can be associated to different portfolios according to their specific dependence structures. Furthermore, copulae may change as time evolves, reflecting the evolution of the dependence between financial assets.

The structure of this paper is as follows. In the next section we give a short review of the copula theory. In the Section 3 we deal with different copula classes used in the calculation. The simulation and estimation techniques are provided in Sections 4 and 5 respectively. The first static problem on the calculation of the Value-at-Risk for the portfolio return has been discussed in Sections 6 and in the beginning of Section 7. Subsections 7.1 and 7.2 deal with the dynamic estimation of the Value-at-Risk for the Profit and Loss function. The paper is finished with summary.

2 Copulae

The description of copulae for measuring and modelling dependence with its main properties is the subject of this section. The term copula goes back to the works of Sklar (1959) where it was first mentioned. There are a lot of different equivalent definitions that could define the copula, but the most general is the following one.

Definition 1 (Copula) *A d -dimensional copula is a d -dimensional distribution with all uniform marginal distributions.*

Note that by considering random variables X_1, \dots, X_d with univariate distribution functions F_{X_1}, \dots, F_{X_d} and the random variables $U_i = F_{X_i}(X_i)$, $i = 1, \dots, d$ uniformly distributed in $[0, 1]$, a copula may be interpreted as *the joint distribution of the marginal distributions.*

Copulae gained popularity through Sklar's (1959) work where the term was first coined. However, many results had already been proved by Hoeffding (1940) and Hoeffding (1941), who could have been the founder of a copula theory, if he had considered the stochastically more intuitive dependency over the unit cube $[0, 1]^2$ rather than over $[-1/2, 1/2]^2$ as he had done. Copulae allow marginal distributions to be separated from the dependency structure. Sklar's theorem connects copulae with distribution functions such that from the one side every distribution function can be "decomposed" into its marginal distribution and (at least) one copula and from the other side a (unique) copula is obtained from "decoupling" every (continuous) multivariate distribution function from its marginal distributions.

Theorem 1 (Sklar's theorem) *Let F be a multivariate distribution function with margins F_1, \dots, F_d , then a copula C exists such that*

$$F(x_1, \dots, x_d) = C\{F_1(x_1), \dots, F_d(x_d)\}, \quad x_1, \dots, x_d \in \overline{\mathbb{R}}.$$

If F_i are continuous for $i = 1, \dots, d$ then C is unique. Otherwise C is uniquely determined on $F_1(\overline{\mathbb{R}}) \times \dots \times F_d(\overline{\mathbb{R}})$.

Conversely, if C is a copula and F_1, \dots, F_d are univariate distribution functions, then the function F defined above is a multivariate distribution function with margins F_1, \dots, F_d .

The representation in Sklar's Theorem can be used to construct new multivariate distributions by changing either the copula function or marginal distributions. For an arbitrary continuous multivariate distribution we can determine its copula from the transformation

$$C(u_1, \dots, u_d) = F\{F_1^{-1}(u_1), \dots, F_d^{-1}(u_d)\}, \quad u_1, \dots, u_d \in [0, 1], \quad (1)$$

where F_i^{-1} are inverse marginal distribution functions.

Since the copula function is a multivariate distribution with uniform margins, it follows that the copula density can be determined in the usual way

$$c(u_1, \dots, u_d) = \frac{\partial^d C(u_1, \dots, u_d)}{\partial u_1 \dots \partial u_d}, \quad u_1, \dots, u_d \in [0, 1],$$

Being armed with Theorem 1 and (??) we can write the density function $f(\cdot)$ of the d -variate distribution F in terms of copula as follows

$$f(x_1, \dots, x_d) = c\{F_1(x_1), \dots, F_d(x_d)\} \prod_{i=1}^d f_i(x_i), \quad x_1, \dots, x_d \in \overline{\mathbb{R}}.$$

A detailed discussion with proofs and deep mathematical treatment can be found in Joe (1997) and Nelsen (2006). A practical introduction is given in Deutsch and Eller (1999). Embrechts, McNeil and Straumann (1999b) discuss restrictions of the copula technique and their relation to the classical correlation analysis.

3 Copula Classes

Since there are plenty of functions satisfying the assumption of Theorem 1 they should be classified by construction and properties. Here we consider several main classes, like *simplest*, *elliptical*, *Archimedean copulae* and *hierarchical Archimedean copulae*.

3.1 Simplest Copulae

Special cases, like independence and perfect positive or negative dependence can be represented by copulae. If d random variables X_1, \dots, X_d are stochastically independent from Theorem 1, then the structure of such a relationship is given by the product copula

$$\Pi(u_1, \dots, u_d) = \prod_{j=1}^d u_j. \quad (2)$$

Copulae are bounded, this means that for all $u = (u_1, \dots, u_d)^\top \in [0, 1]^d$:

$$W(u_1, \dots, u_d) \leq C(u_1, \dots, u_d) \leq M(u_1, \dots, u_d)$$

where

$$M(u_1, \dots, u_d) = \min(u_1, \dots, u_d)$$

is called the *Fréchet-Hoeffding lower bound* and

$$W(u_1, \dots, u_d) = \max\left(\sum_{i=1}^d u_i - d + 1, 0\right)$$

is the *Fréchet-Hoeffding upper bound*. While M is not a copula for $d > 2$, W is a copula for all d . Both structures represent the perfect negative and perfect positive dependence. From this observation we may conclude that an arbitrary copula C reflects dependence which lies between the perfect negative and positive one.

3.2 Elliptical Copulae

The elliptical copulae are derived from the elliptical distributions using Theorem 1. In the bivariate case one has that a bivariate copula is elliptical if, and only if, it is equal to its associated copula

$$\begin{aligned} C(u_1, u_2, \theta) &= \bar{C}(u_1, u_2, \theta) \\ &= u_1 + u_2 - 1 + C(1 - u_1, 1 - u_2, \theta), \quad u_1, u_2 \in [0, 1]. \end{aligned}$$

The most prominent examples of elliptical copulae are Gaussian and t -copula.

Gaussian Copula

The Gaussian copula represents the *dependence structure* of the multivariate normal distribution, that means that *normal* marginal distributions are combined with a Gaussian copula to form multivariate normal distributions. The combination of *non-normal* marginal distributions with a Gaussian copula results in *meta-Gaussian* distributions, i.e., distributions where *only* the dependence structure is Gaussian.

To obtain the Gaussian copula, let $X = (X_1, \dots, X_d)^\top \sim N_d(\mu, \Sigma)$ with $X_j \sim N(\mu_j, \sigma_j)$ for $j = 1, \dots, d$. A copula C exists:

$$F(x_1, \dots, x_d) = C\{F_1(x_1), \dots, F_d(x_d)\},$$

where F_j is the distribution function of X_j and F the distribution function of X . Let $Y_j = T_j(X_j)$, $T_j(x) = (x - \mu_j)/\sigma_j$. Then $Y_j \sim N(0, 1)$ and $Y = (Y_1, \dots, Y_d)^\top \sim N_d(0, \Psi)$ where Ψ is the correlation matrix associated with Σ . A copula C_Ψ^{Ga} , called *Gaussian copula* exists as follows:

$$F_Y(y_1, \dots, y_d) = C_\Psi^{Ga}\{\Phi(y_1), \dots, \Phi(y_d)\}. \quad (3)$$

An explicit expression for the Gaussian copula is obtained by rewriting (3) with $u_j = \Phi(y_j)$:

$$\begin{aligned} C_\Psi^{Ga}(u_1, \dots, u_d) &= F_Y\{\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_d)\} \\ &= \int_{-\infty}^{\Phi^{-1}(u_1)} \dots \int_{-\infty}^{\Phi^{-1}(u_d)} (2\pi)^{-\frac{d}{2}} |\Psi|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}r^\top \Psi^{-1}r\right) dr_1 \dots dr_d. \end{aligned}$$

The density of the Gaussian copula is given by

$$c_\Psi^{Ga}(u_1, \dots, u_d) = |\Psi|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}\zeta^\top (\Psi^{-1} - I_d)\zeta\right\}. \quad (4)$$

Student's t -Copula

The t -copula, containing the dependence structure from the multivariate t -distribution, may be obtained in a similar way.

Let $X = (X_1, \dots, X_d)^\top \sim t_d(\nu, \mu, \Sigma)$ and $Y = (Y_1, \dots, Y_d)^\top \sim t_d(\nu, 0, \Psi)$ where Ψ is the correlation matrix associated with Σ . The unique copula from Y is the *Student's t -copula* $C_{\nu, \Psi}^t$. For $u = (u_1, \dots, u_d)^\top \in [0, 1]^d$, the *Student's t -copula* is given by

$$C_{\nu, \Psi}^t(u_1, \dots, u_d) = t_{\nu, \Psi}\{t_\nu^{-1}(u_1), \dots, t_\nu^{-1}(u_d)\}$$

where t_ν^{-1} is the quantile function from the univariate t -distribution and $t_{\nu, \Psi}$ the distribution function of Y .

The *density of the t -copula* is given by

$$\begin{aligned} c_{\nu, \Psi}^t(u_1, \dots, u_d) &= \frac{t_{\nu, \Psi}\{t_\nu^{-1}(u_1), \dots, t_\nu^{-1}(u_d)\}}{\prod_{j=1}^d t_{\nu, \Psi}\{t_\nu^{-1}(u_j)\}} \\ &= |\Psi|^{-\frac{1}{2}} \frac{\Gamma(\frac{\nu+d}{2}) \{\Gamma(\frac{\nu}{2})\}^{d-1} \left(1 + \frac{1}{\nu}\zeta^\top \Psi^{-1}\zeta\right)^{-\frac{\nu+d}{2}}}{\{\Gamma(\frac{\nu+1}{2})\}^d \prod_{j=1}^d \left(1 + \frac{1}{\nu}\zeta_j^2\right)^{-\frac{\nu+1}{2}}}. \end{aligned}$$

3.3 Archimedean Copulae

As opposed to elliptical copulae, Archimedean copulae are not constructed using Theorem 1, but are related to Laplace transforms of univariate distribution functions. Let \mathbb{L} denote the class of Laplace transforms which consists of strictly decreasing differentiable functions Joe (1997), i.e.

$$\mathbb{L} = \{\phi : [0; \infty) \rightarrow [0, 1] \mid \phi(0) = 1, \phi(\infty) = 0; (-1)^j \phi^{(j)} \geq 0; j = 1, \dots, \infty\}.$$

The function $C : [0, 1]^d \rightarrow [0, 1]$ defined as

$$C(u_1, \dots, u_d) = \phi\{\phi^{-1}(u_1) + \dots + \phi^{-1}(u_d)\}, \quad u_1, \dots, u_d \in [0, 1]$$

is a d -dimensional Archimedean copula, where $\phi \in \mathbb{L}$ and is called the *generator of the copula*. It is straightforward to show that $C(u_1, \dots, u_d)$ satisfies the conditions of Definition 1.

Some d -dimensional Archimedean copulae are presented below.

Frank (1979) copula, $0 \leq \theta < \infty$.

The first popular Archimedean copula is the so called Frank copula, which is the only elliptical Archimedean copula. Its generator and copula functions are

$$\begin{aligned} \phi(x, \theta) &= \theta^{-1} \log\{1 - (1 - e^{-\theta})e^{-x}\}, \quad 0 \leq \theta < \infty, x \in [0, \infty). \\ C_\theta(u_1, \dots, u_d) &= -\frac{1}{\theta} \log \left[1 + \frac{\prod_{j=1}^d \{\exp(-\theta u_j) - 1\}}{\{\exp(-\theta) - 1\}^{d-1}} \right]. \end{aligned}$$

The dependence becomes maximal when θ tends to infinity and independence is achieved when $\theta = 0$.

Gumbel (1960) copula, $1 \leq \theta < \infty$.

The Gumbel copula is frequently used in financial applications. Its generator and copula functions are

$$\begin{aligned} \phi(x, \theta) &= \exp\{-x^{1/\theta}\}, \quad 1 \leq \theta < \infty, x \in [0, \infty) \\ C_\theta(u_1, \dots, u_d) &= \exp \left[- \left\{ \sum_{j=1}^d (-\log u_j)^\theta \right\}^{\theta^{-1}} \right]. \end{aligned}$$

Consider a bivariate distribution based on the Gumbel copula with univariate extreme value marginal distributions. Genest and Rivest (1989) showed that this distribution is

the only bivariate extreme value distribution based on an Archimedean copula. Moreover, all distributions based on Archimedean copulae belong to its domain of attraction under common regularity conditions. In contrary to the elliptical copulae, the Gumbel copula leads to asymmetric contour diagrams. The Gumbel copula shows stronger linkage between positive values, however, it also shows more variability and more mass in the negative tail.

For $\theta > 1$ this copula allows for the generation of dependence in the upper tail. For $\theta \rightarrow 1$, the Gumbel copula reduces to the product copula and for $\theta \rightarrow \infty$ we obtain the Fréchet-Hoeffding upper bound.

Clayton (1978) copula, $-1 \leq \theta < \infty$, $\theta \neq 0$.

The Clayton copula which, in contrast to the Gumbel copula, has more mass on the lower tail, and less on the upper. The generator and copula function are

$$\begin{aligned}\phi(x, \theta) &= (\theta x + 1)^{-\frac{1}{\theta}}, \quad -1 \leq \theta < \infty, \theta \neq 0, x \in [0, \infty), \\ C_\theta(u_1, \dots, u_d) &= \left\{ \left(\sum_{j=1}^d u_j^{-\theta} \right) - d + 1 \right\}^{-\theta^{-1}}.\end{aligned}$$

The Clayton copula is one of few copulae that has a simple explicit form of density for any dimension

$$c_\theta(u_1, \dots, u_d) = \prod_{j=1}^d \{1 + (j-1)\theta\} u_j^{-(\theta+1)} \left(\sum_{j=1}^d u_j^{-\theta} - d + 1 \right)^{-(\theta^{-1}+d)}.$$

As the parameter θ tends to infinity, dependence becomes maximal and as θ tends to zero, we have independence. As $\theta \rightarrow -1$, the distribution tends to the lower Fréchet bound.

3.4 Hierarchical Archimedean Copulae

A recently developed flexible method is provided by hierarchical Archimedean copulae (HAC). The special, so called fully nested case of the copula function is:

$$\begin{aligned}C(u_1, \dots, u_d) &= \phi_{d-1} \{ \phi_{d-1}^{-1} \circ \phi_{d-2} (\dots [\phi_2^{-1} \circ \phi_1 \{ \phi_1^{-1}(u_1) + \phi_1^{-1}(u_2) \} \\ &\quad + \phi_2^{-1}(u_3)] + \dots + \phi_{d-2}^{-1}(u_{d-1})) + \phi_{d-1}^{-1}(u_d) \} \\ &= \phi_{d-1} [\phi_{d-1}^{-1} \circ C(\{\phi_1, \dots, \phi_{d-2}\})(u_1, \dots, u_{d-1}) + \phi_{d-1}^{-1}(u_d)]\end{aligned}$$

for $\phi_{d-i}^{-1} \circ \phi_{d-j} \in \mathbb{L}^*$, $i < j$, where

$$\begin{aligned}\mathbb{L}^* &= \{ \omega : [0; \infty) \rightarrow [0, \infty) \mid \omega(0) = 0, \\ &\quad \omega(\infty) = \infty; (-1)^{j-1} \omega^{(j)} \geq 0; j = 1, \dots, \infty \}.\end{aligned}$$

In contrast to the Archimedean copula, the HAC defines the whole dependency structure in a recursive way. At the lowest level the dependency between the first two variables is

modelled by a copula function with the generator ϕ_1 , i.e. $z_1 = C(u_1, u_2) = \phi_1\{\phi_1^{-1}(u_1) + \phi_1^{-1}(u_2)\}$. At the second level another copula function is used to model the dependency between z_1 and u_3 , etc. Note that the generators ϕ_i can come from the same family and they differ only through the parameter or, to introduce more flexibility, they come from different generator families. As an alternative to the fully nested model, we can consider copula functions, with arbitrary chosen combinations at each copula level. Okhrin, Okhrin and Schmid (2009a) provide several methodologies in determining the structure of the HAC from the data. The case of $d = 3$ which we use further in applications is quite a simple one. If τ_{12}, τ_{13} and τ_{23} are Kendall's τ , pairwise rank correlation coefficients, we join together those X_i and X_j such that $\max_{i,j \in \{1,2,3\}, i \neq j} \tau_{ij}$. Next we introduce $z = \widehat{C}\{\widehat{F}_i(X_i), \widehat{F}_i(X_j)\}$. Estimation techniques will be considered later. Variable X_{i^*} , $i^* \in \{1, 2, 3\}/\{i, j\}$ is joined afterwards with the z .

Whelan (2004) provides tools for generating samples from Archimedean copulae, Savu and Trede (2006) derived the density of such copulae and Joe (1997) proves their positive quadrant dependence (see Theorem 4.4). Okhrin et al. (2009a) and Okhrin, Okhrin and Schmid (2009b) considered methods for determining the optimal structure of the HAC, provided asymptotic theory for the estimated parameters and derive theoretical properties of this copula family.

4 Monte Carlo Simulation

The Monte-Carlo simulation is often a single reliable solution to many financial problems. Within the simulation study the random variables are generated from some prescribed distributions. There are numerous methods of simulating from copula-based distributions, see Frees and Valdez (1998), Whelan (2004), Marshall and Olkin (1988), McNeil (2008), Embrechts, McNeil and Straumann (1999), Frey and McNeil (2003), Devroye (1986), etc. Here we focus on two of them, on the conditional inversion method and on the method proposed by Marshall and Olkin (1988) for Archimedean copulae with generalizations to hierarchical Archimedean copulae by McNeil (2008).

4.1 Conditional Inverse Method

The simulation from d pseudo random variables with joint distribution defined by a copula C and d marginal distributions F_j , $j = 1, \dots, d$, may follow different techniques.

Defining the copula j -dimensional marginal distribution C_j for $j = 2, \dots, d-1$ as $C_j(u_1, \dots, u_j) = C(u_1, \dots, u_j, 1, \dots, 1)$ and the derivative of C_j with respect to the first $j-1$ arguments as

$$c_{j-1}^j(u_1, \dots, u_j) = \frac{\partial^{j-1} C_j(u_1, \dots, u_j)}{\partial u_1, \dots, \partial u_{j-1}}$$

the probability $P(U_j \leq u_j, U_1 = u_1, \dots, U_{j-1} = u_{j-1})$ can be written as

$$\begin{aligned} \lim_{\Delta u_1, \dots, \Delta u_{j-1} \rightarrow 0} \frac{C_j(u_1 + \Delta u_1, \dots, u_{j-1} + \Delta u_{j-1}, u_j) - C_j(u_1, \dots, u_j)}{\Delta u_1, \dots, \Delta u_{j-1}} \\ = c_{j-1}^j(u_1, \dots, u_j). \end{aligned}$$

Thus, the conditional distribution $\Lambda(u_j)$ (given fixed u_1, \dots, u_{j-1}) is a function of the ratio of derivatives:

$$\begin{aligned}\Lambda(u_j) &= P(U_j \leq u_j \mid U_1 = u_1, \dots, U_{j-1} = u_{j-1}) \\ &= \frac{c_{j-1}^j(u_1, \dots, u_j)}{c_{j-1}^{j-1}(u_1, \dots, u_{j-1})}.\end{aligned}$$

The generation of d pseudo random numbers with given marginal distributions F_j , $j = 1, \dots, d$ and dependence structure given by the copula C follows the steps:

1. generate iid $v_1, \dots, v_d \sim U[0, 1]$.
2. for $j = 1, \dots, d$ calculate $u_j = \Lambda^{-1}(v_j)$.
3. set $x_j = F_j^{-1}(u_j)$.

4.2 Marshal-Olkin Method

The Marshal-Olkin method is developed for the simulations only from Archimedean copulae. The idea this approach is based on the fact that the Archimedean copulae are derived from Laplace transforms. Let M be a univariate cdf of a positive random variable (so that $M(0) = 0$) and ϕ be the Laplace transform of M , i.e.

$$\phi(s) = \int_0^\infty \exp\{-sw\} dM(w), \quad s \geq 0.$$

For any univariate distribution function F , a unique distribution G exists:

$$F(x) = \int_0^\infty G^\alpha(x) dM(\alpha) = \phi\{-\log G(x)\}.$$

Considering d different univariate distributions F_1, \dots, F_d , we obtain

$$C(u_1, \dots, u_d) = \int_0^\infty \prod_{i=1}^d G_i^\alpha dM(\alpha) = \phi \left[\sum_{i=1}^d \phi^{-1}\{F_i(u_i)\} \right]$$

which is a multivariate distribution function. By replacing the product of univariate distributions G_i for $i = 1, \dots, d$ with an arbitrary copula function R we get:

$$C(u_1, \dots, u_d) = \int_0^\infty \dots \int_0^\infty R(G_1^\alpha, \dots, G_d^\alpha) dM(\alpha).$$

Note that for the classical Archimedean copula R is equal to a product copula.

One proceeds with the following three steps to make a draw from a distribution described by an Archimedean copula:

1. generate an observation u from M ;
2. generate an observations (v_1, \dots, v_d) from R ;

3. the generated vector is computed by $x_j = G_j^{-1}(v_j^{1/u})$.

This method works faster than the conditional inverse technique. The drawback is that the distribution M can be determined explicitly only for a few generator functions ϕ like, for example for the Frank, Gumbel and Clayton families. The same problem arises in the case of hierarchical copulae, where $\phi_i \circ \phi_{i+1}^{-1}$ should satisfy the properties of generator functions.

5 Copula Estimation

The estimation of a copula based multivariate distribution involves both the estimation of the copula parameters θ and the estimation of the margins F_j , $j = 1, \dots, d$, however all the parameters from the copula and from the margins could be also estimated in one step. The properties and goodness of the estimator of θ heavily depend on the estimators of F_j , $j = 1, \dots, d$. We distinguish between a parametric and a nonparametric specification of the margins. If we are interested only in the dependency structure, the estimator of $\{\delta_1, \dots, \delta_d, \theta\}$ should be independent of any parametric models for the margins. In practical applications, however, we are interested in a complete distribution model and, therefore, parametric models for margins are preferred.

For nonparametrically estimated margins, one may show the consistency and asymptotic normality of maximum-likelihood (ML) estimators and derive the moments of the asymptotic distribution. The ML estimation can be performed simultaneously for the parameters of the margins and of the copula function. Alternatively, a two-stage procedure can be applied, where we estimate the parameters of margins at the first stage and the copula parameters at the second stage.

Let X be a d -dimensional random variable with parametric univariate marginal distributions $F_j(x_j; \delta_j)$, $j = 1, \dots, d$. Further let a copula belong to a parametric family $\mathcal{C} = \{C_\theta, \theta \in \Theta\}$. The distribution of X can be expressed as

$$F(x_1, \dots, x_d) = C\{F_1(x_1; \delta_1), \dots, F_d(x_d; \delta_d); \theta\}$$

and its density as

$$f(x_1, \dots, x_d; \delta_1, \dots, \delta_d, \theta) = c\{F_1(x_1; \delta_1), \dots, F_d(x_d; \delta_d); \theta\} \prod_{j=1}^d f_j(x_j; \delta_j)$$

where $c(\cdot)$ is the copula density (??). For a sample of observations $\{x_t\}_{t=1}^T$, $x_t = (x_{1,t}, \dots, x_{d,t})^\top$ and a vector of parameters $\alpha = (\delta_1, \dots, \delta_d, \theta)^\top \in \mathbb{R}^{d+1}$ the likelihood function is given by

$$L(\alpha; x_1, \dots, x_T) = \prod_{t=1}^T f(x_{1,t}, \dots, x_{d,t}; \delta_1, \dots, \delta_d, \theta)$$

and the log-likelihood function by

$$\begin{aligned} \ell(\alpha; x_1, \dots, x_T) &= \sum_{t=1}^T \log c\{F_1(x_{1,t}; \delta_1), \dots, F_d(x_{d,t}; \delta_d); \theta\} \\ &+ \sum_{t=1}^T \sum_{j=1}^d \log f_j(x_{j,t}; \delta_j). \end{aligned}$$

The vector of parameters $\alpha = (\delta_1, \dots, \delta_d, \theta)^\top$ contains d parameters δ_j from the marginals and the copula parameter θ . All these parameters can be estimated *in one step*. For practical applications, however, a two step estimation procedure is more efficient.

5.1 FML – Full Maximum Likelihood Estimation

In the Maximum Likelihood estimation method (also called *full maximum likelihood*), the vector of parameters α is estimated in one single step through

$$\tilde{\alpha}_{FML} = \arg \max_{\alpha} \ell(\alpha)$$

The estimates $\tilde{\alpha}_{FML} = (\tilde{\delta}_1, \dots, \tilde{\delta}_d, \tilde{\theta})^\top$ solve

$$(\partial \ell / \partial \delta_1, \dots, \partial \ell / \partial \delta_d, \partial \ell / \partial \theta) = 0.$$

Following the standard theory on ML estimation it is efficient and asymptotically normal. However, it is often computationally demanding to solve the system simultaneously.

5.2 IFM – Inference for Margins

In the IFM (*inference for margins*) method, the parameters δ_j from the marginal distributions are estimated in the first step and used to estimate the dependence parameter θ in the second step:

1. for $j = 1, \dots, d$ the log-likelihood function for each of the marginal distributions are

$$\ell_j(\delta_j) = \sum_{t=1}^T \log f_j(x_{j,t}; \delta_j)$$

and the estimated parameters

$$\hat{\delta}_j = \arg \max_{\delta} \ell_j(\delta_j)$$

2. the *pseudo log-likelihood* function

$$\ell(\theta, \hat{\delta}_1, \dots, \hat{\delta}_d) = \sum_{t=1}^T \log c\{F_1(x_{1,t}; \hat{\delta}_1), \dots, F_d(x_{d,t}; \hat{\delta}_d); \theta\}$$

is maximised over θ to get the dependence parameter estimate $\hat{\theta}$.

The estimates $\hat{\alpha}_{IFM} = (\hat{\delta}_1, \dots, \hat{\delta}_d, \hat{\theta})^\top$ solve

$$(\partial \ell_1 / \partial \delta_1, \dots, \partial \ell_d / \partial \delta_d, \partial \ell / \partial \theta) = 0.$$

Detailed discussion on this method could be found in Joe and Xu (1996) Note, that this procedure does not lead to efficient estimators, however, as argued by Joe (1997) the loss in the efficiency is modest. The advantage of the inference for margins procedure lies in the dramatic reduction of the numerical complexity. Detailed discussion on the inference for margins procedure can be found in Joe and Xu (1996). Note, that this method does not lead to efficient estimators, however, as argued by Joe (1997) the loss in the efficiency is modest.

5.3 CML – Canonical Maximum Likelihood

In the CML (*canonical maximum likelihood*) method, the univariate marginal distributions are estimated through the edf \hat{F} . The asymptotic properties of the multistage estimators of θ do not depend explicitly on the type of the nonparametric estimator, but on its convergence properties. For $j = 1, \dots, d$

$$\hat{F}_j(x) = \frac{1}{T+1} \sum_{t=1}^T \mathbf{I}(x_{j,t} \leq x).$$

The *pseudo log-likelihood* function is

$$\ell(\theta) = \sum_{t=1}^T \log c\{\hat{F}_1(x_{1,t}), \dots, \hat{F}_d(x_{d,t}); \theta\}$$

and the copula parameter estimator $\hat{\theta}_{CML}$ is given by

$$\hat{\theta}_{CML} = \arg \max_{\theta} \ell(\theta).$$

Notice that the first step of the IMF and CML methods estimates the marginal distributions. After marginals are estimated, a *pseudo sample* $\{u_t\}$ of observations transformed in the unit d -cube is obtained and used in the *copula* estimation. As in the IFM, the semi-parametric estimator $\hat{\theta}$ is asymptotically normal under suitable regularity conditions.

6 Asset Allocation

We illustrate the extension of the classical asset allocation problem to copula-based models. We consider an investor with a CRRA utility function $U(x) = (1-\gamma)^{-1}x^{1-\gamma}$ willing to allocate his wealth to d risky assets. We denote the d -dimensional vector of d asset prices by $S_t = (S_{1,t}, \dots, S_{d,t})^\top$ and their continuously compounded asset returns at time $t+1$ by $X_{t+1} = (X_{1,t+1}, \dots, X_{d,t+1})^\top$ where $X_{t+1} = \log S_{t+1} - \log S_t$. The vector of portfolio weights by $w = (w_1, \dots, w_d)^\top$. Let F_{t+1} be the d -dimensional distribution function of X_{t+1} with the mean μ_{t+1} and covariance matrix Σ_{t+1} . The aim is to forecast F_{t+1} for the time period $t+1$ using the data up to time t . The estimator is denoted by \hat{F}_{t+1} with the mean $\hat{\mu}_{t+1}$, the covariance matrix $\hat{\Sigma}_{t+1}$ and the density \hat{f}_{t+1} . The objective of the investor is to maximise the expected utility at the time point $t+1$. This leads to the optimisation problem

$$\max_{w \in \mathcal{W}} \mathbf{E}_{\hat{F}_{t+1}} U(1 + w^\top X_{t+1}). \quad (5)$$

In the case of no short sales constraint we set $\mathcal{W} = \{w \in [0, 1]^d : w^\top \mathbf{1} = 1\}$ else we set $\mathcal{W} = \{w \in \mathbb{R}^d : w^\top \mathbf{1} = 1\}$. The conditional expectation in (5) implies that we integrate the utility with respect to the forecasted distribution \hat{F}_{t+1} . This reduces the problem (5) to the problem

$$\max_{w \in \mathcal{W}} \int \dots \int U(1 + w^\top X_{t+1}) \hat{f}_{t+1}(X_{t+1}) dX_{t+1}.$$

There are several alternative parametric approaches to modelling F_{t+1} . Let $\Sigma_{d,t+1}$ denote the diagonal matrix containing only the main diagonal of Σ_{t+1} . Then $\Sigma_{t+1} = \Sigma_{d,t+1}^{1/2} R_{t+1} \Sigma_{d,t+1}^{1/2}$, where R_{t+1} denotes the correlation matrix. A standard approach is to define the model of the asset returns in the form

$$\Sigma_{d,t}^{-1/2}(X_t - \mu_t) \sim N_d(0, R_t), \quad (6)$$

where the conditional moments μ_t and Σ_t are modelled by a GARCH type process.

To introduce a copula-based distribution into the asset allocation we deviate from the normality assumption and assume that $F = C(F_1, \dots, F_d)$. Thus (7) is replaced by:

$$\Sigma_{d,t}^{-1/2}(X_t - \mu_t) \sim C(F_1, \dots, F_d) \quad (7)$$

with some given functional forms of the copula and the marginal distributions. Similarly as above, the parameters of the conditional moments of the copula and of the marginal distributions are estimated using the ML method.

In Patton (2004) the investor allocates his wealth between small cap and large cap stocks (i.e. $d = 2$). The conditional mean is defined as linear function of the lagged asset returns and additional explanatory variables. The conditional variance is stated in the TAR(1,1) form. The rotated Gumbel copula with skewed t margins are used to construct the bivariate distribution of the residuals. This model reveals the highest likelihood function and the lowest AIC and BIC criterion. It is concluded that unconstrained portfolios derived from the normality assumption performed worse in 9 of 10 different trading strategies compared to the Gumbel model.

7 Value-at-Risk of the Portfolio Returns

If the return of the stock i at time point t is denoted as X_{it} then the portfolio value V at time t is defined recursively as

$$V_t = V_{t-1} \left(1 + \sum_{i=1}^d w_i X_{it} \right),$$

where w_i for $i = 1, \dots, d$ are the corresponding portfolio weights. Ruled with this notation the portfolio return is then given by

$$R_{tp} = \frac{V_t}{V_{t-1}} - 1 = \sum_{i=1}^d X_{it} w_i.$$

In our study we consider the case of equally weighted portfolio, i.e. $w_i = \frac{1}{d}$ for $i = 1, \dots, d$. The portfolio return is the random variable and its distribution strongly depends on the underlying distribution of the indices.

The distribution function of R_p , dropping the time index, is given by

$$F_{R_p}(\xi) = P(R_p \leq \xi). \quad (8)$$

One of the main advantages of copulae is the fact that they allow flexible modelling of the tail behaviour of multivariate distributions. Since the tail behaviour explains the

simultaneous outliers of asset returns, it is of special interest in risk management. The *Value-at-Risk* of a portfolio at level α is defined as the lower α -quantile of the distribution of the portfolio return, i.e.

$$\text{VaR}(\alpha) = F_{R_p}^{-1}(\alpha). \quad (9)$$

The VaR is a reasonable measure of risk if we assume that the returns are elliptically distributed. Moreover, the assumption of ellipticity implies that minimising the variance in the Markowitz problem also minimises the VaR, the expected shortfall and any other coherent measure of risk. However, this statement is false in the non-elliptical case. Moreover, regarding the effect of diversification the variance is the smallest (highest) for perfect negative (positive) correlation of the assets. This also holds for the VaR in the elliptical case, however, not for the non-elliptical distributions. This implies that for copula based distribution the VaR should be used with caution and its computation should be awarded more attention. Detailed description of the VaR estimation procedure at prescribed level α can be found in Giacomini and Härdle (2005).

Our aim is to determine such ξ that $P(R_p \leq \xi) = \alpha$. Note that

$$R_p = w^\top X = \sum_{i=1}^d w_i X_i = \sum_{i=1}^d w_i F_i^{-1}(u_i),$$

where F_i denotes the marginal distributions of individual asset returns, $u_i = F_i(X_i) \sim U[0, 1]$ for all $i = 1, \dots, d$ and $u_1, \dots, u_d \sim C$. The copula C defines the dependency structure between the asset returns. This implies that

$$F_{R_p}(\xi) = P(R_p \leq \xi) = \int_{\mathcal{U}} c(u_1, \dots, u_d) du_1 \dots du_d, \quad (10)$$

with

$$\mathcal{U} = \{[0, 1]^{d-1} \times [0, u_d(\xi)]\}, \quad u_d(\xi) = F_d\left\{\xi/w_d - \sum_{i=1}^{d-1} w_i F_i^{-1}(u_i)/w_d\right\}. \quad (11)$$

For fixed α , the VaR is determined by solving (10) numerically for ξ . Direct multidimensional numerical integration is a tedious task which can be substantially simplified by using the Monte-Carlo integration. For this purpose we have to generate random samples from C using the methods described in Section 4.

In the empirical study we consider four countries Canada, Germany, U.S. and U.K. from the MCSI index and eleven models of the joint multivariate distribution of indices, which include t -copula, Gaussian copula, simple exchangeable Archimedean copula, binary HAC and aggregated binary HAC, with normally and t -distributed margins. As a benchmark we use the empirical VaR, based purely on the real data.

In the cases where margins are t -distributed, we consider t -distribution with three degrees of freedom, while estimated t -distributions for this data are $t_{3.163}$, $t_{3.420}$, $t_{3.023}$, $t_{2.879}$. Multivariate t -copula in this case has eight degrees of freedom. Let us consider the simulation procedure, where on the first stage we estimate the covariance matrix $\widehat{\Sigma} = \{\widehat{\Sigma}_{ij}\}_{i,j=1,\dots,d}$, mean vector $\widehat{\mu} = \{\widehat{\mu}_i\}_{i=1,\dots,d}$ from the real data set and assume, or estimate, the marginal distributions $\widehat{F}_i(\cdot)$ (in our case they are normally or t -distributed), for $i =$

$1, \dots, d$. Next we show how to sample $u_1, \dots, u_d \in \mathcal{U}$ from (11). First we simulate the vector u of a dimension $d - 1$

$$u_1, \dots, u_{d-1} \sim U(0, 1).$$

Based on u we consider $x = \{x_i\}_{i=1, \dots, d-1}$ which for normal margins is equal to

$$x_i = \Phi^{-1}(u_i) \sqrt{\widehat{\Sigma}_{ii}} + \widehat{\mu}_i, \quad i = 1, \dots, d - 1,$$

and for t margins is

$$x_i = t^{-1}(u_i) \sqrt{\frac{\nu_i - 2}{\nu_i} \widehat{\Sigma}_{ii}} + \widehat{\mu}_i, \quad i = 1, \dots, d - 1,$$

where ν_i , $i = 1, \dots, d$ are degrees of freedom for marginal distributions. This transformation returns a normally or t -distributed vector x with the same parameters as the real data set.

Theoretically, in further steps we have to find bounds for the last stock (or index) to gain the portfolio ξ which is the α quantile. Thus, we separate our maximally reachable portfolio return ξ into two parts

$$\xi = \sum_{i=1}^{d-1} \frac{1}{d} X_i + \frac{1}{d} X_d,$$

then the return of the last index given the return of the portfolio is

$$X_d = d\xi - \sum_{i=1}^{d-1} X_i,$$

where the upper bound for our last value in vector u is then

$$u_d^* = \widehat{F}_d \left(d\xi - \sum_{i=1}^{d-1} x_i \right).$$

Value u_d^* is uniformly distributed on $[0, 1]$ and we simulate the last element of the vector $u_d \sim U(0, u_d^*)$.

As mentioned above, the goal is to compute (10) which for this setting is

$$F_{R_p}(\xi) = \int \cdots \int_{[0,1]^{d-1} \times [0, u_d^*]} c(u_1, \dots, u_d) du_1 \dots du_d.$$

Then by solving $F_{R_p}(\xi) = \alpha$ we find $R_\alpha = \text{VaR}(\alpha)$. In our study we solve the equations numerically using the golden section method. The integration is performed using the Monte-Carlo technique

$$P(\widehat{R_p} \leq \xi) = \frac{1}{n_s} \sum_{i=1}^{n_s} c(u_{1i}, \dots, u_{di})$$

where n_s is equal to 10^8 , α is set to be 1% and the values u_{1i}, \dots, u_{di} for $i = 1, \dots, n_s$ are simulated using the method described above. The precision of R is set at 0.00015.

Table 1: VaR for the 4-dimensional data set

	N	t_3
N	-0.0194	-0.0210
t_8	-0.0199	-0.0213
<i>AC</i>	<i>-0.0174</i>	<i>-0.0154</i>
<i>HAC</i> _{binary}	-0.0187	-0.0194
<i>HAC</i> _{binary aggr.}	-0.0188	-0.0194
Empirical	-0.0235	

The final results for all methods are given in Table 1. In the left-hand column we provide the models with normal margins and in the right-hand column with t margins. From top to bottom we have five different copula functions like Gaussian, t , simple Archimedean copula, binary HAC and binary aggregated HAC. The empirical VaR which is at the bottom of the table is derived from the empirical quantile. Bold fonts in the table emphasize those results which are closest in absolute value to the empirical one in each column, and italic fonts the worst cases in absolute value.

As can be seen from Table 1, the results which are the best in absolute value are those returned by the model with t -copula and t margins. The model based on the simple Archimedean copula is the worst one. This is quite natural, since this copula needs exchangeability between variables, which is not observable here (see previous section). HAC with binary as well as aggregated binary structures, unfortunately, give us results that are not much worse compared to t -copula and Gaussian copula. For VaR(0.01) the t -copula with t margins provided the best result.

7.1 VaR of the P&L

This sub-section introduces the main assumptions and steps necessary to estimate the VaR from a Profit and Loss of a linear portfolio using copulae. Static and time-varying methods and their VaR performance evaluation through backtesting are described below.

In this section w is the portfolio, which is represented by the number of assets for a specified stock in the portfolio, $w = \{w_1, \dots, w_d\}$, $w_i \in \mathbb{Z}$. The value V_t of the portfolio w is given non-recursively by

$$V_t = \sum_{j=1}^d w_j S_{j,t} \tag{12}$$

and the random variable

$$\begin{aligned} L_{t+1} &= (V_{t+1} - V_t) \\ &= \sum_{j=1}^d w_j S_{j,t} \{ \exp(X_{j,t+1}) - 1 \}. \end{aligned}$$

also called *profit and loss (P&L) function*, expresses the absolute change in the portfolio value in one period.

Similarly to the previous case, the distribution function of L , dropping the time index, is given by

$$F_L(x) = P(L \leq x). \quad (13)$$

As usual the *Value-at-Risk* at level α from a portfolio w is defined as the α -quantile from F_L :

$$\text{VaR}(\alpha) = F_L^{-1}(\alpha). \quad (14)$$

It follows from (13) that F_L depends on the d -dimensional distribution of log-returns F_X . In general, the *loss distribution* F_L depends on a random process representing the *risk factors* influencing the P&L from a portfolio. In the present case log-returns are a suitable risk factor choice. Thus, modelling their distribution is essential to obtain the quantiles from F_L .

Contrary to the previous section, here log-returns are assumed to be time-dependent, thus a log-returns process $\{X_t\}$ can be modelled as

$$X_{j,t} = \mu_{j,t} + \sigma_{j,t}\varepsilon_{j,t}$$

where $\varepsilon_t = (\varepsilon_{1,t}, \dots, \varepsilon_{d,t})^\top$ are standardised *i.i.d.* innovations with $E[\varepsilon_{j,t}] = 0$ and $E[\varepsilon_{j,t}^2] = 1$ for $j = 1, \dots, d$; \mathcal{F}_t is the available information at time t :

$$\mu_{j,t} = E[X_{j,t} \mid \mathcal{F}_{t-1}]$$

is the conditional mean given \mathcal{F}_{t-1} and

$$\sigma_{j,t}^2 = E[(X_{j,t} - \mu_{j,t})^2 \mid \mathcal{F}_{t-1}]$$

is the conditional variance given \mathcal{F}_{t-1} . The innovations $\varepsilon = (\varepsilon_1, \dots, \varepsilon_d)^\top$ have joint distribution

$$F_\varepsilon(\varepsilon_1, \dots, \varepsilon_d) = C_\theta\{F_1(\varepsilon_1), \dots, F_d(\varepsilon_d)\}, \quad (15)$$

where C_θ is a copula belonging to a parametric family $\mathcal{C} = \{C_\theta, \theta \in \Theta\}$, and F_j , $j = 1, \dots, d$ are continuous marginal distributions of ε_j . To obtain the Value-at-Risk in this set up, the dependence parameter and distribution function from residuals are estimated from a sample of log-returns and used to generate P&L Monte Carlo samples. Their quantiles at different levels are the estimators for the Value-at-Risk.

For a portfolio w on d assets and a sample $\{x_{j,t}\}_{t=1}^T$, $j = 1, \dots, d$ of log-returns, the Value-at-Risk at level α is estimated according to the following steps:

1. Estimation of residuals $\hat{\varepsilon}_t$ from the prespecified time-series model;
2. Specification and estimation of marginal distributions $F_j(\hat{\varepsilon}_j)$;
3. Specification of a parametric copula family \mathcal{C} and estimation of dependence parameter θ ;
4. Generation of Monte Carlo sample of innovations ε and losses L , for the forecast on the one day;

5. Estimation of $\widehat{VaR}(\alpha)$, the empirical α -quantile from the forecasted L .

The application of the (*static*) procedure described above on sliding windows of a time series $\{x_{j,t}\}_{t=1}^T$ delivers a sequence of parameters for a copula family. Hence the denomination *time-varying copulae*.

Using moving windows of size r in time t

$$\{x_t\}_{t=s-w+1}^s$$

for $s = r, \dots, T$, the procedure described in the section above generates the time series $\{\widehat{VaR}_t\}_{t=r}^T$ of Value-at-Risk and $\{\hat{\theta}_t\}_{t=r}^T$ dependence parameters estimates.

Afterwards *Backtesting* is used to evaluate the performance of the specified copula family \mathcal{C} . The estimated values for the VaR are compared with the true realisations $\{l_t\}$ of the P&L function, an *exceedance* occurring for each l_t smaller than $\widehat{VaR}_t(\alpha)$. The ratio of the number of exceedances to the number of observations gives the *exceedances ratio* $\hat{\alpha}$:

$$\hat{\alpha} = \frac{1}{T-r} \sum_{t=r}^T \mathbf{I}\{l_t < \widehat{VaR}_t(\alpha)\}.$$

The estimation methods described before are used on two portfolio, the first composed of 2 positions, the second of 3 positions. Different copulae are used in static and dynamic setups and their VaR performance is compared based on backtesting.

In this section, the Value-at-Risk of portfolios for two companies (Tyssenkrupp (TKA) and Volkswagen (VOW) from 01.12.1997 to 03.07.2007) is computed using different copulae.

Assuming the log-returns $\{X_{j,t}\}$ follow a GARCH(1,1) process we have

$$X_{j,t} = \mu_{j,t} + \sigma_{j,t}\varepsilon_{j,t}$$

where

$$\sigma_{j,t}^2 = \omega_j + \alpha_j \sigma_{j,t-1}^2 + \beta_j (X_{j,t-1} - \mu_{j,t-1})^2$$

and $\omega > 0$, $\alpha_j \geq 0$, $\beta_j \geq 0$, $\alpha_j + \beta_j < 1$.

The fit of a GARCH(1,1) model to the sample of log returns $\{x_t\}_{t=1}^T$, $X_t = (X_{1,t}, X_{2,t})^\top$, $T = 2500$, gives the estimates $\hat{\omega}_j$, $\hat{\alpha}_j$ and $\hat{\beta}_j$, as in Table 2, and empirical residuals $\{\hat{\varepsilon}_t\}_{t=1}^T$, where $\hat{\varepsilon}_t = (\hat{\varepsilon}_{1,t}, \hat{\varepsilon}_{2,t})^\top$. The marginal distributions are specified as normal, i.e., $\hat{\varepsilon}_j \sim N(\hat{\mu}_j, \hat{\sigma}_j)$ with parameters $\hat{\delta}_j = (\hat{\mu}_j, \hat{\sigma}_j)$ estimated from the data.

Figure 1 displays the Kernel density estimator of the residuals and of the normal density, estimated with an Quartic kernel. The dependence parameters are estimated for different copula families (Gaussian, Clayton and Gumbel). Residuals $\hat{\varepsilon}$ and fitted copulae (Gaussian, Clayton and Gumbel) are plotted in Figure 2.

In the dynamic approach, the empirical residuals are sampled in moving windows with a fixed size $r = 250$, $\{\hat{\varepsilon}_t\}_{t=s-r+1}^s$, for $s = r, \dots, T$. The time series from estimated dependence parameters for each copula family are in Figure 3.

The same portfolio compositions as in the static case are used to generate P&L samples. The series of estimated Value-at-Risk and the P&L function for selected portfolios are plotted in Figure 4, 5 and 6.

	$\hat{\mu}_j$	$\hat{\omega}_j$	$\hat{\alpha}_j$	$\hat{\beta}_j$	BL	KS
MRK	7.392e-04 (3.672e-04)	4.588e-06 (1.557e-06)	3.333e-02 (6.225e-03)	9.572e-01 (8.568e-03)	0.1285	1.255e-11
TKA	7.845e-04 (3.308e-04)	3.549e-06 (1.149e-06)	7.087e-02 (9.837e-03)	9.252e-01 (9.915e-03)	0.1360	4.189e-05
VOW	9.720e-04 (3.480e-04)	1.239e-05 (2.699e-06)	9.303e-02 (1.301e-02)	8.830e-01 (1.566e-02)	1.927e-05	3.422e-06

Table 2: Fitting of univariate GARCH(1,1) to asset returns. The standard deviation of the parameters are given in parentheses. The last two columns provide the p -values of the Box-Ljung test (BL) for autocorrelations and Kolmogorov-Smirnov test (KS) for normality applied to the residuals

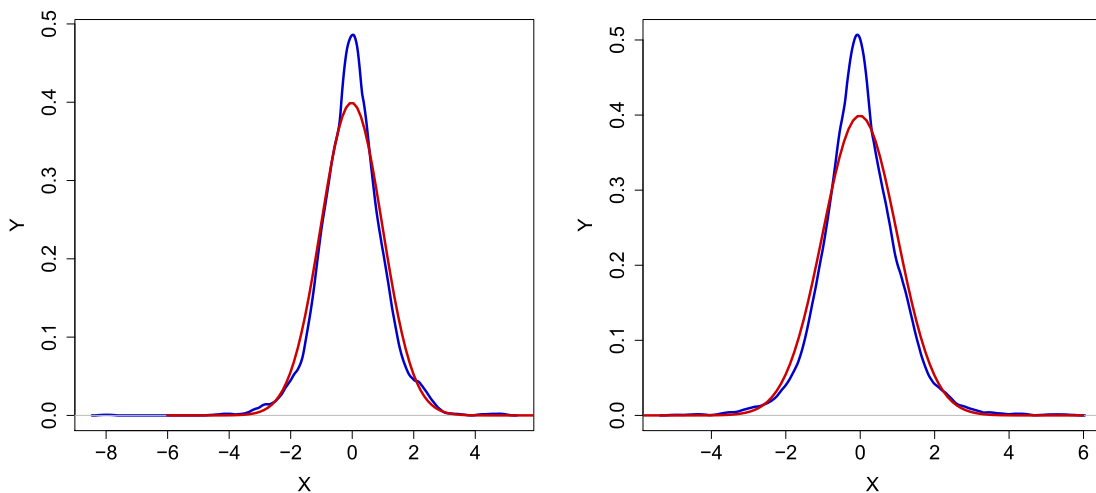


Fig. 1: Kernel density estimator of the residuals and of the normal density from TKA (left) and VOW (right). Quartic kernel, $\hat{h} = 2.78\hat{\sigma}n^{-0.2}$.

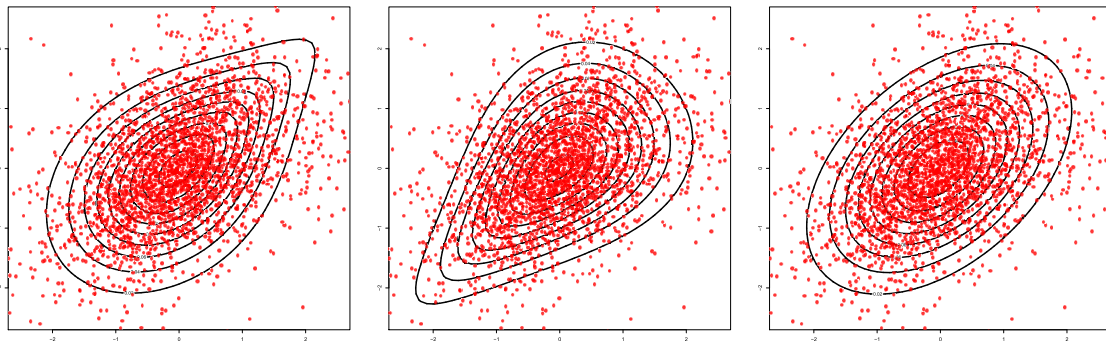


Fig. 2: Residuals $\hat{\varepsilon}$ and fitted copulae: Gaussian ($\hat{\rho} = 0.462$), Clayton ($\hat{\theta} = 0.880$), Gumbel ($\hat{\theta} = 1.439$).

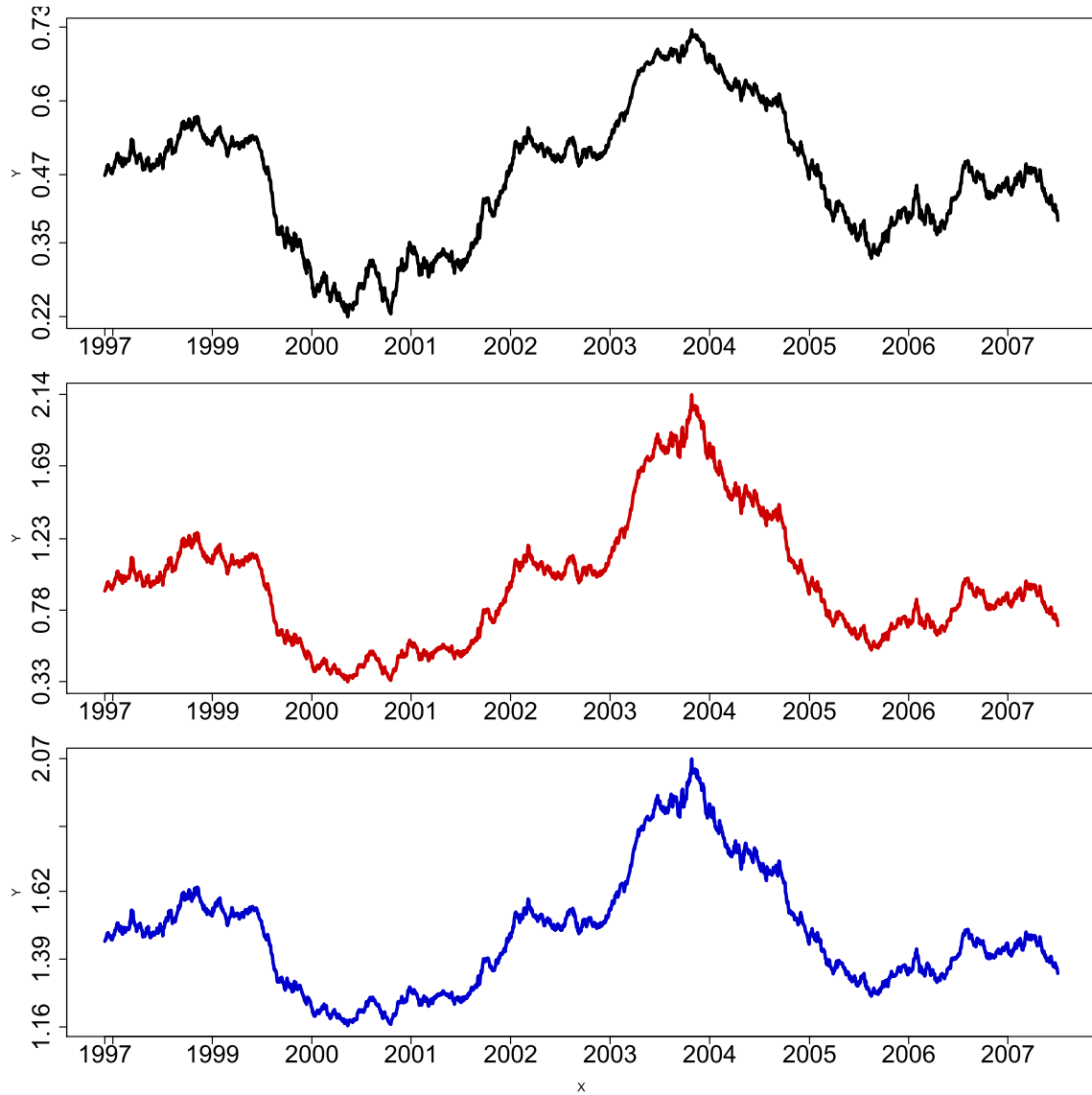


Fig. 3: Dependence parameter $\hat{\theta}$, estimated using the IFM method, Gaussian (upper panel), Gumbel (middle panel) and Clayton (lower panel) copulae, moving window ($w = 250$).

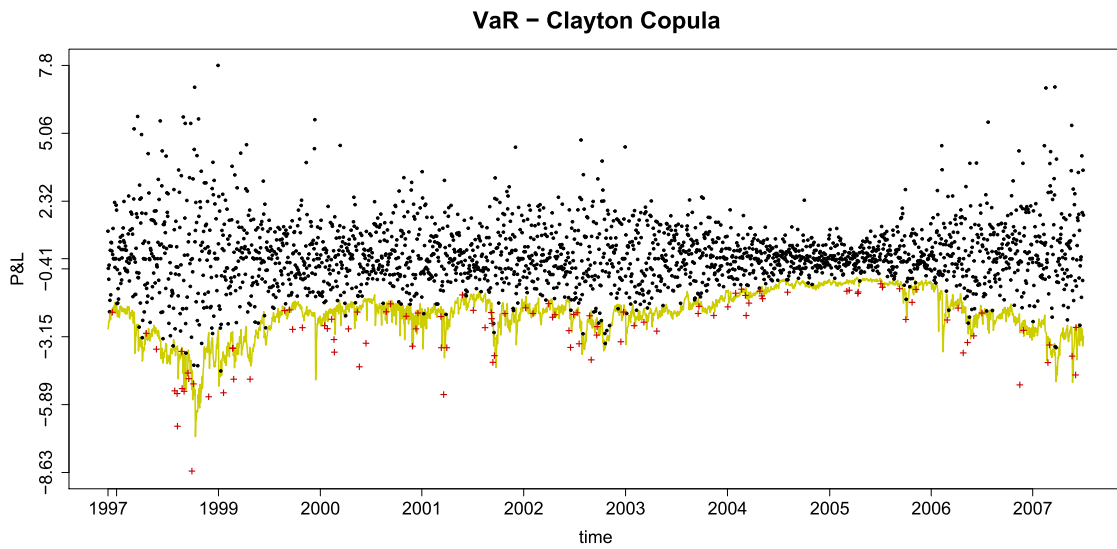


Fig. 4: $\widehat{VaR}(\alpha)$ (solid line), P&L (dots) and exceedances (crosses), $\alpha = 0.05$, $\hat{\alpha} = 0.0424$. P&L samples generated with Clayton copula.

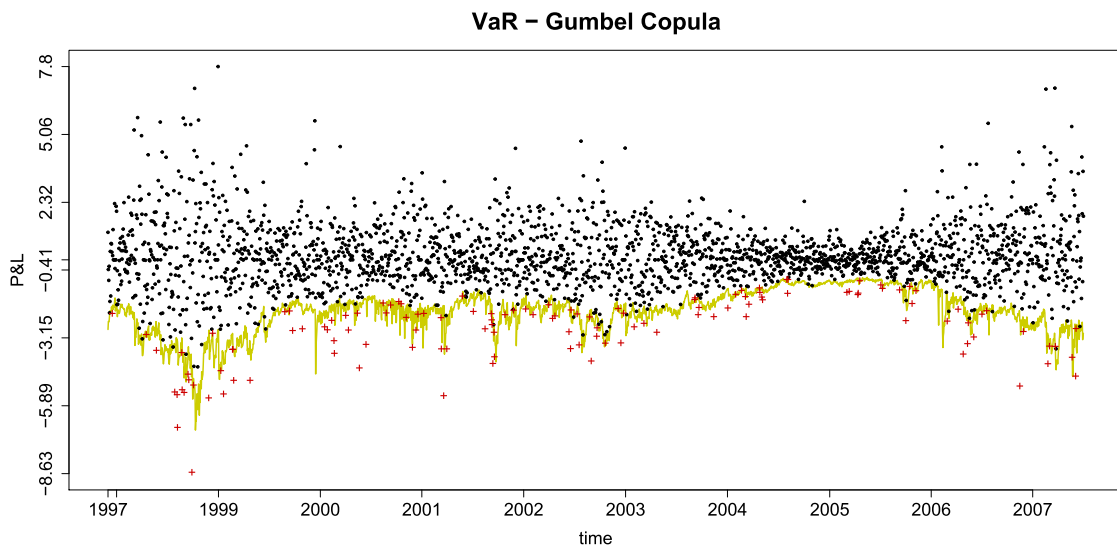


Fig. 5: $\widehat{VaR}(\alpha)$ (solid line), P&L (dots) and exceedances (crosses), $\alpha = 0.05$, $\hat{\alpha} = 0.0508$. P&L samples generated with Gumbel copula.

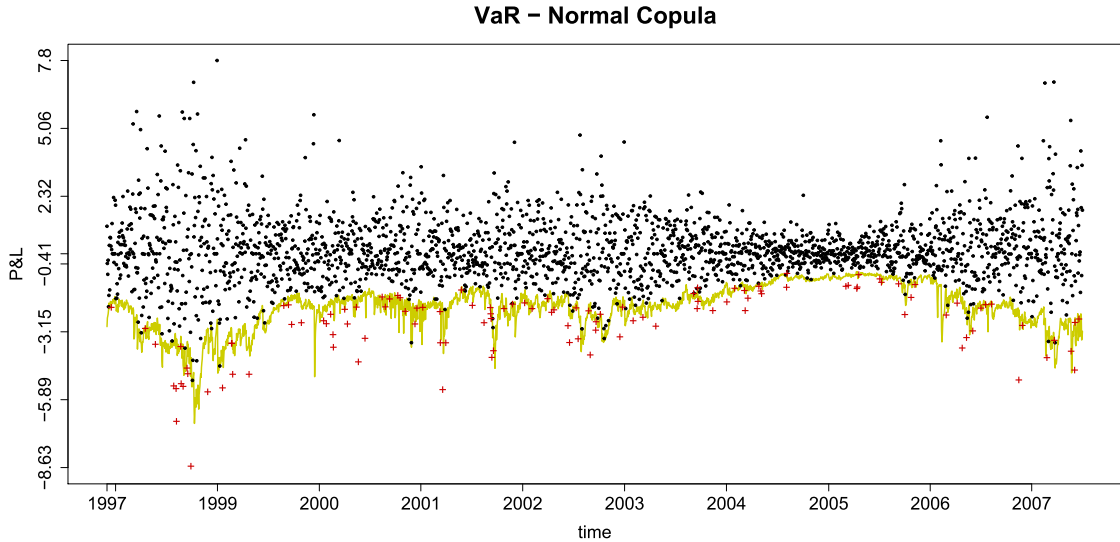


Fig. 6: $\widehat{VaR}(\alpha)$ (solid line), P&L (dots) and exceedances (crosses), $\alpha = 0.05$, $\hat{\alpha} = 0.0464$. P&L samples generated with Gaussian copula.

7.2 3-dimensional Portfolio

In this section, the Value-at-Risk of portfolios composed of 3 positions (Merck (MRK), Tyssenkrupp (TKA) and Volkswagen (VOW) from 01.12.1997 to 03.07.2007) is computed using a time-varying simple Gumbel copula and time-varying hierarchical Archimedean copula with generators from the Gumbel family.

The estimation of the parameters of the 3-dimensional copula was done by the IFM method. Concerning the HAC, we determine the structure under each window and re-estimate the parameters.

The fit of a GARCH(1,1) model to the sample of log returns $\{X_t\}_{t=1}^T$, $X_t = (X_{1,t}, X_{2,t}, X_{3,t})^\top$, $T = 2500$, gives the estimates $\hat{\omega}_j$, $\hat{\alpha}_j$ and $\hat{\beta}_j$, as in Table 2, and empirical residuals $\{\hat{\varepsilon}_t\}_{t=1}^T$, where $\hat{\varepsilon}_t = (\hat{\varepsilon}_{1,t}, \hat{\varepsilon}_{2,t}, \hat{\varepsilon}_{3,t})^\top$, as in upper right part of Figure 8. The marginal distributions are specified as normal, $\hat{\varepsilon}_j \sim N(\hat{\mu}_j, \hat{\sigma}_j)$ with the estimated parameters $\hat{\delta}_j = (\hat{\mu}_j, \hat{\sigma}_j)$.

The estimated Value-at-Risk at level α together with the P&L function are plotted in Figure 9 for the simple Archimedean Copula (AC) and on 10 for the HAC. As can be seen from the backtesting results for different VaR levels, HAC outperforms the simple AC in all levels. This implies the necessity of dependence flexibility in modelling of log-returns.

8 Summary

To conclude, a summary of the main findings of this paper. We calculated the Value-at-Risk for the static and dynamic portfolio constructed by different methods. Three different copulae - Gumbel, Clayton and Gaussian - were used to estimate the Value-at-Risk from the two- (MRK and TKA) and three- (MRK, TKA and VOW) dimensional portfolios. From the time series of estimated dependence parameters, we can verify that the dependence structure is represented in a similar form with all copula families, as in Figure 3.

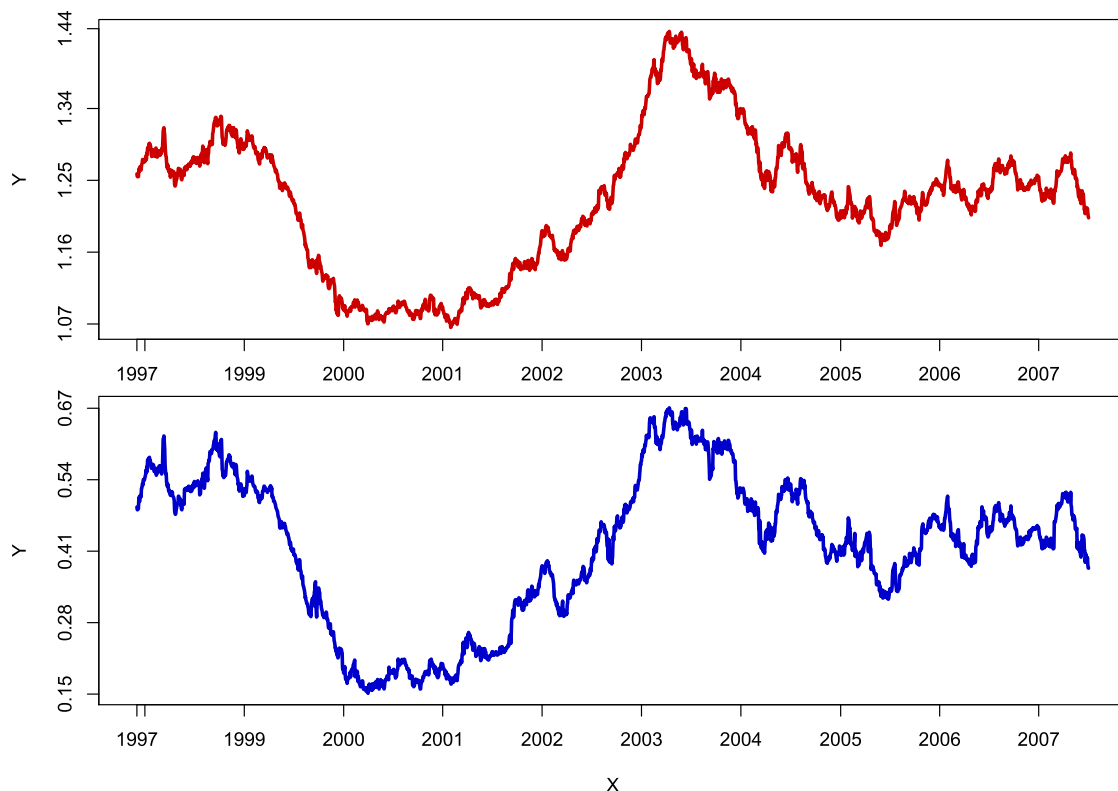


Fig. 7: Dependence parameter $\hat{\theta}$, estimated using the IFM method, Clayton (upper panel) and Gumbel (lower panel) copulae, moving window ($w = 250$).

Using backtesting results to compare the performance in the VaR estimation, we remark that on average the Clayton and Gaussian copulae *overestimate* the VaR. In terms of capital requirement, a financial institution computing VaR with those copulae would be requested to keep *more* capital aside than necessary to guarantee the desired confidence level.

The estimation with Gumbel copula, on another side, produced results close to the desired level. Gumbel copulae seems to represent specific data dependence structures (like lower tail dependencies, relevant to explain simultaneous losses) better than Gaussian and Clayton copulae.

References

- Clayton, D. G. (1978). A model for association in bivariate life tables and its application in epidemiological studies of familial tendency in chronic disease incidence, *Biometrika* **65**: 141–151.
- Deutsch, H. and Eller, R. (1999). *Derivatives and Internal Models*, Macmillan Press.
- Devroye, L. (1986). *Non-uniform Random Variate Generation*, Springer Verlag, New York.
- Embrechts, P., McNeil, A. J. and Straumann, D. (1999). Correlation and dependence in risk management: Properties and pitfalls, *RISK* pp. 69–71.

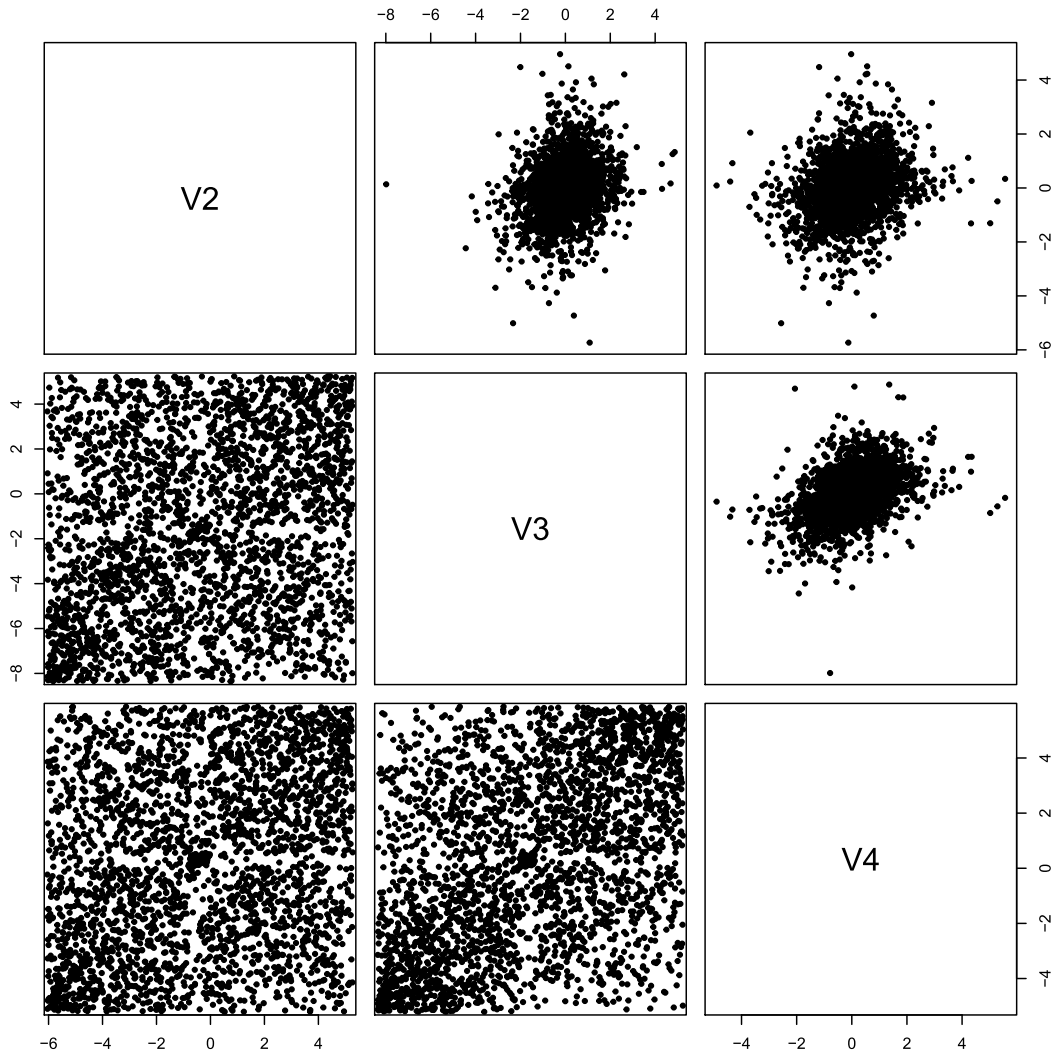


Fig. 8: Scatterplots from GARCH residulas (upper triangular) and from residuals mapped on unit square by the cdf (lower triangular).

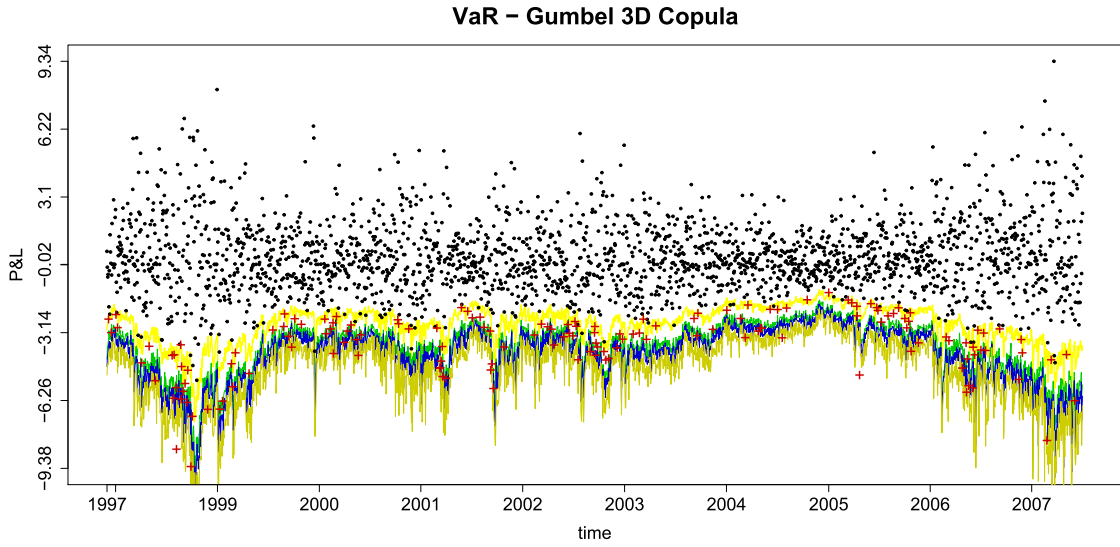


Fig. 9: $\widehat{VaR}(\alpha)$ and P&L (dots), estimated with 3-dimensional simple Gumbel copula, $\alpha_1 = 0.05$ ($\hat{\alpha}_1 = 0.0612$), $\alpha_2 = 0.01$ ($\hat{\alpha}_2 = 0.0232$), $\alpha_3 = 0.005$ ($\hat{\alpha}_3 = 0.016$) and $\alpha_4 = 0.001$ ($\hat{\alpha}_4 = 0.006$).

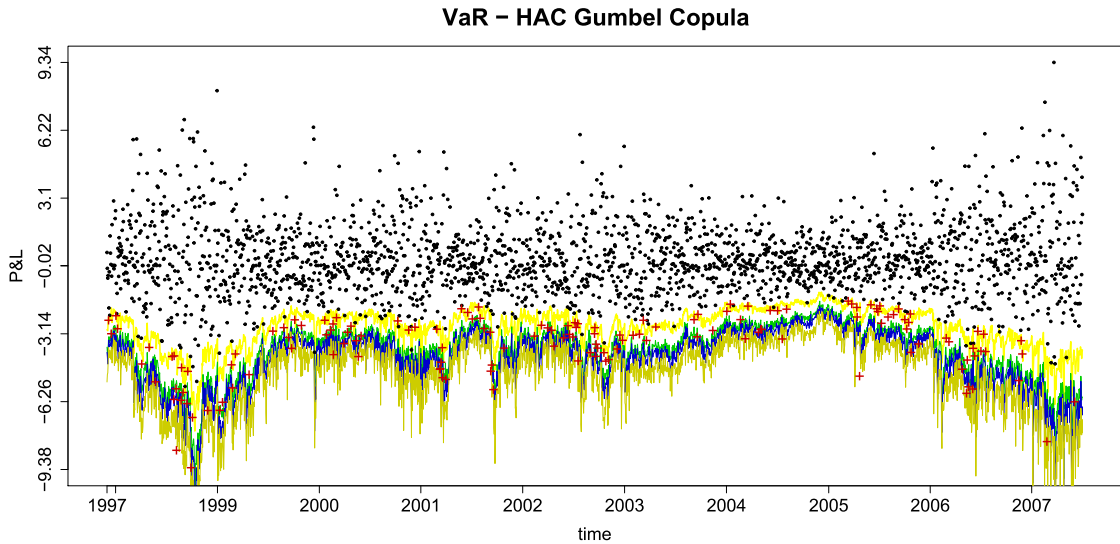


Fig. 10: $\widehat{VaR}(\alpha)$ and P&L (dots), estimated with 3-dimensional HAC with Gumbel generators, $\alpha_1 = 0.05$ ($\hat{\alpha}_1 = 0.0592$), $\alpha_2 = 0.01$ ($\hat{\alpha}_2 = 0.0208$), $\alpha_3 = 0.005$ ($\hat{\alpha}_3 = 0.014$) and $\alpha_4 = 0.001$ ($\hat{\alpha}_4 = 0.004$).

- Embrechts, P., McNeil, A. and Straumann, D. (1999b). Correlation: Pitfalls and alternatives, *RISK May*: 69–71.
- Frank, M. J. (1979). On the simultaneous associativity of $f(x, y)$ and $x + y - f(x, y)$, *Aequationes Mathematicae* **19**: 194–226.
- Frees, E. and Valdez, E. (1998). Understanding relationships using copulas, *North American Actuarial Journal* **2**: 1–125.
- Frey, R. and McNeil, A. J. (2003). Dependent defaults in models of portfolio credit risk, *Journal of Risk* **6**(1): 59–92.
- Genest, C. and Rivest, L.-P. (1989). A characterization of Gumbel family of extreme value distributions, *Statistics and Probability Letters* **8**: 207–211.
- Giacomini, E. and Härdle, W. (2005). Value-at-risk calculations with time varying copulae, *Proceedings 55th International Statistical Institute, Sydney 2005* .
- Gumbel, E. J. (1960). Distributions des valeurs extrêmes en plusieurs dimensions, *Publ. Inst. Statist. Univ. Paris* **9**: 171–173.
- Hoeffding, W. (1940). Masstabinvariante Korrelationstheorie, *Schriften des Mathematischen Instituts und des Instituts für Angewandte Mathematik der Universität Berlin* **5**(3): 179–233.
- Hoeffding, W. (1941). Masstabinvariante Korrelationsmasse für diskontinuierliche Verteilungen, *Archiv für die mathematische Wirtschafts- und Sozialforschung* **7**: 49–70.
- Joe, H. (1997). *Multivariate Models and Dependence Concepts*, Chapman and Hall, London.
- Joe, H. and Xu, J. J. (1996). The estimation method of inference functions for margins for multivariate models, *Technical Report 166*, Department of Statistics, University of British Columbia.
- Marshall, A. W. and Olkin, J. (1988). Families of multivariate distributions, *Journal of the American Statistical Association* **83**: 834–841.
- McNeil, A. J. (2008). Sampling nested Archimedean copulas, *Journal Statistical Computation and Simulation* . forthcoming.
- Nelsen, R. B. (2006). *An Introduction to Copulas*, Springer Verlag, New York.
- Okhrin, O., Okhrin, Y. and Schmid, W. (2009a). On the structure and estimation of hierarchical Archimedean copulas. under revision in *Journal of Econometrics*.
- Okhrin, O., Okhrin, Y. and Schmid, W. (2009b). Properties of Hierarchical Archimedean Copulas, *SFB 649 Discussion Paper 2009-014*, Sonderforschungsbereich 649, Humboldt-Universität zu Berlin, Germany. available at <http://sfb649.wiwi.hu-berlin.de/papers/pdf/SFB649DP2009-014.pdf>.
- Patton, A. J. (2004). On the out-of-sample importance of skewness and asymmetric dependence for asset allocation, *Journal of Financial Econometrics* **2**: 130–168.

- Savu, C. and Tiede, M. (2006). Hierarchical Archimedean copulas, *Discussion paper*, University of Muenster.
- Sklar, A. (1959). Fonctions de répartition à n dimensions et leurs marges, **8**: 229–231.
- Whelan, N. (2004). Sampling from Archimedean copulas, *Quantitative Finance* **4**: 339–352.

CONFIDENCE BANDS IN QUANTILE REGRESSION

WOLFGANG K. HÄRDLE AND SONG SONG
Humboldt-Universität zu Berlin

Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be independent and identically distributed random variables and let $l(x)$ be the unknown p -quantile regression curve of Y conditional on X . A quantile smoother $l_n(x)$ is a localized, nonlinear estimator of $l(x)$. The strong uniform consistency rate is established under general conditions. In many applications it is necessary to know the stochastic fluctuation of the process $\{l_n(x) - l(x)\}$. Using strong approximations of the empirical process and extreme value theory, we consider the asymptotic maximal deviation $\sup_{0 \leq x \leq 1} |l_n(x) - l(x)|$. The derived result helps in the construction of a uniform confidence band for the quantile curve $l(x)$. This confidence band can be applied as a econometric model check. An economic application considers the relation between age and earnings in the labor market by means of parametric model specification tests, which presents a new framework to describe trends in the entire wage distribution in a parsimonious way.

1. INTRODUCTION

In standard regression function estimation, most investigations are concerned with the conditional mean regression. However, new insights about the underlying structures can be gained by considering other aspects of the conditional distribution. The quantile curves are key aspects of inference in various economic problems and are of great interest in practice. These describe the conditional behavior of a response variable (e.g., wage of workers) given the value of an explanatory variable (e.g., education level, experience, occupation of workers) and investigate changes in both tails of the distribution, other than just the mean.

When examining labor markets, economists are concerned with whether discrimination exists, e.g., for different genders, nationalities, union status, etc. To study this question, we need to separate out other effects first, e.g., age, education, etc. The crucial relation between age and earnings or salaries belongs to the most carefully studied subjects in labor economics. The fundamental work in mean regression can be found in Murphy and Welch (1990). Quantile regression estimates could provide more accurate measures. Koenker and Hallock (2001) present a group of important economic applications, including quantile

Financial support from the Deutsche Forschungsgemeinschaft via SFB 649 "Ökonomisches Risiko," Humboldt-Universität zu Berlin, is gratefully acknowledged. We thank the editor and two referees for concrete suggestions on improving the manuscript and restructuring the paper. Their valuable comments and suggestions are gratefully acknowledged. Address correspondence to Song Song, Institute for Statistics and Econometrics, Humboldt-Universität zu Berlin, Spandauer Straße 1, 10178 Berlin, Germany; e-mail: songsong@cms.hu-berlin.de.

Engel curves, and claim that “quantile regression is gradually developing into a comprehensive strategy for completing the regression prediction.” Besides this, it is also well known that a quantile regression model (e.g., the conditional median curve) is more robust to outliers, especially for fat-tailed distributions. For symmetric conditional distributions the quantile regression generates the nonparametric mean regression analysis because the $p = 0.5$ (median) quantile curve coincides with the mean regression.

As first introduced by Koenker and Bassett (1978), one may assume a parametric model for the p -quantile curve and estimate parameters by the interior point method discussed by Koenker and Park (1996) and Portnoy and Koenker (1997). Similarly, we can also adopt nonparametric methods to estimate conditional quantiles. The first one, a more direct approach using a check function such as a robustified local linear smoother, is provided by Fan, Hu, and Troung (1994) and further extended by Yu and Jones (1997, 1998). An alternative procedure is first to estimate the conditional distribution function using the double-kernel local linear technique of Fan, Yao, and Tong (1996) and then to invert the conditional distribution estimator to produce an estimator of a conditional quantile by Yu and Jones (1997, 1998). Beside these, Hall, Wolff, and Yao (1999) proposed a weighted version of the Nadaraya–Watson estimator, which was further studied by Cai (2002). Recently Jeong and Härdle (2008) have developed the conditional quantile causality test. More generally, for an M -regression function that involves quantile regression as a special case, the uniform Bahadur representation and application to the additive model are studied by Kong, Linton, and Xia (2010). An interesting question for parametric fitting, especially from labor economists, would be how well these models fit the data, when compared with the nonparametric estimation method.

Let $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ be a sequence of independent and identically distributed (i.i.d.) bivariate random variables with joint probability density function (pdf) $f(x, y)$, joint cumulative distribution function (cdf) $F(x, y)$, conditional pdf $f(y|x), f(x|y)$, conditional cdf $F(y|x), F(x|y)$ for Y given X and X given Y , respectively, and marginal pdf $f_X(x)$ for X , $f_Y(y)$ for Y where $x \in J$ and J is a possibly infinite interval in \mathbb{R}^d and $y \in \mathbb{R}$. In general, X may be a multivariate covariate, although here we restrict attention to the univariate case and $J = [0, 1]$ for convenience. Let $l(x)$ denote the p -quantile curve, i.e., $l(x) = F_{Y|x}^{-1}(p)$.

Under a “check function,” the quantile regression curve $l(x)$ can be viewed as the minimizer of $L(\theta) \stackrel{\text{def}}{=} E\{\rho_p(y - \theta)|X = x\}$ (with respect to θ) with $\rho_p(u) = pu\mathbf{1}\{u \in (0, \infty)\} - (1 - p)u\mathbf{1}\{u \in (-\infty, 0)\}$, which was originally motivated by an exercise in Ferguson (1967, p. 51) in the literature.

A kernel-based p -quantile curve estimator $l_n(x)$ can naturally be constructed by minimizing:

$$L_n(\theta) = n^{-1} \sum_{i=1}^n \rho_p(Y_i - \theta) K_h(x - X_i) \tag{1}$$

with respect to $\theta \in I$ where I is a possibly infinite, or possibly degenerate, interval in \mathbb{R} and $K_h(u) = h^{-1}K(u/h)$ is a kernel with bandwidth h . The numerical solution of (1) may be found iteratively as in Lejeune and Sarda (1988) and Yu, Lu, and Stander (2003).

In light of the concepts of M -estimation as in Huber (1981), if we define $\psi(u)$ as

$$\begin{aligned} \psi_p(u) &= p\mathbf{1}\{u \in (0, \infty)\} - (1 - p)\mathbf{1}\{u \in (-\infty, 0)\} \\ &= p - \mathbf{1}\{u \in (-\infty, 0)\}, \end{aligned}$$

$l_n(x)$ and $l(x)$ can be treated as a zero (with respect to θ) of the function

$$\tilde{H}_n(\theta, x) \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n K_h(x - X_i)\psi(Y_i - \theta), \tag{2}$$

$$\tilde{H}(\theta, x) \stackrel{\text{def}}{=} \int_{\mathbb{R}} f(x, y)\psi(y - \theta) dy, \tag{3}$$

correspondingly.

To show the uniform consistency of the quantile smoother, we shall reduce the problem of strong convergence of $l_n(x) - l(x)$, uniformly in x , to an application of the strong convergence of $\tilde{H}_n(\theta, x)$ to $\tilde{H}(\theta, x)$, uniformly in x and θ , as given by Theorem 2.2 in Härdle, Janssen, and Serfling (1988). It is shown that under general conditions almost surely (a.s.)

$$\sup_{x \in J} |l_n(x) - l(x)| \leq B^* \max \left\{ (nh/(\log n))^{-1/2}, h^{\tilde{\alpha}} \right\}, \quad \text{as } n \rightarrow \infty,$$

where B^* and $\tilde{\alpha}$ are parameters defined more precisely in Section 2.

Note that without assuming K has compact support (as we do here) under similar assumptions Franke and Mwita (2003) obtain

$$\begin{aligned} l_n(x) &= \hat{F}_{Y|x}^{-1}(p), \\ \hat{F}(y|x) &= \frac{\sum_{i=1}^n K_h(x - X_i)\mathbf{1}(Y_i < y)}{\sum_{i=1}^n K_h(x - X_i)}, \end{aligned}$$

$$\sup_{x \in J} |l_n(x) - l(x)| \leq B^{**} \left\{ (nh/(s_n \log n))^{-1/2} + h^2 \right\}, \quad \text{as } n \rightarrow \infty$$

for α -mixing data where B^{**} is some constant and $s_n, n \geq 1$ is an increasing sequence of positive integers satisfying $1 \leq s_n \leq n/2$ and some other criteria. Thus $\{nh/(\log n)\}^{-1/2} \leq \{nh/(s_n \log n)\}^{-1/2}$.

By employing similar methods to those developed in Härdle (1989) it is shown in this paper that

$$\begin{aligned} &P \left((2\delta \log n)^{1/2} \left[\sup_{x \in J} r(x) |l_n(x) - l(x)| / \lambda(K)^{1/2} - d_n \right] < z \right) \\ &\rightarrow \exp\{-2 \exp(-z)\}, \quad \text{as } n \rightarrow \infty \end{aligned} \tag{4}$$

from the asymptotic Gumbel distribution where $r(x)$, δ , $\lambda(K)$, d_n are suitable scaling parameters. The asymptotic result (4) therefore allows the construction of (asymptotic) uniform confidence bands for $l(x)$ based on specifications of the stochastic fluctuation of $l_n(x)$. The strong approximation with Brownian bridge techniques that we use in this paper is available only for the approximation of the two-dimensional empirical process. The extension to the multivariate covariable can be done by partial linear modeling, which deserves further research.

The plan of the paper is as follows. In Section 2, the stochastic fluctuation of the process $\{l_n(x) - l(x)\}$ and the uniform confidence band are presented through the equivalence of several stochastic processes, with a strong uniform consistency rate of $\{l_n(x) - l(x)\}$ also shown. In Section 3, in a small Monte Carlo study we investigate the behavior of $l_n(x)$ when the data are generated by fat-tailed conditional distributions of $(Y|X = x)$. In Section 4, an application considers a wage-earning relation in the labor market. All proofs are sketched in the Appendix.

2. RESULTS

The following assumptions will be convenient. To make x and X clearly distinguishable, we replace x by t sometimes, but they are essentially the same.

(A1) The kernel $K(\cdot)$ is positive and symmetric, has compact support $[-A, A]$, and is Lipschitz continuously differentiable with bounded derivatives.

(A2) $(nh)^{-1/2}(\log n)^{3/2} \rightarrow 0$, $(n \log n)^{1/2} h^{5/2} \rightarrow 0$, $(nh^3)^{-1}(\log n)^2 \leq M$, where M is a constant.

(A3) $h^{-3}(\log n) \int_{|y|>a_n} f_Y(y)dy = \mathcal{O}(1)$, where $f_Y(y)$ is the marginal density of Y and $\{a_n\}_{n=1}^\infty$ is a sequence of constants tending to infinity as $n \rightarrow \infty$.

(A4) $\inf_{t \in J} |q(t)| \geq q_0 > 0$, where $q(t) = \partial E\{\psi(Y - \theta)|t\} / \partial \theta|_{\theta=l(t)} \cdot f_X(t) = f\{l(t)|t\} f_X(t)$.

(A5) The quantile function $l(t)$ is Lipschitz twice continuously differentiable for all $t \in J$.

(A6) $0 < m_1 \leq f_X(t) \leq M_1 < \infty$, $t \in J$; the conditional densities $f(\cdot|y)$, $y \in \mathbb{R}$, are uniform local Lipschitz continuous of order $\tilde{\alpha}$ (uLL- $\tilde{\alpha}$) on J , uniformly in $y \in \mathbb{R}$, with $0 < \tilde{\alpha} \leq 1$.

Define also

$$\sigma^2(t) = E[\psi^2\{Y - l(t)\}|t] = p(1 - p),$$

$$H_n(t) = (nh)^{-1} \sum_{i=1}^n K\{(t - X_i)/h\} \psi\{Y_i - l(t)\},$$

$$D_n(t) = \partial(nh)^{-1} \sum_{i=1}^n K\{(t - X_i)/h\} \psi\{Y_i - \theta\} / \partial \theta|_{\theta=l(t)}$$

and assume that $\sigma^2(t)$ and $f_X(t)$ are differentiable.

Assumption (A1) on the compact support of the kernel could possibly be relaxed by introducing a cutoff technique as in Csörgö and Hall (1982) for density estimators. Assumption (A2) has purely technical reasons: to keep the bias at a lower rate than the variance and to ensure the vanishing of some nonlinear remainder terms. Assumption (A3) appears in a somewhat modified form also in Johnston (1982). Assumptions (A5) and (A6) are common assumptions in robust estimation as in Huber (1981) and Härdle et al. (1988) that are satisfied by exponential and generalized hyperbolic distributions.

For the uniform strong consistency rate of $l_n(x) - l(x)$, we apply the result of Härdle et al. (1988) by taking $\beta(y) = \psi(y - \theta)$, $y \in \mathbb{R}$, for $\theta \in I = \mathbb{R}$, $q_1 = q_2 = -1$, $\gamma_1(y) = \max\{0, -\psi(y - \theta)\}$, $\gamma_2(y) = \min\{0, -\psi(y - \theta)\}$, and $\lambda = \infty$ to satisfy the representations for the parameters there. Thus from Härdle et al.'s Theorem 2.2 and Remark 2.3(v), we immediately have the following lemma.

LEMMA 2.1. *Let $\tilde{H}_n(\theta, x)$ and $\tilde{H}(\theta, x)$ be given by (2) and (3). Under Assumption (A6) and $(nh/\log n)^{-1/2} \rightarrow \infty$ through Assumption (A2), for some constant A^* not depending on n , we have a.s. as $n \rightarrow \infty$*

$$\sup_{\theta \in I} \sup_{x \in J} |\tilde{H}_n(\theta, x) - \tilde{H}(\theta, x)| \leq A^* \max \left\{ (nh/\log n)^{-1/2}, h^{\tilde{\alpha}} \right\}. \tag{5}$$

For our result on $l_n(\cdot)$, we shall also require

$$\inf_{x \in J} \left| \int \psi \{y - l(x) + \varepsilon\} dF(y|x) \right| \geq \tilde{q}|\varepsilon|, \quad \text{for } |\varepsilon| \leq \delta_1, \tag{6}$$

where δ_1 and \tilde{q} are some positive constants; see also Härdle and Luckhaus (1984). This assumption is satisfied if there exists a constant \tilde{q} such that $f(l(x)|x) > \tilde{q}/p$, $x \in J$.

THEOREM 2.1. *Under the conditions of Lemma 2.1 and also assuming (6), we have a.s. as $n \rightarrow \infty$*

$$\sup_{x \in J} |l_n(x) - l(x)| \leq B^* \max \left\{ (nh/\log n)^{-1/2}, h^{\tilde{\alpha}} \right\} \tag{7}$$

with $B^* = A^*/m_1\tilde{q}$ not depending on n and m_1 a lower bound of $f_X(t)$. If additionally $\tilde{\alpha} \geq \{\log(\sqrt{\log n}) - \log(\sqrt{nh})\}/\log h$, it can be further simplified to

$$\sup_{x \in J} |l_n(x) - l(x)| \leq B^* \{ (nh/\log n)^{-1/2} \}.$$

THEOREM 2.2. *Let $h = n^{-\delta}$, $\frac{1}{5} < \delta < \frac{1}{3}$, $\lambda(K) = \int_{-A}^A K^2(u) du$, and*

$$d_n = (2\delta \log n)^{1/2} + (2\delta \log n)^{-1/2} \left[\log \left\{ c_1(K)/\pi^{1/2} \right\} + \frac{1}{2} \{ \log \delta + \log \log n \} \right],$$

if $c_1(K) = \{K^2(A) + K^2(-A)\}/\{2\lambda(K)\} > 0$;

$$d_n = (2\delta \log n)^{1/2} + (2\delta \log n)^{-1/2} \log\{c_2(K)/2\pi\}$$

otherwise with $c_2(K) = \int_{-A}^A \{K'(u)\}^2 du / \{2\lambda(K)\}$. Then (4) holds with

$$r(x) = (nh)^{1/2} f\{l(x)|x\} \{f_X(x)/p(1-p)\}^{1/2}.$$

This theorem can be used to construct uniform confidence intervals for the regression function as stated in the following corollary.

COROLLARY 2.1. *Under the assumptions of Theorem 2.2, an approximate $(1-\alpha) \times 100\%$ confidence band over $[0, 1]$ is*

$$l_n(t) \pm (nh)^{-1/2} \left\{ p(1-p)\lambda(K)/\hat{f}_X(t) \right\}^{1/2} \hat{f}^{-1}\{l(t)|t\} \left\{ d_n + c(\alpha)(2\delta \log n)^{-1/2} \right\},$$

where $c(\alpha) = \log 2 - \log |\log(1-\alpha)|$ and $\hat{f}_X(t)$, $\hat{f}\{l(t)|t\}$ are consistent estimates for $f_X(t)$, $f\{l(t)|t\}$.

In the literature, according to Fan et al. (1994, 1996), Yu and Jones (1997, 1998), Hall et al. (1999), Cai (2002), and others, asymptotic normality at interior points for various nonparametric smoothers, e.g., local constant, local linear, reweighted Nadaraya–Watson methods, etc., has been shown:

$$\sqrt{nh}\{l_n(t) - l(t)\} \sim N(0, \tau^2(t))$$

with $\tau^2(t) = \lambda(K)p(1-p)/[f_X(t)f^2\{l(t)|t\}]$. Note that the bias term vanishes here as we adjust h . With $\tau(t)$ introduced, we can further write Corollary 2.1 as

$$l_n(t) \pm (nh)^{-1/2} \left\{ d_n + c(\alpha)(2\delta \log n)^{-1/2} \right\} \hat{\tau}(t).$$

Through minimizing the approximation of asymptotic mean square error, the optimal bandwidth h_p can be computed. In practice, the rule of thumb for h_p is given by Yu and Jones (1998):

1. Use ready-made and sophisticated methods to select optimal bandwidth h_{mean} from conditional mean regression, e.g., Ruppert, Sheather, and Wand (1995);
2. $h_p = [p(1-p)/\varphi^2\{\Phi^{-1}(p)\}]^{1/5} \cdot h_{\text{mean}}$ with φ , Φ as the pdf and cdf of a standard normal distribution

Obviously the further p lies from 0.5, the more smoothing is necessary.

The proof is essentially based on a linearization argument after a Taylor series expansion. The leading linear term will then be approximated in a similar way as in Johnston (1982) and Bickel and Rosenblatt (1973). The main idea behind the proof is a strong approximation of the empirical process of $\{(X_i, Y_i)_{i=1}^n\}$ by a sequence of Brownian bridges as proved by Tusnady (1977).

As $l_n(t)$ is the zero (with respect to θ) of $\tilde{H}_n(\theta, t)$, it follows by applying second-order Taylor expansions to $\tilde{H}_n(\theta, t)$ around $l(t)$ that

$$l_n(t) - l(t) = \{H_n(t) - E H_n(t)\}/q(t) + R_n(t), \tag{8}$$

where $\{H_n(t) - E H_n(t)\}/q(t)$ is the leading linear term and

$$R_n(t) = H_n(t)\{q(t) - D_n(t)\}/\{D_n(t) \cdot q(t)\} + E H_n(t)/q(t) + \frac{1}{2}\{l_n(t) - l(t)\}^2 \cdot \{D_n(t)\}^{-1} \tag{9}$$

$$\cdot (nh)^{-1} \sum_{i=1}^n K\{(x - X_i)/h\} \psi''\{Y_i - l(t) + r_n(t)\}, \tag{10}$$

$$|r_n(t)| < |l_n(t) - l(t)|$$

is the remainder term. In the Appendix it is shown (Lemma A.1) that $\|R_n\| = \sup_{t \in J} |R_n(t)| = \mathcal{O}_p\{(nh \log n)^{-1/2}\}$.

Furthermore, the rescaled linear part

$$Y_n(t) = (nh)^{1/2} \{\sigma^2(t) f_X(t)\}^{-1/2} \{H_n(t) - E H_n(t)\}$$

is approximated by a sequence of Gaussian processes, leading finally to the Gaussian process

$$Y_{5,n}(t) = h^{-1/2} \int K\{(t - x)/h\} dW(x). \tag{11}$$

Drawing upon the result of Bickel and Rosenblatt (1973), we finally obtain asymptotically the Gumbel distribution.

We also need the Rosenblatt (1952) transformation,

$$T(x, y) = \{F_{X|y}(x|y), F_Y(y)\},$$

which transforms (X_i, Y_i) into $T(X_i, Y_i) = (X'_i, Y'_i)$ mutually independent uniform random variables. In the event that x is a d -dimensional covariate, the transformation becomes

$$T(x_1, x_2, \dots, x_d, y) = \{F_{X_1|y}(x_1|y), F_{X_2|y}(x_2|x_1, y), \dots, F_{X_k|x_{d-1}, \dots, x_1, y}(x_k|x_{d-1}, \dots, x_1, y), F_Y(y)\}. \tag{12}$$

With the aid of this transformation, Theorem 1 of Tusnady (1977) may be applied to obtain the following lemma.

LEMMA 2.2. *On a suitable probability space a sequence of Brownian bridges B_n exists such that*

$$\sup_{x \in J, y \in \mathbb{R}} |Z_n(x, y) - B_n\{T(x, y)\}| = \mathcal{O}\left\{n^{-1/2}(\log n)^2\right\} \quad a.s.,$$

where $Z_n(x, y) = n^{1/2}\{F_n(x, y) - F(x, y)\}$ denotes the empirical process of $\{(X_i, Y_i)\}_{i=1}^n$.

For $d > 2$, it is still an open problem that deserves further research.

Before we define the different approximating processes, let us first rewrite (11) as a stochastic integral with respect to the empirical process $Z_n(x, y)$:

$$Y_n(t) = \{hg'(t)\}^{-1/2} \iint K\{(t-x)/h\}\psi\{y-l(t)\}dZ_n(x, y),$$

$$g'(t) = \sigma^2(t)f_X(t).$$

The approximating processes are now

$$Y_{0,n}(t) = \{hg(t)\}^{-1/2} \iint_{\Gamma_n} K\{(t-x)/h\}\psi\{y-l(t)\}dZ_n(x, y), \tag{13}$$

where $\Gamma_n = \{|y| \leq a_n\}$, $g(t) = E[\psi^2\{y-l(t)\} \cdot \mathbf{1}(|y| \leq a_n) | X = t] \cdot f_X(t)$

$$Y_{1,n}(t) = \{hg(t)\}^{-1/2} \iint_{\Gamma_n} K\{(t-x)/h\}\psi\{y-l(t)\}dB_n\{T(x, y)\}, \tag{14}$$

$\{B_n\}$ being the sequence of Brownian bridges from Lemma 2.2.

$$Y_{2,n}(t) = \{hg(t)\}^{-1/2} \iint_{\Gamma_n} K\{(t-x)/h\}\psi\{y-l(t)\}dW_n\{T(x, y)\}, \tag{15}$$

$\{W_n\}$ being the sequence of Wiener processes satisfying

$$B_n(x', y') = W_n(x', y') - x'y'W_n(1, 1),$$

$$Y_{3,n}(t) = \{hg(t)\}^{-1/2} \iint_{\Gamma_n} K\{(t-x)/h\}\psi\{y-l(x)\}dW_n\{T(x, y)\}, \tag{16}$$

$$Y_{4,n}(t) = \{hg(t)\}^{-1/2} \int g(x)^{1/2}K\{(t-x)/h\}dW(x), \tag{17}$$

$$Y_{5,n}(t) = h^{-1/2} \int K\{(t-x)/h\}dW(x), \tag{18}$$

$\{W(\cdot)\}$ being the Wiener process.

Lemmas A.2–A.7 in the Appendix ensure that all these processes have the same limit distributions. The result then follows from the next lemma.

LEMMA 2.3 (Theorem 3.1 in Bickel and Rosenblatt, 1973). *Let $d_n, \lambda(K), \delta$ as in Theorem 2.2. Let*

$$Y_{5,n}(t) = h^{-1/2} \int K\{(t-x)/h\}dW(x).$$

Then, as $n \rightarrow \infty$, the supremum of $Y_{5,n}(t)$ has a Gumbel distribution:

$$P \left\{ (2\delta \log n)^{1/2} \left[\sup_{t \in J} |Y_{5,n}(t)| / \{\lambda(K)\}^{1/2} - d_n \right] < z \right\} \rightarrow \exp\{-2 \exp(-z)\}.$$

3. A MONTE CARLO STUDY

We generate bivariate data $\{(X_i, Y_i)\}_{i=1}^n, n = 500$ with joint pdf:

$$f(x, y) = g \left(y - \sqrt{x + 2.5} \right) \mathbf{1}(x \in [-2.5, 2.5]), \tag{19}$$

$$g(u) = \frac{9}{10} \varphi(u) + \frac{1}{90} \varphi(u/9).$$

The p -quantile curve $l(x)$ can be obtained from a zero (with respect to θ) of

$$9\Phi(\theta) + \Phi(\theta/9) = 10p,$$

with Φ as the cdf of a standard normal distribution. Solving it numerically gives the 0.5-quantile curve $l(x) = \sqrt{x + 2.5}$ and the 0.9-quantile curve $l(x) = 1.5296 + \sqrt{x + 2.5}$. We use the quartic kernel:

$$K(u) = \frac{15}{16} (1 - u^2)^2, \quad |u| \leq 1, \\ = 0, \quad |u| > 1.$$

In Figure 1 the raw data, together with the 0.5-quantile curve, are displayed. The random variables generated with probability $\frac{1}{10}$ from the fat-tailed pdf $\frac{1}{9} \varphi(u/9)$ (see eqn. (19)) are marked as squares whereas the standard normal random variables are shown as stars. We then compute both the Nadaraya–Watson estimator $m_n^*(x)$ and the 0.5-quantile smoother $l_n(x)$. The bandwidth is set to 1.25, which is equivalent to 0.25 after rescaling x to $[0, 1]$ and fulfills the requirements of Theorem 2.2.

In Figure 1 $l(x)$, $m_n^*(x)$, and $l_n(x)$ are shown as a dotted line, dashed-dot line, and solid line, respectively. At first sight $m_n^*(x)$ has clearly more variation and has the expected sensitivity to the fat tails of $f(x, y)$. A closer look reveals that $m_n^*(x)$ for $x \approx 0$ apparently even leaves the 0.5-quantile curve. It may be surprising that this happens at $x \approx 0$ where no outlier is placed, but a closer look at Figure 1 shows that the large negative data values at both $x \approx -0.1$ and $x \approx 0.25$ cause the problem. This data value is inside the window ($h = 1.10$) and therefore distorts $m_n^*(x)$ for $x \approx 0$. The quantile smoother $l_n(x)$ (solid line) is unaffected and stays fairly close to the 0.5-quantile curve. Similar results can be obtained in Figure 2 corresponding to the 0.9 quantile ($h = 1.25$) with the 95% confidence band.

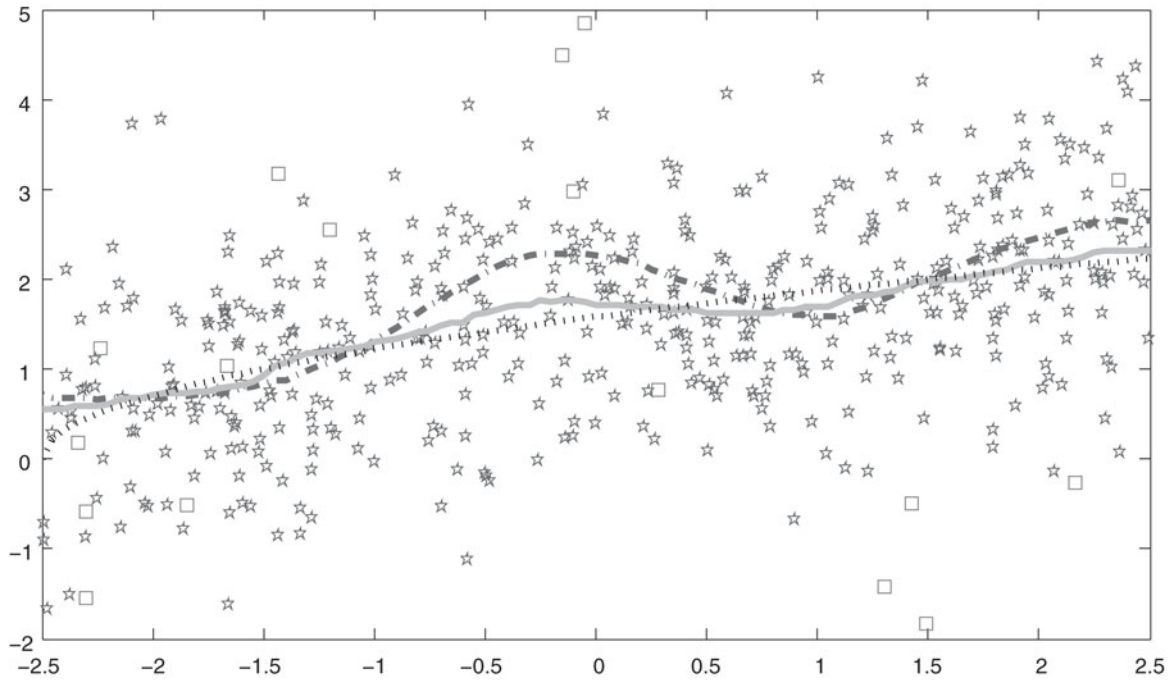


FIGURE 1. The 0.5-quantile curve, the Nadaraya–Watson estimator $m_n^*(x)$, and the 0.5-quantile smoother $l_n(x)$.

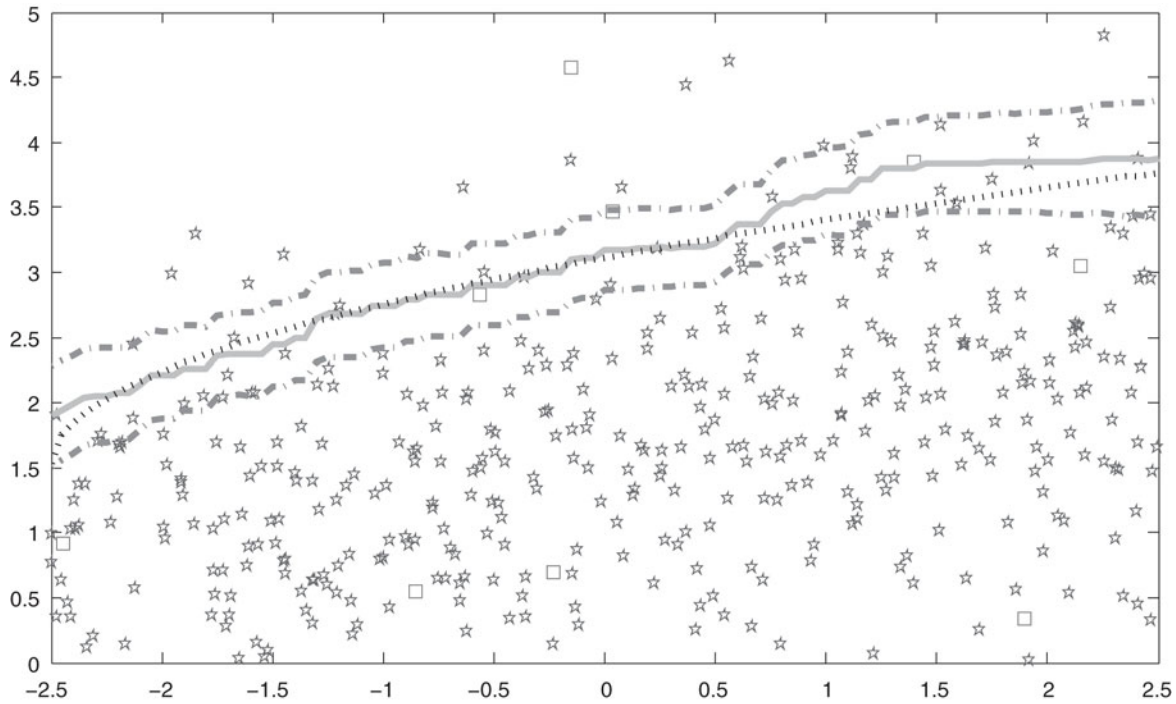


FIGURE 2. The 0.9-quantile curve, the 0.9-quantile smoother, and 95% confidence band.

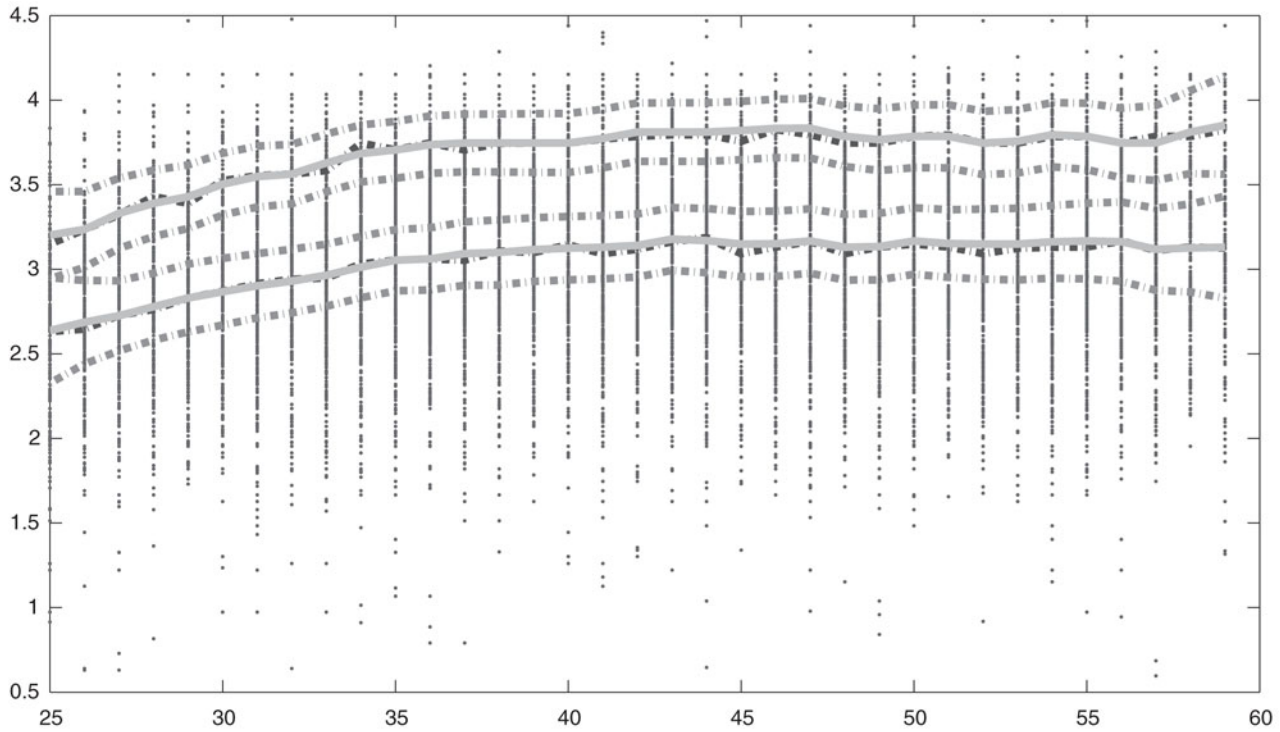


FIGURE 3. The original observations, local quantiles, 0.5- and 0.9-quantile smoothers, and corresponding 95% confidence bands.

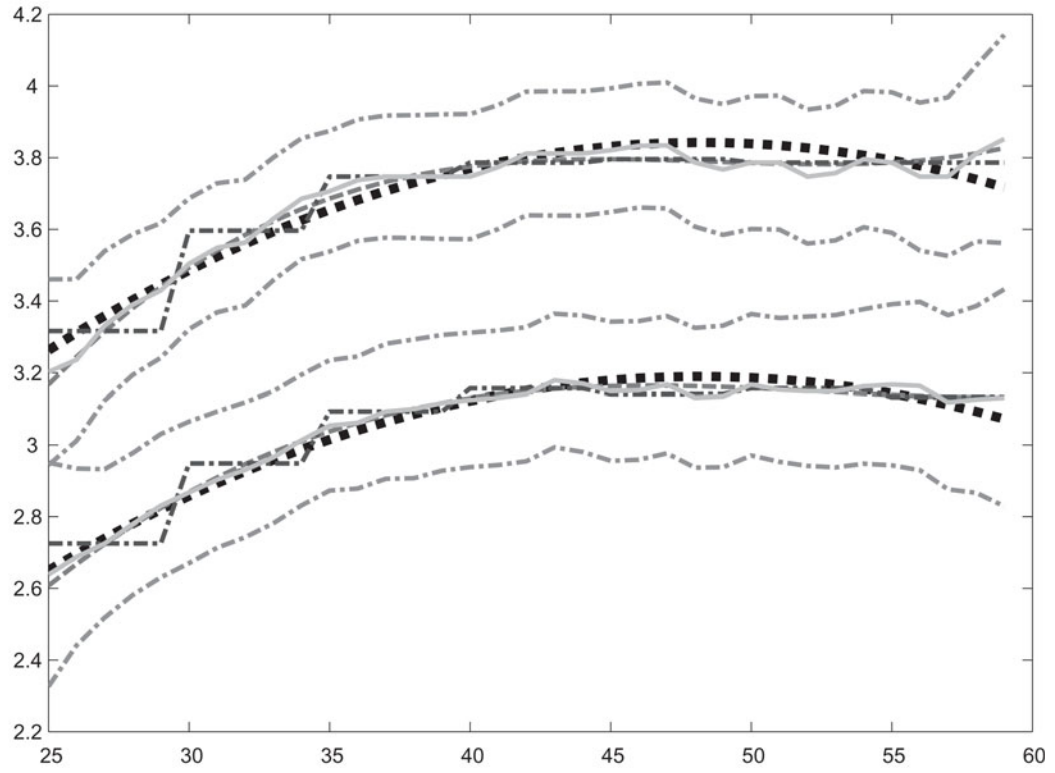


FIGURE 4. Quadratic, quartic, set of dummies (for age groups) estimates, 0.5- and 0.9-quantile smoothers, and their corresponding 95% confidence bands.

4. APPLICATION

Recently there has been great interest in finding out how the financial returns of a job depend on the age of the employee. We use the Current Population Survey (CPS) data from 2005 for the following group: male aged 25–59, full-time employed, and college graduate containing 16,731 observations, for the age-earning estimation. As is usual for wage data, a log transformation to hourly real wages (unit: U.S. dollar) is carried out first. In the CPS all ages (25–59) are reported as integers. We rescaled them into $[0, 1]$ by dividing 40 by bandwidth 0.059 for nonparametric quantile smoothers. This is equivalent to setting bandwidth 2 for the original age data.

In Figure 3 the original observations are displayed as small stars. The local 0.5 and 0.9 quantiles at the integer points of age are shown as dashed lines, whereas the corresponding nonparametric quantile smoothers are displayed as solid lines with corresponding 95% uniform confidence bands shown as dashed-dot lines. A closer look reveals a quadratic relation between age and logged hourly real wages. We use several popular parametric methods to estimate the 0.5 and 0.9 conditional quantiles, e.g., quadratic, quartic, and set of dummies (a dummy variable for each 5-year age group) models; the results are displayed in Figure 4. With the help of the 95% uniform confidence bands, we can conduct the parametric model specification test. At the 5% significance level, we could not reject any model. However, when the confidence level further decreases and the uniform confidence bands get narrower, the “set of dummies” parametric model will be the first one to be rejected. At the 10% significance level, the set of dummies (for age groups) model is rejected whereas the other two are not. As the quadratic model performs quite similarly to the quartic one, for simplicity it is suggested in practice to measure the $\log(\text{wage})$ -earning relation in mean regression, which coincides with the approach of Murphy and Welch (1990).

REFERENCES

- Bickel, P. & M. Rosenblatt (1973) On some global measures of the deviation of density function estimators. *Annals of Statistics* 1, 1071–1095.
- Cai, Z.W. (2002) Regression quantiles for time series. *Econometric Theory* 18, 169–192.
- Csörgő, S. & P. Hall (1982) Upper and lower classes for triangular arrays. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete* 61, 207–222.
- Fan, J., T.C. Hu, & Y.K. Troung (1994) Robust nonparametric function estimation. *Scandinavian Journal of Statistics* 21, 433–446.
- Fan, J., Q. Yao, & H. Tong (1996) Estimation of conditional densities and sensitivity measures in nonlinear dynamical systems. *Biometrika* 83, 189–206.
- Ferguson, T.S. (1967) *Mathematical Statistics: A Decision Theoretic Approach*. Academic Press.
- Franke, J. & P. Mwita (2003) Nonparametric Estimates for Conditional Quantiles of Time Series. Report in *Wirtschaftsmathematik* 87, University of Kaiserslautern.
- Hall, P., R. Wolff, & Q. Yao (1999) Methods for estimating a conditional distribution function. *Journal of the American Statistical Association* 94, 154–163.

- Härdle, W. (1989) Asymptotic maximal deviation of M -smoothers. *Journal of Multivariate Analysis* 29, 163–179.
- Härdle, W., P. Janssen & R. Serfling (1988) Strong uniform consistency rates for estimators of conditional functionals. *Annals of Statistics* 16, 1428–1429.
- Härdle, W. & S. Luckhaus (1984) Uniform consistency of a class of regression function estimators. *Annals of Statistics* 12, 612–623.
- Huber, P. (1981) *Robust Statistics*. Wiley.
- Jeong, K. & W. Härdle. (2008) A Consistent Nonparametric Test for Causality in Quantile. SFB 649 Discussion Paper.
- Johnston, G. (1982) Probabilities of maximal deviations of nonparametric regression function estimates. *Journal of Multivariate Analysis* 12, 402–414.
- Koenker, R. & G.W. Bassett (1978) Regression quantiles. *Econometrica* 46, 33–50.
- Koenker, R. & K.F. Hallock (2001) Quantile regression. *Journal of Econometric Perspectives* 15, 143–156.
- Koenker, R. & B.J. Park (1996) An interior point algorithm for nonlinear quantile regression. *Journal of Econometrics* 71, 265–283.
- Kong, E., O. Linton, & Y. Xia (2010) Uniform Bahadur representation for local polynomial estimates of M -regression and its application to the additive model. *Econometric Theory*, forthcoming.
- Lejeune, M.G. & P. Sarda (1988) Quantile regression: A nonparametric approach. *Computational Statistics and Data Analysis* 6, 229–239.
- Murphy, K. & F. Welch (1990) Empirical age-earnings profiles. *Journal of Labor Economics* 8, 202–229.
- Parzen, M. (1962) On estimation of a probability density function and mode. *Annals of Mathematical Statistics* 32, 1065–1076.
- Portnoy, S. & R. Koenker (1997) The Gaussian hare and the Laplacian tortoise: Computability of squared-error versus absolute-error estimators (with discussion). *Statistical Sciences* 12, 279–300.
- Rosenblatt, M. (1952) Remarks on a multivariate transformation. *Annals of Mathematical Statistics* 23, 470–472.
- Ruppert, D., S.J. Sheather, & M.P. Wand (1995) An effective bandwidth selector for local least squares regression. *Journal of the American Statistical Association* 90, 1257–1270.
- Tusnady, G. (1977) A remark on the approximation of the sample distribution function in the multidimensional case. *Periodica Mathematica Hungarica* 8, 53–55.
- Yu, K. & M.C. Jones (1997) A comparison of local constant and local linear regression quantile estimation. *Computational Statistics and Data Analysis* 25, 159–166.
- Yu, K. & M.C. Jones (1998) Local linear quantile regression. *Journal of the American Statistical Association* 93, 228–237.
- Yu, K., Z. Lu, & J. Stander (2003) Quantile regression: Applications and current research areas. *Journal of the Royal Statistical Society, Series D* 52, 331–350.

APPENDIX

Proof of Theorem 2.1 . By the definition of $l_n(x)$ as a zero of (2), we have, for $\varepsilon > 0$,

$$\text{if } l_n(x) > l(x) + \varepsilon, \quad \text{then } \tilde{H}_n\{l(x) + \varepsilon, x\} > 0. \quad (\text{A.1})$$

Now

$$\tilde{H}_n\{l(x) + \varepsilon, x\} \leq \tilde{H}\{l(x) + \varepsilon, x\} + \sup_{\theta \in I} |\tilde{H}_n(\theta, x) - \tilde{H}(\theta, x)|. \quad (\text{A.2})$$

Also, by the identity $\tilde{H}\{l(x), x\} = 0$, the function $\tilde{H}\{l(x) + \varepsilon, x\}$ is not positive and has a magnitude $\geq m_1 \tilde{q} \varepsilon$ by Assumption (A6) and (6), for $0 < \varepsilon < \delta_1$. That is, for $0 < \varepsilon < \delta_1$, $\tilde{H}\{l(x) + \varepsilon, x\} \leq -m_1 \tilde{q} \varepsilon$. (A.3)

Combining (A.1)–(A.3), we have, for $0 < \varepsilon < \delta_1$,

$$\text{if } l_n(x) > l(x) + \varepsilon, \quad \text{then } \sup_{\theta \in I} \sup_{x \in J} |\tilde{H}_n(\theta, x) - \tilde{H}(\theta, x)| > m_1 \tilde{q} \varepsilon.$$

With a similar inequality proved for the case $l_n(x) < l(x) + \varepsilon$, we obtain, for $0 < \varepsilon < \delta_1$,

$$\text{if } \sup_{x \in J} |l_n(x) - l(x)| > \varepsilon, \quad \text{then } \sup_{\theta \in I} \sup_{x \in J} |\tilde{H}_n(\theta, x) - \tilde{H}(\theta, x)| > m_1 \tilde{q} \varepsilon. \quad \text{(A.4)}$$

It readily follows that (A.4) and (5) imply (7). ■

Subsequently we first show that $\|R_n\|_\infty = \sup_{t \in J} |R_n(t)|$ vanishes asymptotically faster than the rate $(nh \log n)^{-1/2}$; for simplicity we will just use $\|\cdot\|$ to indicate the sup-norm.

LEMMA A.1. *For the remainder term $R_n(t)$ defined in (9) we have*

$$\|R_n\| = \mathcal{O}_p\{(nh \log n)^{-1/2}\}. \quad \text{(A.5)}$$

Proof. First we have by the positivity of the kernel K ,

$$\begin{aligned} \|R_n\| &\leq \left[\inf_{0 \leq t \leq 1} \{ |D_n(t)| \cdot q(t) \} \right]^{-1} \{ \|H_n\| \cdot \|q - D_n\| + \|D_n\| \cdot \|E H_n\| \} \\ &\quad + C_1 \cdot \|l_n - l\|^2 \cdot \left\{ \inf_{0 \leq t \leq 1} |D_n(t)| \right\}^{-1} \cdot \|f_n\|_\infty, \end{aligned}$$

where $f_n(x) = (nh)^{-1} \sum_{i=1}^n K\{(x - X_i)/h\}$.

The desired result, Lemma A.1, will then follow if we prove

$$\|H_n\| = \mathcal{O}_p\{(nh)^{-1/2} (\log n)^{1/2}\}, \quad \text{(A.6)}$$

$$\|q - D_n\| = \mathcal{O}_p\{(nh)^{-1/4} (\log n)^{-1/2}\}, \quad \text{(A.7)}$$

$$\|E H_n\| = \mathcal{O}(h^2), \quad \text{(A.8)}$$

$$\|l_n - l\|^2 = \mathcal{O}_p\{(nh)^{-1/2} (\log n)^{-1/2}\}. \quad \text{(A.9)}$$

Because (A.8) follows from the well-known bias calculation

$$E H_n(t) = h^{-1} \int K\{(t - u)/h\} E[\psi\{y - l(t)\} | X = u] f_X(u) du = \mathcal{O}(h^2),$$

where $\mathcal{O}(h^2)$ is independent of t in Parzen (1962), we have from Assumption (A2) that $\|E H_n\| = \mathcal{O}_p\{(nh)^{-1/2} (\log n)^{-1/2}\}$.

According to Lemma A.3 in Franke and Mwita (2003),

$$\sup_{t \in J} |H_n(t) - E H_n(t)| = \mathcal{O}\{(nh)^{-1/2} (\log n)^{1/2}\}$$

and the following inequality

$$\begin{aligned} \|H_n\| &\leq \|H_n - \mathbf{E} H_n\| + \|\mathbf{E} H_n\| \\ &= \mathcal{O}\left\{(nh)^{-1/2}(\log n)^{1/2}\right\} + \mathcal{O}_p\left\{(nh)^{-1/2}(\log n)^{-1/2}\right\} \\ &= \mathcal{O}\left\{(nh)^{-1/2}(\log n)^{1/2}\right\}, \end{aligned}$$

statement (A.6) thus is obtained.

Statement (A.7) follows in the same way as (A.6) using Assumption (A2) and the Lipschitz continuity properties of K, ψ', l .

According to the uniform consistency of $l_n(t) - l(t)$ shown before, we have

$$\|l_n - l\| = \mathcal{O}_p\{(nh)^{-1/2}(\log n)^{1/2}\},$$

which implies (A.9).

Now the assertion of the lemma follows, because by tightness of $D_n(t)$, $\inf_{0 \leq t \leq 1} |D_n(t)| \geq q_0$ a.s. and thus

$$\|R_n\| = \mathcal{O}_p\{(nh \log n)^{-1/2}\}(1 + \|f_n\|).$$

Finally, by Theorem 3.1 of Bickel and Rosenblatt (1973), $\|f_n\| = \mathcal{O}_p(1)$; thus the desired result $\|R_n\| = \mathcal{O}_p\{(nh \log n)^{-1/2}\}$ follows. ■

We now begin with the subsequent approximations of the processes $Y_{0,n} - Y_{5,n}$.

LEMMA A.2.

$$\|Y_{0,n} - Y_{1,n}\| = \mathcal{O}\left\{(nh)^{-1/2}(\log n)^2\right\} \quad a.s.$$

Proof. Let t be fixed and put $L(y) = \psi\{y - l(t)\}$ still depending on t . Using integration by parts, we obtain

$$\begin{aligned} &\iint_{\Gamma_n} L(y)K\{(t-x)/h\}dZ_n(x,y) \\ &= \int_{u=-A}^A \int_{y=-a_n}^{a_n} L(y)K(u)dZ_n(t-h \cdot u,y) \\ &= - \int_{-A}^A \int_{-a_n}^{a_n} Z_n(t-h \cdot u,y)d\{L(y)K(u)\} \\ &\quad + L(a_n)(a_n) \int_{-A}^A Z_n(t-h \cdot u,a_n)dK(u) \\ &\quad - L(-a_n)(-a_n) \int_{-A}^A Z_n(t-h \cdot u,-a_n)dK(u) \\ &\quad + K(A) \left\{ \int_{-a_n}^{a_n} Z_n(t-h \cdot A,y)dL(y) \right. \\ &\quad \left. + L(a_n)(a_n)Z_{n_a}(t-h \cdot A,a_n) - L(-a_n)(-a_n)Z_n(t-h \cdot A,-a_n) \right\} \end{aligned}$$

$$-K(-A) \left\{ \int_{-a_n}^{a_n} Z_n(t+h \cdot A, y) dL(y) + L(a_n)(a_n)Z_n(t+h \cdot A, a_n) - L(-a_n)(-a_n)Z_n(t+h \cdot A, -a_n) \right\}.$$

If we apply the same operation to $Y_{1,n}$ with $B_n\{T(x, y)\}$ instead of $Z_n(x, y)$ and use Lemma 2.2, we finally obtain

$$\sup_{0 \leq t \leq 1} h^{1/2} g(t)^{1/2} |Y_{0,n}(t) - Y_{1,n}(t)| = \mathcal{O} \left\{ n^{-1/2} (\log n)^2 \right\} \quad \text{a.s.} \quad \blacksquare$$

LEMMA A.3. $\|Y_{1,n} - Y_{2,n}\| = \mathcal{O}_p(h^{1/2})$.

Proof. Note that the Jacobian of $T(x, y)$ is $f(x, y)$. Hence

$$Y_{1,n}(t) - Y_{2,n}(t) = \left| \{g(t)h\}^{-1/2} \iint_{\Gamma_n} \psi\{y-l(t)\} K\{(t-x)/h\} f(x, y) dx dy \right| \cdot |W_n(1, 1)|.$$

It follows that

$$h^{-1/2} \|Y_{1,n} - Y_{2,n}\| \leq |W_n(1, 1)| \cdot \|g^{-1/2}\| \cdot \sup_{0 \leq t \leq 1} h^{-1} \iint_{\Gamma_n} |\psi\{y-l(t)\} K\{(t-x)/h\}| f(x, y) dx dy.$$

Because $\|g^{-1/2}\|$ is bounded by assumption, we have

$$h^{-1/2} \|Y_{1,n} - Y_{2,n}\| \leq |W_n(1, 1)| \cdot C_4 \cdot h^{-1} \int K\{(t-x)/h\} dx = \mathcal{O}_p(1). \quad \blacksquare$$

LEMMA A.4. $\|Y_{2,n} - Y_{3,n}\| = \mathcal{O}_p(h^{1/2})$.

Proof. The difference $|Y_{2,n}(t) - Y_{3,n}(t)|$ may be written as

$$\left| \{g(t)h\}^{-1/2} \iint_{\Gamma_n} [\psi\{y-l(t)\} - \psi\{y-l(x)\}] K\{(t-x)/h\} dW_n\{T(x, y)\} \right|.$$

If we use the fact that l is uniformly continuous, this is smaller than

$$h^{-1/2} |g(t)|^{-1/2} \cdot \mathcal{O}_p(h),$$

and the lemma thus follows. \blacksquare

LEMMA A.5. $\|Y_{4,n} - Y_{5,n}\| = \mathcal{O}_p(h^{1/2})$.

Proof.

$$\begin{aligned} |Y_{4,n}(t) - Y_{5,n}(t)| &= h^{-1/2} \left| \int \left[\left\{ \frac{g(x)}{g(t)} \right\}^{1/2} - 1 \right] K\{(t-x)/h\} dW(x) \right| \\ &\leq h^{-1/2} \left| \int_{-A}^A W(t-hu) \frac{\partial}{\partial u} \left[\left\{ \frac{g(t-hu)}{g(t)} \right\}^{1/2} - 1 \right] K(u) du \right| \end{aligned}$$

$$\begin{aligned}
 &+ h^{-1/2} \left| K(A)W(t-hA) \left[\left\{ \frac{g(t-Ah)}{g(t)} \right\}^{1/2} - 1 \right] \right| \\
 &+ h^{-1/2} \left| K(-A)W(t+hA) \left[\left\{ \frac{g(t+Ah)}{g(t)} \right\}^{1/2} - 1 \right] \right| \\
 &S_{1,n}(t) + S_{2,n}(t) + S_{3,n}(t), \quad \text{say.}
 \end{aligned}$$

The second term can be estimated by

$$h^{-1/2} \|S_{2,n}\| \leq K(A) \cdot \sup_{0 \leq t \leq 1} |W(t-Ah)| \cdot \sup_{0 \leq t \leq 1} h^{-1} \left| \left[\left\{ \frac{g(t-Ah)}{g(t)} \right\}^{1/2} - 1 \right] \right|.$$

By the mean value theorem it follows that

$$h^{-1/2} \|S_{2,n}\| = \mathcal{O}_p(1).$$

The first term $S_{1,n}$ is estimated as

$$\begin{aligned}
 h^{-1/2} S_{1,n}(t) &= \left| h^{-1} \int_{-A}^A W(t-uh)K'(u) \left[\left\{ \frac{g(t-uh)}{g(t)} \right\}^{1/2} - 1 \right] du \right. \\
 &\quad \left. \cdot \frac{1}{2} \int_{-A}^A W(t-uh)K(u) \left\{ \frac{g(t-uh)}{g(t)} \right\}^{1/2} \left\{ \frac{g'(t-uh)}{g(t)} \right\} du \right| \\
 &= |T_{1,n}(t) - T_{2,n}(t)|, \quad \text{say;}
 \end{aligned}$$

$\|T_{2,n}\| \leq C_5 \cdot \int_{-A}^A |W(t-hu)|du = \mathcal{O}_p(1)$ by assumption on $g(t) = \sigma^2(t) \cdot f_X(t)$. To estimate $T_{1,n}$ we again use the mean value theorem to conclude that

$$\sup_{0 \leq t \leq 1} h^{-1} \left| \left\{ \frac{g(t-uh)}{g(t)} \right\}^{1/2} - 1 \right| < C_6 \cdot |u|;$$

hence

$$\|T_{1,n}\| \leq C_6 \cdot \sup_{0 \leq t \leq 1} \int_{-A}^A |W(t-hu)|K'(u)u/du = \mathcal{O}_p(1).$$

Because $S_{3,n}(t)$ is estimated as $S_{2,n}(t)$, we finally obtain the desired result. ■

The next lemma shows that the truncation introduced through $\{a_n\}$ does not affect the limiting distribution.

LEMMA A.6. $\|Y_n - Y_{0,n}\| = \mathcal{O}_p\{(\log n)^{-1/2}\}$.

Proof. We shall only show that $g'(t)^{-1/2}h^{-1/2} \iint_{\mathbb{R}-\Gamma_n} \psi\{y-l(t)\}K\{(t-x)/h\}dZ_n(x, y)$ fulfills the lemma. The replacement of $g'(t)$ by $g(t)$ may be proved as in Lemma A.4 of Johnston (1982). The preceding quantity is less than $h^{-1/2}\|g^{-1/2}\| \cdot \iint_{\{|y|>a_n\}} \psi\{y -$

$l(\cdot)\}K\{(\cdot - x)/h\}dZ(x, y)\|$. It remains to be shown that the last factor tends to zero at a rate $\mathcal{O}_p\{(\log n)^{-1/2}\}$. We show first that

$$V_n(t) = (\log n)^{1/2}h^{-1/2} \iint_{\{|y|>a_n\}} \psi\{y - l(t)\}K\{(t - x)/h\}dZ_n(x, y)$$

$$\xrightarrow{p} 0 \quad \text{for all } t,$$

and then we show tightness of $V_n(t)$. The result then follows:

$$V_n(t) = (\log n)^{1/2}(nh)^{-1/2} \sum_{i=1}^n [\psi\{Y_i - l(t)\}\mathbf{1}(|Y_i| > a_n)K\{(t - X_i)/h\}$$

$$- \mathbb{E} \psi\{Y_i - l(t)\}\mathbf{1}(|Y_i| > a_n)K\{(t - X_i)/h\}]$$

$$= \sum_{i=1}^n X_{n,t}(t),$$

where $\{X_{n,t}(t)\}_{i=1}^n$ are i.i.d. for each n with $\mathbb{E} X_{n,t}(t) = 0$ for all $t \in [0, 1]$. We then have

$$\mathbb{E} X_{n,t}^2(t) \leq (\log n)(nh)^{-1} \mathbb{E} \psi^2\{Y_i - l(t)\}\mathbf{1}(|Y_i| > a_n)K^2\{(t - X_i)/h\}$$

$$\leq \sup_{-A \leq u \leq A} K^2(u) \cdot (\log n)(nh)^{-1} \mathbb{E} \psi^2\{Y_i - l(t)\}\mathbf{1}(|Y_i| > a_n).$$

Hence

$$\text{Var}\{V_n(t)\} = \mathbb{E} \left\{ \sum_{i=1}^n X_{n,t}(t) \right\}^2 = n \cdot \mathbb{E} X_{n,t}^2(t)$$

$$\leq \sup_{-A \leq u \leq A} K^2(u)h^{-1}(\log n) \int_{\{|y|>a_n\}} f_y(y) dy \cdot M_\psi,$$

where M_ψ denotes an upper bound for ψ^2 . This term tends to zero by Assumption (A3). Thus by Markov’s inequality we conclude that

$$V_n(t) \xrightarrow{p} 0 \quad \text{for all } t \in [0, 1].$$

To prove tightness of $\{V_n(t)\}$ we refer again to the following moment condition as stated in Lemma A.1:

$$\mathbb{E}\{|V_n(t) - V_n(t_1)| \cdot |V_n(t_2) - V_n(t)|\} \leq C' \cdot (t_2 - t_1)^2$$

C' denoting a constant, $t \in [t_1, t_2]$.

We again estimate the left-hand side by Schwarz’s inequality and estimate each factor separately:

$$\mathbb{E}\{V_n(t) - V_n(t_1)\}^2 = (\log n)(nh)^{-1} \mathbb{E} \left[\sum_{i=1}^n \Psi_n(t, t_1, X_i, Y_i) \cdot \mathbf{1}(|Y_i| > a_n) \right. \\ \left. - \mathbb{E}\{\Psi_n(t, t_1, X_i, Y_i) \cdot \mathbf{1}(|Y_i| > a_n)\} \right]^2,$$

where $\Psi_n(t, t_1, X_i, Y_i) = \psi\{Y_i - l(t)\}K\{(t - X_i)/h\} - \psi\{Y_i - l(t_1)\}K\{(t_1 - X_i)/h\}$. Because ψ, K are Lipschitz continuous except at one point and the expectation is taken afterward, it follows that

$$[E\{V_n(t) - V_n(t_1)\}^2]^{1/2} \leq C_7 \cdot (\log n)^{1/2} h^{-3/2} |t - t_1| \cdot \left\{ \int_{\{|y| > a_n\}} f_y(y) dy \right\}^{1/2}.$$

If we apply the same estimation to $V_n(t_2) - V_n(t_1)$ we finally have

$$\begin{aligned} E\{|V_n(t) - V_n(t_1)| \cdot |V_n(t_2) - V_n(t_1)|\} &\leq C_7^2 (\log n) h^{-3} |t - t_1| |t_2 - t_1| \times \int_{\{|y| > a_n\}} f_y(y) dy \\ &\leq C' \cdot |t_2 - t_1|^2 \quad \text{because } t \in [t_1, t_2] \quad \text{by Assumption (A3)}. \end{aligned}$$

LEMMA A.7. Let $\lambda(K) = \int K^2(u) du$ and let $\{d_n\}$ be as in Theorem 2.2. Then

$$(2\delta \log n)^{1/2} [\|Y_{3,n}\| / \{\lambda(K)\}^{1/2} - d_n]$$

has the same asymptotic distribution as

$$(2\delta \log n)^{1/2} [\|Y_{4,n}\| / \{\lambda(K)\}^{1/2} - d_n].$$

Proof. $Y_{3,n}(t)$ is a Gaussian process with

$$E Y_{3,n}(t) = 0$$

and covariance function

$$\begin{aligned} r_3(t_1, t_2) &= E Y_{3,n}(t_1) Y_{3,n}(t_2) \\ &= \{g(t_1)g(t_2)\}^{-1/2} h^{-1} \iint_{\Gamma_n} \psi^2\{y - l(x)\} K\{(t_1 - x)/h\} \\ &\quad \times K\{(t_2 - x)/h\} f(x, y) dx dy \\ &= \{g(t_1)g(t_2)\}^{-1/2} h^{-1} \iint_{\Gamma_n} \psi^2\{y - l(x)\} f(y|x) dy K\{(t_1 - x)/h\} \\ &\quad \times K\{(t_2 - x)/h\} f_X(x) dx \\ &= \{g(t_1)g(t_2)\}^{-1/2} h^{-1} \int g(x) K\{(t_1 - x)/h\} K\{(t_2 - x)/h\} dx \\ &= r_4(t_1, t_2), \end{aligned}$$

where $r_4(t_1, t_2)$ is the covariance function of the Gaussian process $Y_{4,n}(t)$, which proves the lemma. ■

INVESTORS' PREFERENCE: ESTIMATING AND DEMIXING OF THE WEIGHT FUNCTION IN SEMIPARAMETRIC MODELS FOR BIASED SAMPLES

Ya'acov Ritov and Wolfgang K. Härdle

The Hebrew University of Jerusalem and Humboldt-Universität zu Berlin

Abstract: We consider a semiparametric model for the weight function in a biased sample model. The object of our interest parametrizes the weight function, and it is non-Euclidean. The model discussed is motivated by the estimation of the mixing distribution of individual utility functions in the DAX market. We discuss the estimation rate of different functionals of the weight functions.

Key words and phrases: Empirical pricing kernel, exponential mixture, inverse problem, mixture distribution, risk aversion.

1. Introduction

A sample X_1, \dots, X_n is considered biased if it is sampled from a density p which is represented as

$$p(x) = \frac{q(x)w(x)}{\int q(u)w(u)du}. \quad (1.1)$$

Here q is some 'natural' pdf (probability density function) for the problem, representing the 'true' underlying distribution, while w is a given weight function that biases the sample. In a standard example, X represents the severity of the disease, and q is the density of X among patients at admission to the hospital. However, it may be more convenient to take a random sample from the population of patients who are in the hospital at a given time. If the time of hospitalization is proportional to the severity of the case, then the sample is taken from the density p , which is equal to q 'length biased' with $w(x) \equiv x$. Vardi (1985) was the first to systematically analyze these models; asymptotic theory was developed in Gill, Vardi and Wellner (1988); Gilbert, Lele and Vardi (1999) extended the model to the situation where the weight function depends on some parameter, $w(x) = w(x; f)$; the large sample properties were discussed in Gilbert (2000). Equation (1.1) has some similarities to the classical choice-based sample problem, Manski and Lerman (1977), or retrospective case-control studies, Mantel (1973). In fact one can consider the situation as if one has an infinite

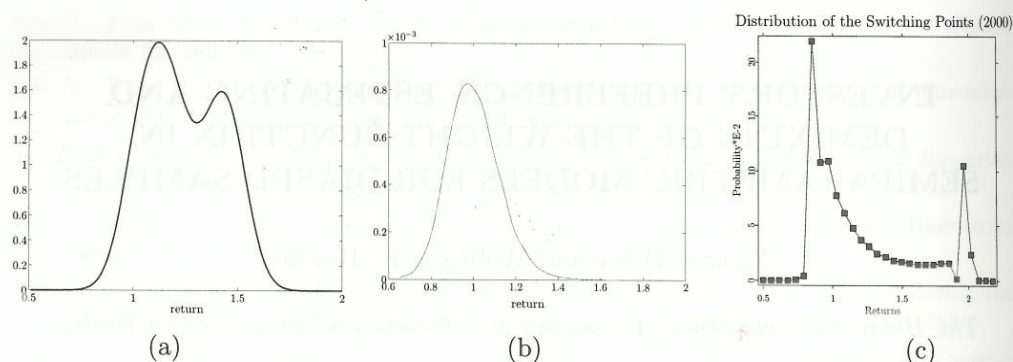
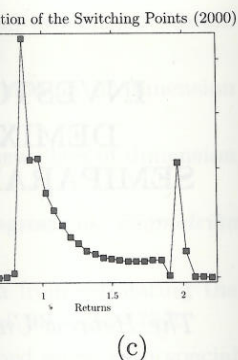


Figure 1. The DAX data, 24/03/2000 half a year look ahead: (a) p , the historical density; (b) q , the risk neutral density; (c) The estimate of f , the mixing density. Figures are taken from DHM.

sample from the control group, and hence q is known, and a finite sample from the control, the biased sample. The likelihood ratio between the two is the given $w(x; f)$. The main difficulty we face in this paper is the particular form of $w(x; f)$ we have.

Technically speaking, our paper is about estimating f , the parameter of the weight function, $w(x) = w(x; f)$. In the model we consider, q is taken as known, while the weight function is parametrized by a non-Euclidean parameter. This brings us to an inverse problem of estimating and demixing the weight function.

In subject matter, our model is motivated by the research on risk aversion and proclivity, and more precisely on the empirical pricing kernel (EPK), see Detlefsen, Härdle and Moro (2007) (hereafter DHM). The EPK describes the apparent utility behavior as function of the individual investors utility function. In this model q is the risk neutral density of asset pricing, and is derived from theoretical considerations. The density p on the other hand is the density of the empirical (historical) prices. See parts (a) and (b) of Figure 1 for an example. In asset pricing the EPK links a risk neutral investor's behavior to individual utilities, which gives in our notation a semiparametric modeling of the weight function w . The integral function of the pricing kernel q/p is the utility function used by a representing individual. Knowing p and q yields the exact form of the utility function, cf. Ait-Sahalia and Lo (2000), and Rosenberg and Engle (2002). The risk neutral (state price) density (SPD) q can be calculated from market data on European options. There are more than 5,000 observations each day for maturity from one week to two years. The SPD can therefore be estimated very precisely. Much empirical research work has demonstrated the so called EPK paradox: the resulting utility function is partially concave and partially convex, more precisely of the Friedman and Savage type, Friedman and Savage (1948).



(c)
ad: (a) p , the
mate of f , the

finite sample from
the two is the given
ular form of $w(x; f)$

the parameter of the
is taken as known,
an parameter. This
the weight function.
ch on risk aversion
kernel (EPK), see
EPK describes the
ors utility function.
and is derived from
s the density of the
e 1 for an example.
avior to individual
deling of the weight
the utility function
the exact form of the
g and Engle (2002).
culated from market
uations each day for
e be estimated very
the so called EPK
nd partially convex,
and Savage (1948).

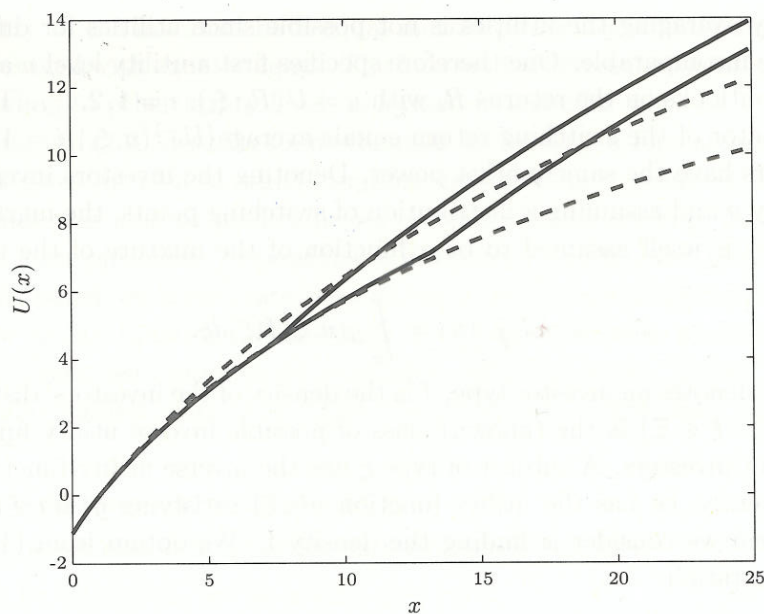


Figure 2. The utility function $U(\cdot; \xi)$ of (3.5) ($\alpha_1 = 2, \alpha_2 = 2.25, c = 2$) for two different values of ξ (solid lines), and of (3.8) for two values broken lines.

This so called risk aversion puzzle has also been recently discussed in Chabi-Yo, Garcia and Renault (2008); a recursive utility approach to dynamic pricing kernel estimation is published in Gallant and Hong (2007); a fundamental reference on asset pricing theory is the book by Cochrane (2005).

It is assumed in DHM that the observed density of the DAX value has density of the form $p(x) = cq(x)w(x; f)$, where $q \in \{q_\nu, \nu \in N \subseteq \mathbb{R}^d\}$ is the theoretical derived risk neutral density, assumed to follow a given parametric function, and c is a normalization factor, that is, of the type (1.1). The weight function is theoretically derived as

$$w(x; f) = \frac{1}{U'}(x), \tag{1.2}$$

where U is the market utility function, and prime denotes derivative. The market utility is estimated for option data and available historical data, and it also showed the risk aversion puzzle for the DAX stock market. In DHM an aggregation mechanism was proposed that similarly to Chabi-Yo, Garcia and Renault (2008) uses a switching point ξ . This point characterizes the investors switch from a bearish (low return) to a bullish (high return) risk aversion pattern. A graph of two different utility functions $u(\cdot; \xi)$ with switching points $\xi_1 < \xi_2$ is presented in Figure 2.

Simply averaging the utilities is not possible since utilities for different investors are incomparable. One therefore specifies first a utility level u and aggregates the outlooks on the returns R_i with $u = U(R_i; \xi_i)$, $i = 1, 2, \dots$. The aggregate estimator of the switching return equals average $\{U^{-1}(u, \xi_i), i = 1, 2, \dots\}$ if all investors have the same market power. Denoting the investors inverse utility function by g and assuming a distribution of switching points, the market utility function U_f is itself assumed to be a function of the mixture of the individual investors:

$$x = U_f^{-1}(u) = \int_{\Xi} g(u; \xi) f(\xi) d\xi. \quad (1.3)$$

Here $\xi \in \Xi$ denotes an investor type, f is the density of the investors' distribution, and $\{g(\cdot; \xi) : \xi \in \Xi\}$ is the (known) class of possible inverse utility functions of the different investors. A subject of type ξ has the inverse utility function $g(\cdot; \xi)$ or, equivalently, he has the utility function $u(\cdot; \xi)$ satisfying $g\{u(x; \xi); \xi\} \equiv x$. The problem we consider is finding the density f . We obtain from (1.1)–(1.3) the representation:

$$p(x) = cq(x) \int \frac{\partial}{\partial u} g(u; \xi) f(\xi) d\xi,$$

where u solves

$$x = \int g(u; \xi) f(\xi) d\xi. \quad (1.4)$$

See Figure 1 for an example taken from DHM of estimates of p , q , and f . See also Figure 2 for an example of $g^{-1}(\cdot; \xi)$.

Aggregation problem (1.3) is a way of aggregating preferences that is not based on the equilibrium theory usually associated with Walras (1874). The situation considered here is of a different type and is hypothetical when applied to real markets. The DAX market data were mentioned as suitable for testing the disaggregation techniques described in the paper.

Aggregation procedure (1.3) relates to the situation where the price of an asset is obtained as the result of a survey of investors (or experts) before they made trades. Thus, this price should be considered as a forecast for the next period, not a reflection of the struggle for limited resources in the market between investors with different preferences and endowments.

The survey proceeds as following. Each market participant is asked what the price will be if the conditions in the market are, for example, extremely good. Extremely good corresponds to some utility level \tilde{u}_1 in the minds of investors. In this way all investors agree that they are discussing an economic situation with the same utility level. As the next step, each investor forms his forecast about how high the prices would be in such a situation. Those forecasted prices are recorded and averaged to produce an aggregate opinion of all market participants

(or experts). If the investors have equal market power, their individual opinions will be averaged with equal weights. The forecast for different economic situations corresponding to other utility levels is formed in a similar way.

To sum up, (1.3) describes a mechanism for forming a forecast about future prices. It gives an idea of which opinions prevailed in a group of investors or experts that was able to predict prices correctly before trading, for example if they were more optimistic or pessimistic investors (experts), and to what degree.

In this paper we investigate the estimation of the non-Euclidean parameter f of a few utility functions. The result is typical for inverse problems, in that slightly different assumption yield completely different results. In fact, we present three similar models, similar to those investigated in DHM, that exhibit these behaviors:

- (i) there is no consistent estimator of f ;
- (ii) f can be estimated at a regular nonparametric rate of $n^{-\alpha}$;
- (iii) f can be estimated, but at a very slow rate.

Interestingly, there is a sort of uncertainty principle: the better we can estimate the function $U^{-1}(u)$, the worse we can demix it and estimate f . This is not unexpected. We cannot estimate f well when large differences in f have only minor impact on $\int g(\cdot; \xi) f(\xi) d\xi$.

The structure of the rest of the paper is as follows. In Section 2, we suggest an algorithm for calculating the generalized maximum-likelihood estimator (GMLE) for the semiparametric weight function of the model suggested by DHM. Rates of convergence of the demixing estimator for the DHM's model are discussed in Section 3, as well as of estimates of the mixture itself.

2. EPK: Model and an EM estimator

We consider the EPK problem. We start from (1.4) and we assume that q is known. In practice, it is assumed only to belong to some parametric family $\{q_\nu\}$. However, we deal in the following with rates that are much slower than the parametric \sqrt{n} rate, and the estimate of ν is based on a much larger sample than the estimates of the rest of the parameters. Therefore, the assumption that ν is known considerably simplifies the discussion without impacting the results.

Rewrite (1.4) as

$$p \left\{ \int g(u; \xi) f(\xi) d\mu(\xi) \right\} \int \frac{\partial}{\partial u} g(u; \xi) f(\xi) d\mu(\xi) \\ = cq \left\{ \int g(u; \xi) f(\xi) d\mu(\xi) \right\} \left\{ \int \frac{\partial}{\partial u} g(u; \xi) f(\xi) d\mu(\xi) \right\}^2, \quad (2.1)$$

where μ is some dominating measure (e.g., Lebesgue or the counting measure). Noting that the LHS of (2.1) integrates to 1, c can be found to yield

$$p \left\{ \int g(u; \xi) f(\xi) d\mu(\xi) \right\} = \frac{q \left\{ \int g(u; \xi) f(\xi) d\mu(\xi) \right\} \int \frac{\partial}{\partial u} g(u; \xi) f(\xi) d\mu(\xi)}{\int q \left\{ \int g(v; \xi) f(\xi) d\mu(\xi) \right\} \left\{ \int \frac{\partial}{\partial u} g(v; \xi) f(\xi) d\mu(\xi) \right\}^2 dv}.$$

The market utility $U(x) = U(x; f)$ is given by

$$x \equiv \int g \left\{ U(x; f); \xi \right\} f(\xi) d\mu(\xi) \equiv \psi_f \left\{ U(x; f) \right\}.$$

We obtain

$$p(x) = \frac{q(x) \int \frac{\partial}{\partial u} g(U(x; f); \xi) f(\xi) d\mu(\xi)}{\int q(y) \int \frac{\partial}{\partial u} g(U(y; f); \xi) f(\xi) d\mu(\xi) dy} = \frac{q(x) \psi'_f \left\{ \psi_f^{-1}(x) \right\}}{\int q(y) \psi'_f \left\{ \psi_f^{-1}(y) \right\} dy}. \quad (2.2)$$

The statistical model assumed by DHM is that we obtain a simple random sample from p , where p is parametrized in (2.2) by the non-Euclidean parameter f . A natural approach is to estimate f by the MLE or a variant of it, which we develop now. Note that $\nabla_f \psi_f(u) = g(u; \cdot)$, and by taking the gradient of $x \equiv \int g \left\{ \psi_f^{-1}(x); \xi \right\} f(\xi) d\mu(\xi)$ we obtain

$$0 = g \left\{ \psi_f^{-1}(x); \cdot \right\} + \psi'_f \left\{ \psi_f^{-1}(x) \right\} \nabla_f \psi_f^{-1}(x).$$

The derivative of the log-likelihood is given therefore by

$$\begin{aligned} \dot{\ell}_f(\xi) &= \sum_{i=1}^n \frac{1}{\psi'_f \left\{ \psi_f^{-1}(X_i) \right\}} \left[\frac{\partial}{\partial u} g \left\{ \psi_f^{-1}(X_i); \xi \right\} - \frac{\psi''_f \left\{ \psi_f^{-1}(X_i) \right\}}{\psi'_f \left\{ \psi_f^{-1}(X_i) \right\}} g \left\{ \psi_f^{-1}(X_i); \xi \right\} \right] \\ &\quad - n A_f(\xi), \\ &= \sum_{i=1}^n \frac{1}{\psi'_f \left\{ U_i \right\}} \left\{ \frac{\partial}{\partial u} g \left\{ U_i; \xi \right\} - \frac{\psi''_f(U_i)}{\psi'_f(U_i)} g(U_i; \xi) \right\} - n A_f(\xi), \end{aligned}$$

with $U_i = \psi_f^{-1}(X_i)$, and for all $\xi \in \text{supp} f$, where $A_f(\xi)$ is the mean of the first term under f . Since the density of U_i is given by

$$r_f(u) = p \left\{ \psi_f(u) \right\} \psi'_f(u) = \frac{q \left\{ \psi_f(u) \right\} \left\{ \psi'_f(u) \right\}^2}{\int q \left\{ \psi_f(v) \right\} \left\{ \psi'_f(v) \right\}^2 dv},$$

we obtain that

$$A_f(\xi) = \frac{\int q \left\{ \psi_f(u) \right\} \left\{ \psi'_f(u) \right\} \frac{\partial}{\partial u} g(u; \xi) - \psi''_f(u) g(u; \xi) du}{\int q \left\{ \psi_f(v) \right\} \left\{ \psi'_f(v) \right\}^2 dv}.$$

We discuss now how a GMLE can be constructed, and suggest a pseudo-EM algorithm, that is justified as being the limiting result of proper EM algorithms

applied in approximate models. To be clear, the approximation introduced in the following is needed only as a justification for an algorithm applied to the formal model. The algorithm itself is "exact" and maximizes the exact likelihood. The technical problem we want to circumvent is the exact functional dependency of X_i and U_i which affects the EM. As an intermediate step we weaken the functional dependency into a proper statistical dependency.

The model of a random sample from the density p can be well-approximated as $\sigma \rightarrow 0$ by a $X_i = \psi_f(U_i) + \varepsilon_i, i = 1, \dots, n$, where $\varepsilon_1, \dots, \varepsilon_n$ is a random sample from $N(0, \sigma^2)$ independent from the random sample U_1, \dots, U_n taken from the density r_f . Now, the log-likelihood of the joint density is given by

$$\ell_f = \sum_{i=1}^n \left[\log q\{\psi_f(U_i)\} + 2 \log\{\psi'_f(U_i)\} \right] - nC_f - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \psi_f(U_i))^2, \tag{2.2}$$

where $C_f = \log \int q\{\psi_f(v)\}\{\psi'_f(v)\}^2 dv$. By a well-known formula for the Bayes estimator in the Gaussian measurement error model, here the distribution of $\psi_f(U_i) - X_i$, given X_i , is normal with mean $\sigma^2 f'_X(X_i)/f_X(X_i)$ and second moment $\sigma^4 f''_X(X_i)/f_X(X_i) + \sigma^2$, where f_X is the marginal density of X_i . At the limit as $\sigma^2 \rightarrow 0$, the conditional expectation of the log-likelihood, given the X_i 's, amounts therefore to replacing U_i by $\psi_f^{-1}(X_i)$. We conclude that the limiting EM algorithm iterates therefore between the following steps.

The E step:

$$U_i \leftarrow \psi_f^{-1}(X_i), \quad i = 1, \dots, n, \tag{2.3}$$

The M step:

$$f \leftarrow \operatorname{argmax} \left[\sum_{i=1}^n \left\{ \log q\{\psi_f(U_i)\} + 2 \log\{\psi'_f(U_i)\} \right\} - nC_f \right].$$

Let $\mathbf{U} = (U_1, \dots, U_n)$, $\mathbf{X} = (X_1, \dots, X_n)$, and denote the E-step by $\mathbf{U} = \psi_f^{-1}(\mathbf{X})$. The M-step can be accomplished by solving the likelihood equation:

$$0 = \dot{\ell}_f^M(\xi; \mathbf{U}) = \sum_{i=1}^n \left[\frac{q'\{\psi_f(U_i)\}}{q\{\psi_f(U_i)\}} g(U_i; \xi) + \frac{2}{\psi'_f(U_i)} \frac{\partial}{\partial u} g(U_i, \xi) - \dot{C}_f(\xi) \right], \tag{2.4}$$

for all $\xi \in \operatorname{supp} f$, where

$$\begin{aligned} \dot{C}_f(\xi) &= \frac{\int [(q'\{\psi_f(v)\}/q\{\psi_f(v)\})g(v; \xi) + (2/\psi'_f(v)) \frac{\partial}{\partial u} g(v, \xi)] q\{\psi_f(v)\}\{\psi'_f(v)\}^2 dv}{\int q\{\psi_f(v)\}\{\psi'_f(v)\}^2 dv} \\ &= E_f \left[\frac{q'\{\psi_f(U)\}}{q\{\psi_f(U)\}} g(U; \xi) + \frac{2}{\psi'_f(U)} \frac{\partial}{\partial u} g(U, \xi) \right] \\ &= E_f \{T_f(U; \xi)\}, \quad \text{say.} \end{aligned}$$

g measure).

$$\frac{d\mu(\xi)}{\mu(\xi)^2 dv}$$

$$\frac{1}{dy} \tag{2.2}$$

dom sample
meter f . A
, which we
ient of $x \equiv$

$(X_i; \xi)$

of the first

pseudo-EM
I algorithms

However, there is no need in the M-step to find the exact maximizer of the log-likelihood. All that is needed is that the likelihood be strictly increasing (if possible at all) at every M-step. Therefore, the exact M-step given above can be replaced by an approximate M-step, that is obtained by considering an approximate Newton-Raphson solution of (2.4), where the $\mathcal{O}_p(\sqrt{n})$ terms in the Hessian of the log-likelihood are discarded. That is the term

$$\sum_{i=1}^n \left\{ \nabla_f T_f(U_i; \xi) - E_f \nabla_f T_f(U; \xi) \right\}.$$

We consider therefore the Newton-Raphson EM (NR-EM) algorithm:

$$f_{i+1} = \begin{cases} \tilde{f}_i \triangleq f_i + H_{f_i}^{-1} \ell_{f_i}^M \{ \cdot; \psi_{f_i}^{-1}(\mathbf{X}) \} & \ell_{\tilde{f}_i} > \ell_{f_i} \\ \text{the solution of (2.3)} & \text{otherwise,} \end{cases}$$

where $H_f : L_2(\mu) \rightarrow L_2(\mu)$ is the operator $H_f(\xi, \zeta) = \text{Cov}_f \{ T_f(U; \xi), T_f(U; \zeta) \}$.

3. EPK: Rates of Convergence

In the previous section we considered the MLE estimate of f . In this section we consider simple estimators of the type suggested by DHM. Using these estimators we will be able to discuss possible minimax rates of convergence. In essence, we start with a naive nonparametric estimator of the mixture, and in the second step we improve it or demix it for f .

One simple method for demixing the EPK is to start with (1.4) which can be written as

$$1 = c \int \frac{\partial}{\partial u} g(u; \xi) f(\xi) d\xi \frac{q}{p} \left\{ \int g(u; \xi) f(\xi) d\xi \right\} = c \frac{\partial}{\partial u} \frac{q}{p} \left\{ \int g(u; \xi) f(\xi) d\xi \right\}.$$

Hence $q/p \{ \int g(u; \xi) f(\xi) d\xi \} = \alpha + \beta u$ for some α and β , or

$$\int g(u; \xi) f(\xi) d\xi = \left(\frac{p}{q} \right)^{-1} (\alpha + \beta u). \quad (3.1)$$

The utility function of an individual is defined up to affine transformation. To assure that it is well defined, we assume that at the return of 1 the value of the utility is 0, and that of the derivative is 1. In terms of the inverse utility function this translates to $g(0, \xi) \equiv \frac{\partial}{\partial u} g(0, \xi) \equiv 1$. Hence

$$\begin{aligned} \alpha &= \frac{p(1)}{q(1)} \\ \beta &= \frac{p'(1)}{q(1)} - \frac{p(1) q'(1)}{q(1) q(1)}. \end{aligned} \quad (3.2)$$

The parameter f is therefore the solution of

$$\int g(u; \xi) f(\xi) d\xi = \psi(u) \tag{3.3}$$

for some ψ given explicitly by (3.1) and (3.2). Since g is estimated as a parametric density (based on a much larger sample), and p can be estimated at a standard non-parametric rate based on a direct sample from p , ψ can as well be estimated at a regular density estimation rate.

The analysis of this section starts with (3.3). We assume that ψ and its relevant derivatives can be estimated at a polynomial rate $\|\hat{\psi}^{(i)} - \psi^{(i)}\|_\infty = \mathcal{O}_p(n^{-\alpha_i})$ for some $\alpha_i > 0$. The natural estimator suggested by DHM is given by the inverse function of a weighed density estimator. Under strict monotonicity and boundness, the inverse function inherits most properties from the density kernel estimator.

Note that model (3.3) looks like a linear model. For example, if f is approximated by a finite distribution with point mass at ξ_1, \dots, ξ_m , and (3.3) is considered at the k points u_1, \dots, u_k , then it can be written as

$$\hat{\psi}(u_i) = \sum_{j=1}^m \beta_j g(u_i; \xi_j) + \varepsilon_i, \quad i = 1, \dots, k. \tag{3.4}$$

(3.4) looks like a standard linear model and, indeed, we suggest estimating f by solving it. However, it is not. Most linear model assumptions are violated, e.g., $\varepsilon_1, \dots, \varepsilon_k$ are not i.i.d. and they are not independent of the random u_1, \dots, u_k .

The basic idea of this section is as follow. We assume that we have some naive nonparametric estimator of ψ . We then proceed to use the pseudo linear model (3.4) to to estimate the mixing distribution and to improve the estimate of ψ itself. We show that this method yields the minimax rates.

How fast can f be estimated? In the rest of the section we present simple examples following DHM. These examples show that in a very similar models very different types of behavior can be obtained. It can be that (i) There is no consistent estimator of f ; (ii) f can be estimated at a regular nonparametric rate of $n^{-\alpha}$; (iii) f can be estimated but at a very slow rate. Thus one can suspect that any optimistic result of demixing depends too heavily on assumptions, and are *a priori* not robust (at least in the minimax sense). In particular, any result should be checked to stand against different changes in the model.

3.1. Switching between two utilities

Following DHM assume that for $x, \xi > 0$,

$$U(x; \xi) = \alpha_2(1 - c)^{1-1/\alpha_2} \left\{ [x - \xi]_+^{1/\alpha_1} \vee (x - c)^{1/\alpha_2} \right\} - \alpha_2(1 - c), \tag{3.5}$$

where $\alpha_2 > \alpha_1 > 1$ are given, $c < 0$, and $[x]_+ = x1(x > 0)$. See Figure 2. Then

$$g(u; \xi) = \min \left\{ \beta^{\alpha_2} \{u + \alpha_2(1 - c)\}^{\alpha_2} + c, \beta^{\alpha_1} \{u + \alpha_2(1 - c)\}^{\alpha_1} + \xi \right\},$$

where $\beta = \alpha_2^{-1}(1 - c)^{-1+1/\alpha_2}$. To simplify the notation and generalize the discussion, we consider a slightly more general case.

Theorem 3.1. *Suppose q is known and bounded away from 0 on a open interval, p has $s > 2$ bounded derivatives, and*

$$g(u; \xi) = \begin{cases} g_2(u) & -\infty < u \leq h(\xi) \\ g_1(u) + \xi & \infty > u > h(\xi) \end{cases}, \quad \xi > 0,$$

where g_1, g_2 are continuous with bounded derivatives, and h given by

$$h^{-1} = g_2 - g_1 \tag{3.6}$$

is a strictly increasing function. Then, f can be estimated with an $\mathcal{O}_p(n^{-(s-2)/(2s+1)})$ error.

Proof. Note that $g(u; \xi)$ is continuous in ξ . Equation (3.3) can be translated to

$$\psi(u) = \int^{h^{-1}(u)} \xi f(\xi) d\xi + g_2(u)F\{h^{-1}(u)\} + g_2(u)\{1 - F\{h^{-1}(u)\}\},$$

where F is the cdf corresponding to the pdf f . Changing variables and considering (3.6),

$$\psi\{h(s)\} = \int^s \xi f(\xi) d\xi - sF(s) + g_2\{h(s)\}.$$

Taking a derivative gives $F(s) = h'(s)\{g_2'\{h(s)\} - \psi'\{h(s)\}\}$. Hence estimating F at s is equivalent to the estimation of ψ' at $h(s)$. In other words, $f(\cdot)$ can be estimated at the same rate as the rate of the estimation of second derivative of ψ , which in turn is governed by the rate of estimation of the second derivative of p . Since, by assumption, p has s bounded derivatives, f can be estimated with an $\mathcal{O}_p(n^{-(s-2)/(2s+1)})$ error, cf. Silverman (1986).

3.2. Polynomial and exponential inverse utility function

Theorem 3.1 described a relatively optimistic example. However, modest changes in the inverse utility function may create situations in which f can hardly be estimated, or even not at all.

Here is a pessimistic example:

Theorem 3.2. *Suppose the CRRA (constant relative risk aversion) utility*

$$g(u; \zeta) = (\alpha \zeta^{\alpha-1})^{-1} \left\{ (u + \zeta)^\alpha - \zeta^\alpha \right\} + 1, \quad u \in \mathbb{R}, \zeta \in \mathbb{R}^+, \quad (3.7)$$

where α is a known integer. Then there is no consistent estimator of f .

Note that g in (3.7) is scaled such that both its value and its derivative at zero are equal to 1, that is, it represents one branch of (3.5). The proof of Theorem 3.2 is simple. Since α is an integer, $\psi(\cdot)$ is a function of f only through its first α moments. Hence, these moments can be estimated, but no other aspects of f can be estimated or identified.

Seemingly, more and more moments are revealed as $\alpha \rightarrow \infty$, and therefore, by the above argument, f is going to be identified at the limit. However, it is not clear that the high moments can be estimated effectively. We consider the limiting case explicitly. The limiting form of the inverse utility function, as $\alpha \rightarrow \infty$ and $\alpha/\zeta \rightarrow \xi$, is given by

$$g(u; \xi) \equiv \xi^{-1} (e^{u\xi} - 1) + 1. \quad (3.8)$$

The density f is now identified. For example, all its moments can be estimated, e.g., by $\int \xi^i f(\xi) d\xi = \psi^{(i+1)}(0)$. We are now going to analyze this model in some detail. We will argue that if $f(\cdot)$ is assumed to have two bounded derivatives, then its value at a point can indeed be estimated, but this can be done only at a very slow convergence rate, slower than any polynomial rate.

Theorem 3.3. *Assume that g is given by (3.8) and f is bounded and has two bounded derivatives. Suppose the minimax rate of estimation of ψ is n^γ , $\gamma \in (0, 1/2)$. Then there is an estimator \hat{f} such that $\hat{f}(s) - f(s) = \mathcal{O}_p(n^{-\alpha \log \log n / \log n})$ for some α , and for any $\alpha > 0$ there is no estimator $\tilde{f}(s)$ such that $\tilde{f}(s) - f(s) = \mathcal{O}_p(n^{-\alpha / \log \log n})$.*

The proof is given in the on-line supplement, see <http://www.stat.sinica.edu.tw/statistica>.

3.3. Smoothing the empirical estimate and an uncertainty principle

We start, as in the previous subsections, with a nonparametric $\hat{\psi}$. The purpose of this subsection is to show that this initial estimator can be improved considerably by a simple projection.

We argued in Subsection 3.2 that there is no reasonable estimator of f for g given in (3.8). If (3.8) is believed to be true, does this mean that there is nothing to do? The surprising answer is no. Although f cannot be estimated per-se, many

of its functionals can be estimated quite easily and quite well. For example, as mentioned in Subsection 3.2, its moments. Similarly $\psi(u)$, another functional of f , can be estimated quite easily, considered as a simple linear functional.

Suppose that f is supported on some compact interval $[a, b]$. Then one can approximate $\psi(u) = \sum_{i=1}^m \beta_i u^i + R_m(u)$, where, for some $\tilde{u} \in (0, u)$;

$$0 \leq R_m(u) = \frac{1}{(m+1)!} \psi^{m+1}(\tilde{u}) = \frac{1}{(m+1)!} \int_a^b \xi^m e^{\tilde{u}\xi} f(\xi) d\xi \leq \frac{b^m e^{ub}}{(m+1)!}. \quad (3.9)$$

Generally speaking, the faster the coefficients β converge to 0, the easier it is to estimate ψ and the harder it is to estimate the mixing density g . As (3.9) shows, we need only a few terms to approximate ψ quite well. In fact we show that in this smooth case, where as on the one hand f can be hardly estimated, ψ can be estimated almost at the parametric rate. This is not an accident — these are two faces of one phenomena. The shape of the observable ψ hardly depends on the fine details of f , and essentially depends only on a few aspects of f . These aspects can be estimated well (and hence ψ can be estimated quite precisely). The other aspects can hardly be estimated and hence f cannot be estimated in a reasonable rate. This yields an uncertainty principle — the more you are certain about ψ the less certain you are about f .

Recall that a function g is called completely monotone if $(-1)^k g^{(k)} \geq 0$, and it is called a Bernstein function if its first derivative is completely monotone. It is well-known (Feller (1966)) that g is completely monotone if, and only if, $g(u) = \int_0^\infty e^{-u\xi} dF(\xi)$. In other words, ψ is a Bernstein function. Nonparametric maximum likelihood estimation for an exponential mixture (and hence completely monotone density) was discussed in Jewell (1982). Balabdaoui and Wellner (2007) discussed the estimation of a k -monotone density.

We assume that there is an estimate $\hat{\psi} = \hat{\psi}_n$ at our disposal. For any $u_1, \dots, u_k > 0$, let $\Sigma(u_1, \dots, u_k) \in \mathbb{R}^{k \times k}$, where $\Sigma_{ij}(u_1, \dots, u_k) = \text{Cov}\{\hat{\psi}(u_i), \hat{\psi}(u_j)\}$. Consider the following assumption:

Assumptions 1. For any n there is $k = k_n$ and $u_1, \dots, u_k \in (c, d)$, $0 < c < d$, such that the spectral radius of $\Sigma(u_1, \dots, u_k)$ is $\mathcal{O}(k/n)$, and $\max_i |\mathbb{E}\psi(u_i) - \psi(u_i)|^2 = \mathcal{O}(\log n/n)$.

Assumption 1 is satisfied by many nonparametric density and regression estimators, when they strictly under-smooth. We care much more about bias than about variance of the original estimator $\hat{\psi}$. Thus, we have in mind a kernel estimator with bandwidth of order $n^{-1/4+\epsilon}$. The spectral radius is based on the assumptions that the estimator at points that are a multiple of the bandwidth apart are (almost) independent, for example this is trivially the case with kernel estimators having a compact support. The relationships in the assumption

obtain when the bias of the estimator is $\mathcal{O}(\sigma^2)$, the variance is $\mathcal{O}(1/n\sigma)$, and $k = \mathcal{O}(\sigma^{-1})$.

Consider now the least squares regression of $Y = \{\hat{\psi}(u_1), \dots, \hat{\psi}(u_k)\}^\top$ on the design matrix $Z \in \mathbb{R}^{k \times m}$, $Z_{ij} = u_i^j$. That is, $\hat{\beta} = (Z^\top Z)^{-1} Z^\top Y$, where $\hat{\beta} \in \mathbb{R}^m$. Finally let $\tilde{\psi}(u) = \sum_{j=1}^m \hat{\beta}_j u^j$, $u > 0$. We argue that the error achieved by $\tilde{\psi}$ is almost the parametric rate even though $\hat{\beta}$ can be estimated at a strictly lower rate.

Theorem 3.4. *Suppose $g(u; \xi) \equiv \xi^{-1}(e^{u\xi} - 1)$ and that f is supported on a compact interval. Assume 1 holds and $m = m_n = \log n / \log \log n$. Then $k^{-1} \sum_{i=1}^k \{\tilde{\psi}(u_i) - \psi(u_i)\}^2 = \mathcal{O}_p\{(\log n)^2/n\}$.*

Proof. Let β^0 be the true value $\beta_j^0 = \int \xi^{j-1} f(\xi) d\xi / j!$. Write $Y = Z\beta + \varepsilon$, where ε includes both the random error and the bias terms due to both the estimator and the truncation. The latter term is given in (3.9). By standard least squares results,

$$\begin{aligned} k^{-1} \mathbb{E} \sum_{i=1}^k \{\tilde{\psi}(u_i) - \psi(u_i)\}^2 &= k^{-1} \mathbb{E} \left\{ \varepsilon^\top Z (Z^\top Z)^{-1} Z^\top \varepsilon \right\} \\ &= k^{-1} \text{trace} \left\{ Z (Z^\top Z)^{-1} Z^\top \mathbb{E}(\varepsilon \varepsilon^\top) \right\}. \end{aligned}$$

Since $Z(Z^\top Z)^{-1}Z^\top$ is a projection matrix on a m -dimensional space, the RHS is bounded by the largest eigenvalue of $\mathbb{E}(\varepsilon \varepsilon^\top)$ times m/k . This has three components (variance and two biases) and hence

$$k^{-1} \mathbb{E} \sum_{i=1}^k \{\tilde{\psi}(u_i) - \psi(u_i)\}^2 = \mathcal{O} \left[\frac{m}{k} \left\{ \frac{k}{n} + k \frac{\log n}{n} + k \left(\frac{b^m}{m!} \right)^2 \right\} \right].$$

The factor k before the last two terms is due to the norm of the unit vector in \mathbb{R}^k , and, the last term is by (3.9). The theorem follows by taking $m = \log n / \log \log n$.

A more general result can be based on an assumption like the following.

Assumptions 2. For some c, d and each ε there are $h_{\varepsilon,1}, \dots, h_{\varepsilon, M(\varepsilon)}$ such that

$$\sup_{\xi} \min_{\gamma} \max_{c < u < d} \left| g(u; \xi) - \sum_{j=1}^{M(\varepsilon)} \gamma_j h_j(u) \right| < \varepsilon.$$

Note that clearly the assumption ensures the existence of $\gamma(\cdot)$ such that $\max_{c < u < d} |g(u; \xi) - \sum_{j=1}^{M(\varepsilon)} \gamma_j(\xi) h_j(u)| < \varepsilon$, but then there are also $\beta_j = \int \gamma_j(\xi) f(\xi) d\xi$, $j = 1, \dots, M(\varepsilon)$, such that $\max_{c < u < d} |\psi(u) - \sum_{j=1}^{M(\varepsilon)} \beta_j h_j(u)| < \varepsilon$.

The following theorem can be proved similarly to Theorem 3.4:

Theorem 3.5. *Suppose Assumptions 1 and 2 hold. Let $\varepsilon_n = \operatorname{argmin}_\varepsilon \{M(\varepsilon) / (n + \varepsilon)\}$, and let $\tilde{\psi}$ be the least squares estimate of the regression of $\hat{\psi}$ on $h_{\varepsilon_n, 1}, \dots, h_{\varepsilon_n, M(\varepsilon_n)}$. Then $k^{-1} \sum_{i=1}^k \{\tilde{\psi}(u_i) - \psi(u_i)\}^2 = \mathcal{O}_p(\varepsilon_n)$.*

In practice, Theorems 3.4 and 3.5 may seem to be of limited use — a knowledge of the structure of the span of the individual utility functions is needed, and the regression is based on an identified efficient base, which may not be natural. For example, we used a polynomial base for the exponential utility function. The practical approach is a histogram or discrete approximation of f . Does such a procedure yield an effective estimator, an estimator which is both statistically speaking efficient, but at the same time easy to compute and can be used in off-the-shelf manner?

This is indeed the case. Let $\xi_1, \dots, \xi_{M(\varepsilon)}$ be reasonably spaced points in the support of f . With the notation introduced after Assumption 2, and by a similar argument, for a vector β on the simplex

$$\sup_u \left| \sum_{j=1}^{M(\varepsilon)} \beta_j g(u; \xi_j) - \sum_{j=1}^{M(\varepsilon)} \beta_j \sum_{l=1}^{M(\varepsilon)} \gamma_l(\xi_j) h_l(u) \right| \leq \varepsilon.$$

Hence, one can use the base function $g(\cdot; \xi_1), \dots, g(\cdot; \xi_{M(\varepsilon)})$ as well.

References

- Ait-Sahalia, Y. and Lo, A. (2000). Nonparametric risk-management and implied risk aversion. *J. Econometrics* **94**.
- Balabdaoui, F. and Wellner, J. A. (2007). Estimation of a k-monotone density: limit distribution theory and the spline connection. Manuscript.
- Chabi-Yo, F., Garcia, R. M. and Renault, R. (2008). State dependence can explain the risk aversion puzzle. *Rev. Finan. Stud.* **21**, 973-1011.
- Cochrane, J. H. (2005). *Asset Pricing (Revised)*. Princeton University Press, Princeton.
- Detlefsen, K., Härdle, W. K. and Moro, R. A. (2007). Empirical pricing kernels and investor preferences. SFB649 Discussion paper 2007-017, http://sfb649.wiwi.hu-berlin.de/fedc/discussionPapers_de.php.
- Feller, W. (1966). *An Introduction to Probability Theory and its Applications, Vol. II*. Wiley, New-York.
- Friedman, M. and Savage, L. P. (1948). The utility analysis of choices involving risk. *J. Polit. Economy* **56**, 279-304.
- Gallant, A. R. and Hong, H. (2007). A statistical inquiry into the plausibility of Epstein-Zin-Weil Utility. *J. Finan. Econom.* **5**, 523-559.
- Gilbert, P. B. (2000). Large sample theory of maximum likelihood estimates in semiparametric biased sampling models. *Ann. Statist.* **28**, 151-194.
- Gilbert, P. B., Lele, S. R. and Vardi, Y. (1999). Maximum likelihood estimation in semiparametric selection bias models with application to AIDS vaccine trials. *Biometrika* **86**, 27-43.

Gill, R. I.
in b
Jewell, N
Manski,
base
Mantel, I
Rosenber
Silverman
Vardi, Y.
Walras, M
Departme
E-mail: y
CASE - C
Humbold
E-mail: h

- Gill, R. D., Vardi, Y. and Wellner, J. A. (1988). Large sample theory of empirical distributions in biased sampling models. *Ann. Statist.* **16**, 1069-1112.
- Jewell, N. P. (1982). Mixtures of exponential distributions. *Ann. Statist.* **10**, 479-482.
- Manski, C. F. and Lerman, S. R. (1977). The estimation of choice probabilities from choice based samples. *Econometrica* **45**, 1977-1988.
- Mantel, N. (1973). Synthetic retrospective studies and related topics. *Biometrics* **29**, 479-486.
- Rosenberg, J. and Engle, R. (2002). Empirical pricing kernels. *J. Finan. Econom.* **64**, 341-372.
- Silverman, B., (1986). *Density Estimation*. Chapman and Hall, London.
- Vardi, Y. (1985). Empirical distributions in selection bias models. *Ann. Statist.* **13**, 178-203.
- Walras, M.-E. L. (1874). *Éléments d'économie politique pure, ou théorie de la richesse sociale*.

Department of Statistics, The Hebrew University of Jerusalem 91905, Jerusalem, Israel.

E-mail: yaacov.ritov@gmail.com

CASE - Center for Applied Statistics and Economics, Institute for Statistics and Econometrics, Humboldt-Universität zu, 10178 Berlin, Germany.

E-mail: haerdle@wiwi.hu-berlin.de.

(Received February 2008; accepted February 2009)

The Bayesian Additive Classification Tree Applied to Credit Risk Modelling

Junni L. Zhang¹, Wolfgang K. Härdle²

¹Department of Business Statistics and Econometrics, Guanghua School of Management, Peking University, Beijing 100871, P. R. China; email: zjn@gsm.pku.edu.cn.

²Center for Applied Statistics and Economics, Wirtschaftswissenschaftliche Fakultät, Humboldt-Universität zu Berlin, Spandauer Straße 1, 10178, Berlin, Germany; email: haerdle@wiwi.hu-berlin.de.

Abstract: We propose a new nonlinear classification method based on a Bayesian “sum-of-trees” model, the Bayesian Additive Classification Tree (BACT), which extends the Bayesian Additive Regression Tree (BART) method into the classification context. Like BART, the BACT is a Bayesian nonparametric additive model specified by a prior and a likelihood in which the additive components are trees, and it is fitted by an iterative MCMC algorithm. Each of the trees learns a different part of the underlying function relating the dependent variable to the input variables, but the sum of the trees offers a flexible and robust model. Through several benchmark examples, we show that the BACT has excellent performance. We apply the BACT technique to classify whether firms would be insolvent. This practical example is very important for banks to construct their risk profile and operate successfully. We use the German Creditreform database and classify the solvency status of German firms based on financial statement information. We show that the BACT outperforms the logit model, CART and

the Support Vector Machine in identifying insolvent firms.

Key words and phrases: Classification and Regression Tree, Financial Ratio, Misclassification Rate, Accuracy Ratio

JEL-Codes: C14, C11, C45, C01

1 Introduction

Classification techniques have been popularly used in many fields. Standard classification tools include linear and quadratic discriminant analysis and the logistic model. The support vector machine (SVM) (Vapnik, 1995, 1997) recently arises as an important nonlinear classification tool. It maps the input space nonlinearly into a high dimensional feature space, and tries to find linear separating hyperplanes for the classes in the feature space, penalizing the distances of misclassified cases to the hyperplanes. The SVM has been widely and successfully applied to classification problems in many domains and often shown to have excellent performance compared to other classification methods.

Decision trees compose an important category of nonlinear classification methods. Ever since the introduction of the classification and regression tree (CART) by Breiman et al. (1984), it has attracted strong interest from researchers and practitioners. Figure 1 shows an example of a classification tree, where the root node (t_1) contains all training observations, and the training data are recursively partitioned by values of the input variables (x 's) until reaching the leaf (terminal) nodes (t_3 , t_4 , t_6 and t_7) where the classification decision (for y) is made for all observations contained therein. For regression problems in which the dependent variable is continuous, a predicted value for the dependent variable would be assigned for all observations contained in each leaf node.

Traditional search methods for CART models use locally greedy algorithms to find the partitions. The Bayesian approaches for CART models (Chipman et al., 1998; Denison et al., 1998; Wu et al., 2007) specify a formal prior distribution for trees and other parameters and use Markov Chain Monte Carlo methods to sample them from the posterior distribution.

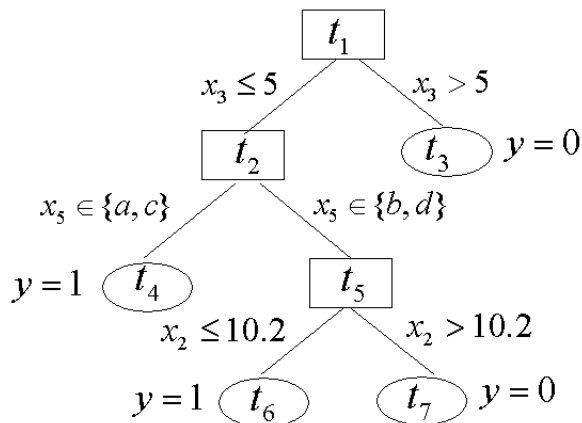


Figure 1: Example of a classification tree.

Chipman et al. (2006) proposed the Bayesian Additive Regression Tree (BART), in which the mean of a continuous dependent variable is approximated by a sum of trees rather than a single tree. This “sum-of-trees” model is defined by a prior and a likelihood, and fitted by iterative MCMC algorithm. Each individual tree explains a different portion of the underlying mean function, but the sum of these trees turns out to be a flexible and adaptive model. Chipman et al. (2006) showed that BART outperforms several competitive models, including LASSO (Efron et al., 2004), gradient boosting (Friedman, 2001), random forests (Breiman, 2001), and neural networks with one layer of hidden units. We will extend BART into the classification context, and therefore term the resulting classification technique as the Bayesian Additive Classification Tree (BACT).

To investigate the differences among the logit model, SVM, CART and BACT, we plot in Figure 2 the contours of these models trained to classify the solvency status of German firms using the German Creditreform database based on only two variables — the ratio of operating income to total assets (x_3 in Figure 2) and the ratio of accounts payable to

total sales (x_{24} in Figure 2). Details of this application will be discussed in Section 4. The contours for the logit model are linear, thus making it inflexible for complex applications. The SVM finds flexible smooth curves in the input space (linear hyperplanes in the feature space) that can separate the classes. The CART is based on a single tree which recursively partitions the observations by the input variables, and hence the contours are piecewise linear. The BACT is based on the sum of many trees, so the contours are not constrained to be piecewise linear as in CART; although these contours are not as smooth as in SVM, they are quite flexible in explaining complex structure.

The rest of this paper is organized as follows. Section 2 will describe the BACT in detail. Section 3 will use several benchmark examples from the UCI Machine Learning Repository to compare the performance of the BACT with the logit model and the SVM. Section 4 will discuss our application to classification of solvency status of Germany firms using the German Creditreform database. Section 5 then concludes.

2 The Bayesian Additive Classification Tree (BACT)

2.1 The Model

Consider a binary classification problem in which an dependent variable $Y \in \{1, 0\}$ needs to be predicted based on a set of input variables $\mathbf{x} = (x_1, \dots, x_p)^\top$. The majority of classification models assume that there is a latent continuous variable Y^* that determines

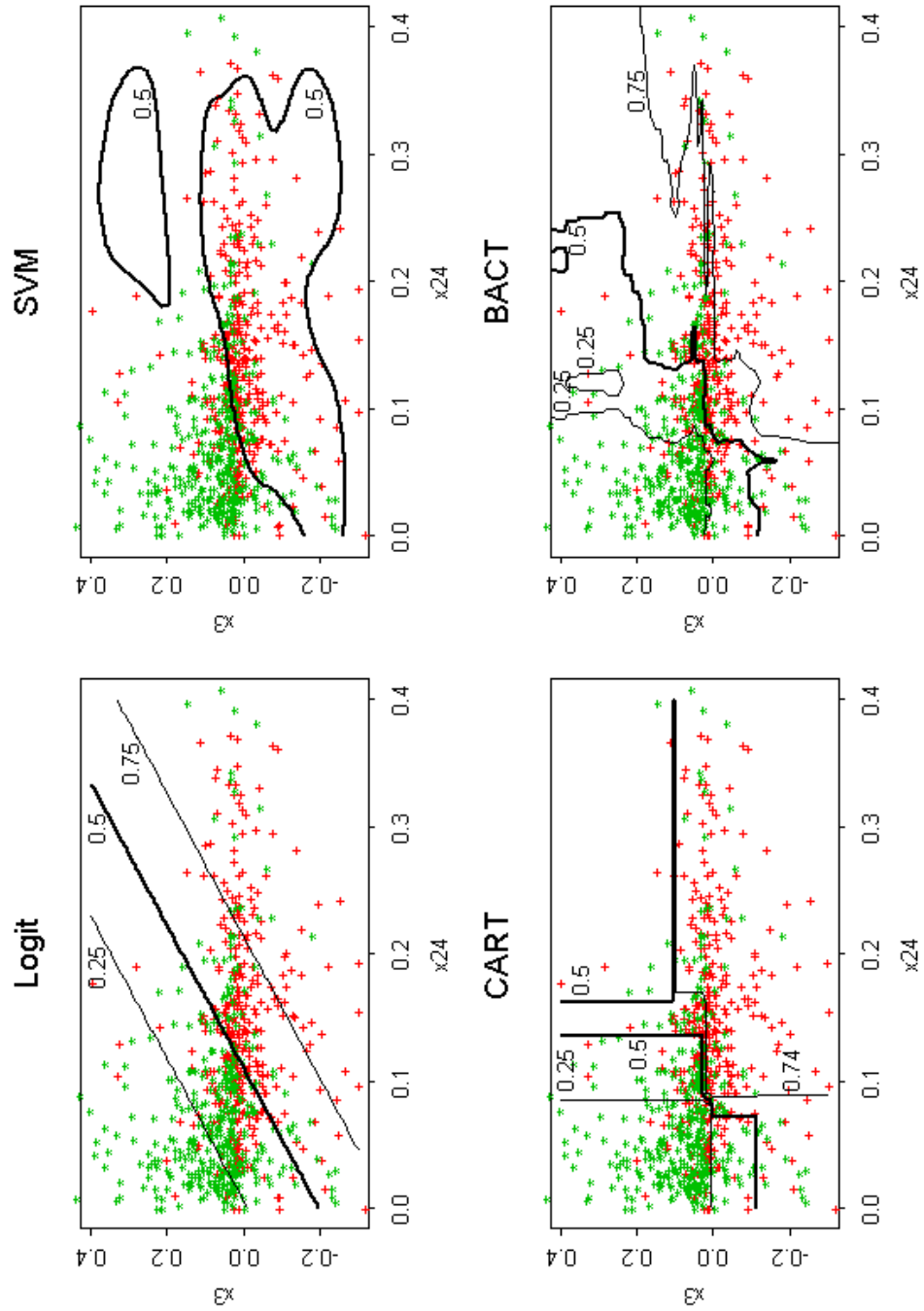


Figure 2: The contour plots for the logit model, SVM, CART, BACT. The pluses and stars represent insolvent firms and solvent firms respectively. The numbers by the contours indicate the probabilities of insolvency.

the value of Y as follows

$$\begin{cases} Y = 1 & \text{if } Y^* \geq 0 \\ Y = 0 & \text{if } Y^* < 0 \end{cases} \quad (1)$$

In the context of generalized linear models (GLM), the relationship of Y^* and \mathbf{x} is

$$Y^* = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \varepsilon,$$

where the distribution of ε determines the link function, e.g. logit or probit. The generalized additive models (GAM, Hastie and Tibshirani (1990)) replace each linear term in the GLM by a more generalized functional form and relate Y^* to \mathbf{x} by

$$Y^* = \beta_0 + f_1(x_1) + \cdots + f_p(x_p) + \varepsilon,$$

where each f_j is an unspecified smooth function.

Following the idea of the BART in Chipman et al. (2006), we assume that Y^* is related to \mathbf{x} through an additive model, where each additive component is a tree based on all input variables (rather than a flexible function based on a single input variable as in GAM). In order to formally introduce the model, we first introduce some notation. Let m denote the number of trees to be used. For $j = 1, \dots, m$, let T_j denote the j 'th tree with a set of partition rules based on the input variables, and let L_j denote the number of leaf nodes in T_j ; for $l = 1, \dots, L_j$, let μ_{jl} denote the (continuous) predicted value associated with the l 'th leaf node in T_j , and let $M_j = \{\mu_{j1}, \mu_{j2}, \dots, \mu_{jL_j}\}$. For a given value of \mathbf{x} , let $g(\mathbf{x}, T_j, M_j)$ denote the predicted value associated with the leaf node that an observation with input variables being \mathbf{x} would land in based on the partition rules for T_j . Thus Y^* is formally modelled as

$$Y^* = g(\mathbf{x}; T_1, M_1) + g(\mathbf{x}; T_2, M_2) + \cdots + g(\mathbf{x}; T_m, M_m) + \varepsilon, \quad (2)$$

and we further assume that $\varepsilon \sim N(0, 1)$, using a probit-like link.

2.2 Prior Specification

In order to make inferences from the model given by (1) and (2) in a Bayesian way, we need to specify a joint prior distribution for the unknown tree structures and leaf nodes parameters. We assume a priori that the tree structures and the leaf node parameters have independent distributions, so the full prior distribution can be written as

$$p\{(T_1, M_1), (T_2, M_2), \dots, (T_m, M_m)\} = \prod_{j=1}^m p(T_j) \prod_{j=1}^m \prod_{l=1}^{L_j} p(\mu_{jl}).$$

We further assume that every tree follows the same prior distribution, and every μ_{jl} follows the same prior distribution. So the task of prior specification is reduced to specifying the prior distribution for a single tree T and that for a single μ_{jl} parameter.

For a single tree T , we need to specify the prior distributions for its partition rules, including whether to further split a node or leave it as a leaf node, and if a further split is needed, which input variable and what values to be used for that split. We use the prior distribution for a single tree T as in Chipman et al. (2006). The prior probability of splitting any node n in tree T is

$$p_{split}(n, T) \propto \alpha(1 + d_n)^{-\beta},$$

where d_n is the depth of node n in tree T (the depth of node n is the length of the path from the root node to node n ; e.g., in Figure 1, the node t_1 has depth 0, and the nodes t_2 and t_3 have depth 1). α and β here are positive hyperparameters, hence the deeper a node is, the smaller probability there is to further split it, or the larger probability that this node becomes a leaf node. It turns out that the performance of BACT is not very sensitive to the

Table 1: Prior distribution on number of terminal nodes based on different values of α and β .

	Setting 1	Setting 2	Setting 3
α	0.5	0.95	0.95
β	2	2	0.1
prior probability of trees with 1 terminal node	0.5	0.05	0.05
prior probability of trees with 2 terminal nodes	0.383	0.552	0.012
prior probability of trees with 3 terminal nodes	0.098	0.275	0.004
prior probability of trees with 4 terminal nodes	0.017	0.092	0.002
prior probability of trees with ≥ 5 terminal nodes	0.003	0.031	0.932

choice of *alpha* and *beta*. We tried three different settings listed in Table 1 where a priori the trees range from small size to large size, and the resulting performance was quite similar. So we just pick $\alpha = .95$ and $\beta = 2$ as in Chipman et al. (2006). If a node needs to be split, the prior for the associated splitting rules assigns equal probability to each available input variable and equal probability on each available rule given the variable.

The prior distribution of μ_{jl} is taken to be a conjugate normal distribution $\mu_{jl} \sim N(0, \sigma_\mu^2)$ (conjugate because ε in (2) follows a normal distribution). From (2), we can see that the expected value of Y^* is equal to the sum of m different μ_{jl} parameters (recall that $g(\mathbf{x}, T_j, M_j)$ is the μ_{jl} parameter associated with the leaf node that an observation with input variables being \mathbf{x} would land in based on the partition rules for T_j); because of the a priori independence of μ_{jl} 's, the prior distribution for the expected value of Y^* is $N(0, m\sigma_\mu^2)$. Combining this with (1), it can be inferred that a priori each observation has probability 0.5 belonging to class 1 and probability 0.5 belonging to class 0.

To specify σ_μ^2 , we use the following procedure. We first estimate the range of Y^* (to be explained soon), and then choose σ_μ^2 such that there is at least 95% prior probability that the

expected value of Y^* is in the estimated range. Let the training data be $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, where N is the number of observations in the training data. We first randomly sample y_i^* for each observation i in the training data from truncated standard normal distributions such that the relationship in (1) holds between y_i^* and the observed y_i . Suppose that the sampled values are $\mathbf{y}^{*(0)} = \{y_i^{*(0)}\}_{i=1}^N$, and denote the minimum and maximum values of $y_i^{*(0)}$ as $\min(\mathbf{y}^{*(0)})$ and $\max(\mathbf{y}^{*(0)})$ respectively. Then $[\min(\mathbf{y}^{*(0)}), \max(\mathbf{y}^{*(0)})]$ is a very rough estimate of the range of Y^* . We choose an initial $\sigma_\mu^{2(0)}$ such that there is at least 95% prior probability that the expected value of Y^* is in this interval, i.e., $[-2\sqrt{m}\sigma_\mu^{2(0)}, 2\sqrt{m}\sigma_\mu^{2(0)}]$ covers $[\min(\mathbf{y}^{*(0)}), \max(\mathbf{y}^{*(0)})]$ and therefore $\sigma_\mu^{2(0)} = \max\{-\min(\mathbf{y}^{*(0)})/2\sqrt{m}, \max(\mathbf{y}^{*(0)})/2\sqrt{m}\}$. We then run the Markov Chain Monte Carlo (MCMC) algorithm to be described in Section 2.3 to generate posterior samples of y_i^* , and suppose that we obtain one posterior draw of $\mathbf{y}^{*(1)} = \{y_i^{*(1)}\}_{i=1}^N$ after dropping the first B_1 posterior draws used to reach convergence. We assume this set of y_i^* can be used to estimate reasonably the range of the true underlying Y^* , and choose the value of σ_μ^2 for further analysis such that there is at least 95% prior probability that the expected value of Y^* is in the interval $[\min(\mathbf{y}^{*(1)}), \max(\mathbf{y}^{*(1)})]$, i.e., $\sigma_\mu^2 = \max\{-\min(\mathbf{y}^{*(1)})/2\sqrt{m}, \max(\mathbf{y}^{*(1)})/2\sqrt{m}\}$.

2.3 Generation of Posterior Samples and Inference

We use the data augmentation method (Tanner and Wong, 1987) by treating $\mathbf{y}^* = \{y_i^*\}_{i=1}^N$ as missing data, and then use the Gibbs sampler to generate samples from the posterior distribution $p\{(T_1, M_1), (T_2, M_2), \dots, (T_m, M_m), \mathbf{y}^* | \mathcal{D}\}$.

Let $T_{(j)}$ denote the $m - 1$ trees other than T_j , and let $M_{(j)}$ denote the parameters

associated with the leaf nodes in $T_{(j)}$. The Gibbs sampler composes of drawing m successive draws of (T_j, M_j) for $j = 1, \dots, m$ from $p\{(T_j, M_j)|T_{(j)}, M_{(j)}, \mathbf{y}^*, \mathcal{D}\}$ followed by draw of \mathbf{y}^* from $p\{\mathbf{y}^*|(T_1, M_1), (T_2, M_2), \dots, (T_m, M_m), \mathcal{D}\}$. The draws of (T_j, M_j) can be generated similar to Chipman et al. (2006). Let $\hat{y}_i^* = \sum_{j=1}^m g(\mathbf{x}_i; T_j, M_j)$ denote the fitted value for observation i from the m trees. Then y_i^* ($i = 1, \dots, N$) can be independently generated from truncated normal distributions:

$$\begin{cases} y_i^* \sim N(\hat{y}_i^*, 1) \text{ and } y_i^* \geq 0 & \text{if } y_i = 1 \\ y_i^* \sim N(\hat{y}_i^*, 1) \text{ and } y_i^* < 0 & \text{if } y_i = 0 \end{cases}$$

After σ_μ^2 has been chosen according to the procedure described in Section 2.2, we can drop the first B_2 posterior draws used to reach convergence, and use subsequent S posterior draws for inference. Denote these S posterior draws as $\{(T_1^{(s)}, M_1^{(s)}), \dots, (T_m^{(s)}, M_m^{(s)})\}_{s=1}^S$. Given the s 'th draw, the probability that an observation with input variables \mathbf{x} belongs to class 1 is $\Phi\left\{\sum_{j=1}^m g(\mathbf{x}, T_j^{(s)}, M_j^{(s)})\right\}$, where Φ is the cumulative distribution function of standard normal distribution. Therefore, the posterior average probability that an observation with input variables \mathbf{x} belongs to class 1 can be estimated as

$$\frac{1}{S} \sum_{s=1}^S \Phi\left\{\sum_{j=1}^m g(\mathbf{x}, T_j^{(s)}, M_j^{(s)})\right\}. \quad (3)$$

We can use (3) to classify observations in training data or other data: if the probability calculated from (3) is larger than 0.5, then the observation is classified into class 1; otherwise it is classified into class 0.

Table 2: For five benchmark data sets from the UCI Machine Learning Repository, the number of cases, the number of variables, and the average misclassification rates for the test data using the logit model, the SVM and the BACT.

Data Set	# Cases	# Variables	Logit	SVM	BACT
breast cancer	683	9	3.8%	2.8%	3.3%
ionosphere	351	34	12.8%	4.5%	7.2%
diabetes	768	8	21.8%	25.2%	24.8%
sonar	208	60	29.8%	19.4%	17.2%
German credit	1000	30	23.6%	27.3%	23.6%

3 Benchmark Examples

To compare the performance of the BACT with the logit model and SVM (in which radial basis function is used as the kernel, and the parameters are chosen by cross-validation), we use five data sets for binary classification from the UCI Machine Learning Repository (Asuncion and Newman, 2007): breast cancer, ionosphere, diabetes, sonar, and German credit. Columns 2-3 in Table 2 summarize the number of cases and the number of variables for these data sets. Throughout the rest of the paper, in the BACT method, we fix $m = 200$, $B_1 = 500$, $B_2 = 1000$ and $S = 1000$.

We partition each data set randomly into 80% of training data and 20% of test data. The training data is used to fit the models, and misclassification rate on the test data is calculated. Such procedure is repeated for 20 times, and columns 4-6 in Table 2 report the average misclassification rates on the test data using the logit model, the SVM and the BACT. We can see that the BACT has comparable performance with the SVM, and has no worse performance than the logit model except for the “diabetes” data set.

4 Classification of Solvency Status of German Firms

We use the German Creditreform database, which contains financial statement information on 20,000 solvent and 1,000 insolvent firms in Germany and spans the period from 1996 to 2002. Information on the insolvent firms were collected two years prior to insolvency. Chen et al. (2007); Härdle et al. (2008) applied SVM to classify the solvency status of German firms, with the former using the German Creditreform database. We will preprocess the data set in the same way as Chen et al. (2007) do, and compare the results of our BACT with those of the logit model, CART and SVM.

Following Chen et al. (2007), we clean the data of firms whose characteristics are very different from the others. We first eliminate firms within industries with small percentage in the industry composition and are left with 949 insolvent firms and 16583 solvent firms in four main industries — Construction, Manufacturing, Wholesale & Retail Trade and Real Estate. We then exclude those firms whose asset size is less than 10^5 EUR or greater than 10^8 EUR, because the credit quality of small firms often depends as much on the finances of a key individual as on the firm itself and largest firms rarely go bankrupt in Germany. We further exclude the solvent firms in 1996 due to lack of insolvent firms in that year. We also eliminate firms with zero value for some variables used as denominators in calculating financial ratios to be used in classification. Several apparent outliers are then deleted and we end up with a data set with 783 insolvent firms and 9,575 solvent firms (due to slightly different ways of deleting outliers, our remaining solvent firms differ a little from the 9,583 solvent firms in Chen et al. (2007)).

We adopt the same set of financial variables to be used for classification as in Chen et al.

(2007) and list them in Table 3. The five number summary of these financial variables are listed in Table 4 for insolvent firms and solvent firms separately. In order to avoid sensitivity to outliers in applying the SVM, Chen et al. (2007) truncated each financial variable to be between its 5% quantile and 95% quantile. The BACT, however, only uses the ordering of values of the input variables in the partition rules, so there is no need to do such truncation.

We use the data from 1997 to 1999 to train the model, and use the data from 2000 to 2002 to test the resulting model. The training set contains 387 insolvent firms and 3535 solvent firms, and the test set contains 396 insolvent firms and 6040 solvent firms. Because the density of insolvent firms is rather low, we need to oversample the insolvent firms in order for the models to pick up the patterns predictive of insolvency (e.g., Berry and Linoff (2000), chap. 5). This is done through the bootstrap technique (Efron and Tibshirani, 1993; Sobehart et al., 2001). For each bootstrap sample, a training subset is constructed as follows. We use all 387 insolvent firms in the training set and randomly sample 387 solvent firms from the training set. This subset of 774 firm with 50% being insolvent is then used to train the model. When training the CART model, the training subset is further randomly partitioned into two parts stratified by the solvency status of the firms. The first part comprises of 80% of the training subset and is used to grow the tree, and the second part comprises of the remaining 20% of the training subset and is used to prune the tree. Performance measures are then evaluated using all observations (396 insolvent firms and 6040 solvent firms) in the test set. The average performance measures over 30 bootstrap samples are then calculated. We can compare average performance measures across different models.

We consider two performance measures: Accuracy Ratio (AR) (Sobehart and Keenan,

Table 3: Definition of financial variables to be used for classification for the Creditreform data.

Var.	Definition
x1	Net Income/Total Assets
x2	Net Income/Total Sales
x3	Operating Income/Total Assets
x4	Operating Income/Total Sales
x5	Earnings before Interest and Tax/Total Assets
x6	Earnings Before Interest, Tax, Depreciation and Amortization/Total Assets
x7	Earnings before Interest and Tax/Total Sales
x8	Own Funds/Total Assets
x9	(Own Funds – Intangible Assets) /(Total Assets – Intangible Assets – Cash and Cash Equivalents – Lands and Buildings)
x10	Current Liabilities/Total Assets
x11	(Current Liabilities – Cash and Cash Equivalents)/Total Assets
x12	Total Liabilities/Total Assets
x13	Debt/Total Assets
x14	Earnings before Interest and Tax/Interest Expense
x15	Cash and Cash Equivalents/Total Assets
x16	Cash and Cash Equivalents/Current Liabilities
x17	(Cash and Cash Equivalents – Inventories)/Current Liabilities
x18	Current Assets/Current Liabilities
x19	(Current Assets – Current Liabilities)/Total Assets
x20	Current Liabilities/Total Liabilities
x21	Total Assets/Total Sales
x22	Inventories/Total Sales
x23	Accounts Receivable/Total Sales
x24	Accounts Payable/Total Sales
x25	$\log(\text{Total Assets})$
x26	Increase (Decrease) in Inventories/Inventories
x27	Increase (Decrease) in Liabilities/Total Liabilities
x28	Increase (Decrease) in Cash Flow/Cash and Cash Equivalents

Table 4: Five number summary (minimum, lower quartile, median, upper quartile, maximum) of the financial variables for insolvent firms and solvent firms.

Var.	Insolvent Firms					Solvent Firms				
	min	Q1	mdn.	Q3	max	min	Q1	mdn.	Q3	max
x1	-1.51	-0.02	0.00	0.02	1.13	-4.82	0.00	0.02	0.06	5.92
x2	-5.41	-0.02	0.00	0.01	6.10	-17.13	0.00	0.01	0.03	15.91
x3	-0.97	-0.04	0.00	0.03	1.14	-4.82	0.00	0.03	0.09	5.97
x4	-3.38	-0.02	0.00	0.02	10.15	-44.81	0.00	0.02	0.04	20.39
x5	-0.99	-0.01	0.02	0.05	1.15	-1.51	0.02	0.05	0.11	5.95
x6	-0.91	0.03	0.07	0.11	1.17	-1.46	0.06	0.11	0.18	5.95
x7	-3.55	-0.01	0.01	0.04	10.27	-39.63	0.01	0.02	0.05	14.53
x8	0.00	0.00	0.05	0.14	0.96	0.00	0.05	0.14	0.28	0.99
x9	-0.86	0.00	0.05	0.17	2.31	-2.68	0.05	0.16	0.37	49.18
x10	0.01	0.37	0.52	0.73	1.00	0.00	0.25	0.42	0.64	4.13
x11	-0.35	0.33	0.49	0.69	0.99	-0.86	0.17	0.36	0.58	4.12
x12	0.01	0.54	0.76	0.89	1.00	0.00	0.42	0.65	0.82	4.37
x13	0.00	0.09	0.21	0.37	0.91	0.00	0.02	0.15	0.33	0.98
x14	-17658.06	-0.56	1.05	1.92	433.40	-22796.04	0.86	2.16	6.55	516896.73
x15	0.00	0.00	0.02	0.06	0.44	0.00	0.01	0.03	0.11	0.90
x16	0.00	0.01	0.03	0.12	25.01	0.00	0.01	0.08	0.30	40.61
x17	0.01	0.43	0.68	0.97	57.44	0.00	0.59	0.94	1.58	238.37
x18	0.03	1.00	1.26	1.84	62.63	0.06	1.11	1.58	2.67	989.76
x19	-0.69	0.00	0.15	0.36	0.92	-3.45	0.06	0.25	0.47	0.98
x20	0.07	0.62	0.84	0.99	1.18	0.01	0.56	0.85	1.00	1.00
x21	0.07	0.40	0.61	0.94	97.26	0.02	0.32	0.48	0.74	828.76
x22	0.00	0.08	0.16	0.34	89.96	-0.14	0.05	0.11	0.21	451.09
x23	0.00	0.07	0.12	0.18	0.87	0.00	0.05	0.09	0.14	21.85
x24	0.00	0.09	0.14	0.19	43.96	0.00	0.04	0.07	0.11	61.29
x25	11.72	14.07	14.87	15.76	18.25	11.51	14.25	15.41	16.62	18.42
x26	-46.89	-0.09	0.00	0.26	2.83	-282.51	-0.01	0.00	0.06	145.12
x27	-12.75	-0.04	0.00	0.11	1.00	-28.91	-0.04	0.00	0.10	1.00
x28	-1283.20	-0.61	0.00	0.18	1.00	-2513.39	-0.27	0.00	0.26	1.75

2001; Engelmann et al., 2003) and misclassification rate. AR is calculated using the Cumulative Accuracy Profiles (CAP) (Sobehart and Keenan, 2001; Engelmann et al., 2003) curve. To obtain the CAP curve, the firms are first ordered by risk scores from riskiest to safest. For BACT and the Logit model, the risk score is simply the predicted probability of insolvency; for SVM, the risk score can be calculated as distance to the separating hyperplane. The higher the risk score is, the riskier the firm is. For a given fraction q of the total number of firms, the CAP curve is constructed by calculating the fraction $r(q)$ of the insolvent firms whose risk scores are equal to or larger than the minimum score at fraction q .

Figure 3 plots the CAP curve for the test set of the Creditreform data where the scoring model is the BACT model trained using one bootstrap training subset. In the ideal case, the insolvent firms will be assigned the highest risk scores, and therefore the CAP curve would be increasing linearly and then stay at one. For a random model without any discriminative power, the fraction q of all firms with the highest risk scores will contain fraction q of all insolvent firms, and therefore the corresponding CAP curve will be a straight line connecting the points $(0,0)$ and $(1,1)$. AR is defined as the ratio of the area between the CAP curve for a scoring model and that for the random model to the area between the CAP curve for the ideal case and that for the random model. The value of AR lies between zero and one, with zero indicating no discriminative power of the scoring model and one indicating perfect discriminative power. Mathematically, AR is defined as

$$AR \equiv \frac{\int_0^1 r_{model}(q) dq - \frac{1}{2}}{\int_0^1 r_{ideal}(q) dq - \frac{1}{2}}, \quad (4)$$

where $r_{model}(q)$ and $r_{ideal}(q)$ indicate $r(q)$ for the scoring model and the ideal case respectively, and the integrals can be approximated by $\frac{1}{N} \sum_{i=1}^N r(i/N)$ where N is the number of

observations in the test set.

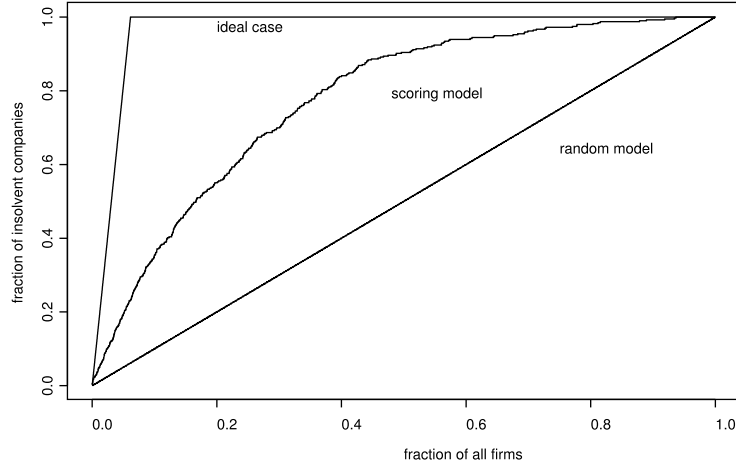


Figure 3: The CAP curve for the test set of the Creditreform data where the scoring model is the BACT model trained using one bootstrap training subset.

We also consider three types of misclassification rates: the overall misclassification rate, the type I misclassification rate and type II misclassification rate. Here type I misclassification refers to the case when the firm is in fact insolvent, but the model classifies the firm as solvent; whereas type II misclassification refers to the case when the firm is in fact solvent, but the model classifies the firm as insolvent. Financial institutions usually seek to keep either type of misclassification rate as low as possible (Sobehart et al., 2001).

Table 5 reports the average values of AR in (4) and the three types of misclassification rates for the Logit model, CART and BACT. Apparently, BACT outperforms the Logit model and CART in all aspects except for average Type I misclassification rate for which BACT is slightly worse than CART.

Table 5: The average values of AR and the three types of misclassification rates for the Logit model, CART and BACT.

Performance Measure	Logit	CART	BACT
AR	52.1%	58.7%	60.4%
Overall Misclassification Rate	30.2%	33.8%	26.6%
Type I Misclassification Rate	28.3%	27.2%	27.6%
Type II Misclassification Rate	30.3%	34.3%	26.5%

Rather than using all data from 2000 to 2002 as the test set, Chen et al. (2007) used a test subset for each bootstrap sample, which comprises of all insolvent firms and a random sample of the same number of solvent firms in the test set. They reported that the median AR value for 30 bootstrap samples was 60.5%, using $\frac{1}{10} \sum_{i=1}^{10} p(i/10)$ to approximate the integrals in calculating the AR value. The median overall misclassification rate was calculated as 28.2%. If we adopt the same procedure, BACT yields a median AR value of 66.5% and median overall classification rate as 27.2%. So BACT also outperforms SVM in identifying the insolvent firms.

5 Concluding Remarks

In this paper, we propose the Bayesian Additive Classification Tree as a general nonlinear classification method. We show that, based on the sum of many trees, the BACT can yield flexible class boundaries, and that it has excellent performance compared with the logit model, CART and SVM, as demonstrated through several benchmark examples and a real application to credit risk modelling.

Because the partitions in each tree depend only on the ordering of the values of the

input variables rather than the values themselves, the BACT is robust to extreme values in the input variables, and the results do not change with monotone transformation of any input variable. Hence little data processing is needed when using the BACT technique. Another thing to note is that although we only discuss binary classification in this paper, extension to multi-class classification is straightforward and left as future research.

Acknowledgement

This work was supported by Deutsche Forschungsgemeinschaft through the SFB 649 “Economic Risk”. Junni L. Zhang’s research was also sponsored by Chinese NSF grant 10401003 and USA NIH 1 R03 TW007197-01A2.

References

- Asuncion, A. and Newman, D. (2007), “UCI Machine Learning Repository,” [Http://www.ics.uci.edu/~mllearn/MLRepository.html](http://www.ics.uci.edu/~mllearn/MLRepository.html), University of California, Irvine, School of Information and Computer Sciences.
- Berry, M. and Linoff, G. (2000), *Mastering Data Mining*, John Wiley and Sons.
- Breiman, L. (2001), “Random forests,” *Machine Learning*, 5–32.
- Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984), *Classification and Regression Trees*, CRC Press.
- Chen, S., Härdle, W. K., and Moro, R. A. (2007), “Modeling Default Risk with Support Vector Machines,” To appear in *Journal of Quantitative Finance*.
- Chipman, H. A., George, E. I., and McCulloch, R. E. (1998), “Bayesian CART model search,” *Journal of the American Statistical Association*, 935–948.
- (2006), “BART: Bayesian Additive Regression Trees,” Technical Report, Graduate School of Business, University of Chicago.
- Denison, D., Mallick, B., and Smith, A. (1998), “A Bayesian CART Algorithm,” *Biometrika*, 363–377.
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004), “Least angle regression,” *Annals of Statistics*, 407–499.
- Efron, B. and Tibshirani, R. J. (1993), *An introduction to the bootstrap*, Chapman and Hall.

- Engelmann, B., Hayden, E., and Tasche, D. (2003), “Testing rating accuracy,” *Risk*, 82–86.
- Friedman, J. H. (2001), “Greedy function approximation: A gradient boosting machine,” *The Annals of Statistics*, 1189–1232.
- Härdle, W. K., Moro, R. A., and Schäfer, D. (2008), “Estimating Probabilities of Default With Support Vector Machines,” *to appear in Journal of Banking and Finance*.
- Hastie, T. J. and Tibshirani, R. J. (1990), *Generalized Additive Models*, Chapman and Hall.
- Sobehart, J. and Keenan, S. (2001), “Measuring default risk accurately,” *Risk*.
- Sobehart, J., Keenan, S., and Stein, R. (2001), “Benchmarking Quantitative Default Risk Models: A Validation Methodology,” *Algo Research Quarterly*.
- Tanner, M. A. and Wong, W. H. (1987), “The calculation of posterior distributions by data augmentation (with discussion),” *Journal of American Statistical Association*, 528–550.
- Vapnik, V. (1995), *The Nature of Statistical Learning Theory*, Springer, New York, NY.
- (1997), *Statistical Learning Theory*, Wiley, New York, NY.
- Wu, Y., Tjelmeland, H., and West, M. (2007), “Bayesian CART: prior specification and posterior simulation,” *Journal of Computational and Graphical Statistics*, in press.

Forecasting Volatility with Support Vector Machine-Based GARCH Model

SHIYI CHEN,^{1*} WOLFGANG K. HÄRDLE² AND
KIHO JEONG³

¹ *China Center for Economic Studies, School of Economics, Fudan University, Shanghai, China*

² *Center for Applied Statistics and Economics, Humboldt University, Berlin, Germany*

³ *School of Economics and Trade, Kyungpook National University, Daegu, Republic of Korea*

ABSTRACT

Recently, support vector machine (SVM), a novel artificial neural network (ANN), has been successfully used for financial forecasting. This paper deals with the application of SVM in volatility forecasting under the GARCH framework, the performance of which is compared with simple moving average, standard GARCH, nonlinear EGARCH and traditional ANN-GARCH models by using two evaluation measures and robust Diebold–Mariano tests. The real data used in this study are daily GBP exchange rates and NYSE composite index. Empirical results from both simulation and real data reveal that, under a recursive forecasting scheme, SVM-GARCH models significantly outperform the competing models in most situations of one-period-ahead volatility forecasting, which confirms the theoretical advantage of SVM. The standard GARCH model also performs well in the case of normality and large sample size, while EGARCH model is good at forecasting volatility under the high skewed distribution. The sensitivity analysis to choose SVM parameters and cross-validation to determine the stopping point of the recurrent SVM procedure are also examined in this study. Copyright © 2009 John Wiley & Sons, Ltd.

KEY WORDS (recurrent) support vector machine; GARCH model; volatility forecasting; Diebold–Mariano test

INTRODUCTION

Volatility is important in financial markets since it is a key variable in portfolio optimization, securities valuation and risk management. Much attention of academics and practitioners has been focused on modeling and forecasting volatility in the last few decades (see Franses and McAleer, 2002, and Poon and Granger, 2003, for a comprehensive review). So far in the literature, the predominant model of the past is the GARCH model by Bollerslev (1986), who generalizes the seminal idea on

*Correspondence to: Shiyi Chen, China Center for Economic Studies, School of Economics, Fudan University, Guoquan Road 600, Shanghai, China 200433. E-mail: shiyichen@fudan.edu.cn

ARCH by Engle (1982), and its various extensions; see Li *et al.* (2002) for recent surveys of the models. The GARCH family models, together with the simplest historical price model prevalent in the pre-GARCH era¹ and stochastic volatility model studied a decade later than GARCH development,² comprise one of the two broad categories of methods widely used in volatility forecasting, the so-called time series volatility model; another is the market determined option implied volatility model.³ This paper limits itself mainly to the analysis within the GARCH framework.

The popularity of the GARCH model is due to its ability to capture volatility persistence or clustering, supported by many studies (Akgriray, 1989; Bollerslev *et al.*, 1992; West and Cho, 1995; Andersen and Bollerslev, 1998; Marcucci, 2005). However, some empirical studies report that the GARCH model provides poor forecasting performance (Jorion, 1995, 1996; Brailsford and Faff, 1996; Figlewski, 1997; McMillan *et al.*, 2000; Choudhry and Wu, 2008). To improve the forecasting ability of the GARCH model, some alternative approaches have been advocated by innovating the model specification and estimation,⁴ by using different evaluation metrics and definitions of realized volatility,⁵ or by enriching the informational content of the model.⁶

As for GARCH model specification and estimation, for example, many financial returns are skewed distributed and nonlinearly dependent such that the linear GARCH model cannot cope with them and therefore forecast of symmetric GARCH model would be biased (Pagan and Schwert, 1990; Bollerslev *et al.*, 1992). To deal with this problem the regime-switching (RS) volatility model is proposed to detect nonlinear behavior in the variance by various tests for asymmetry or threshold

¹This includes simple moving average method, exponential smoothing method, random walk model, ARMA model, exponentially weighted moving average (EWMA) method and its current extension of Riskmetrics™ model, etc.

²The stochastic volatility (SV) model has an additional innovative term in the volatility dynamics (Taylor, 1986). For a detailed discussion on the SV model and its relation to the GARCH class models, see the survey articles by Ghysels *et al.* (1996) and Chib *et al.* (2002), among others.

³The time series volatility model is based on historical price information only, while the option implied volatility (IV) model uses market traded option information alone or in addition to historical price sets to forecast volatility. Many studies examine the relative performance of the IV model to forecasting volatility (Day and Lewis, 1992; Lamoureux and Lastrapes, 1993; Pong *et al.*, 2004; Dotsis *et al.*, 2007; Becker *et al.*, 2009; Neely, 2009). This paper limits itself mainly to analysis within the GARCH framework.

⁴Except for the introduction below, other relatively sophisticated GARCH models and estimations include the multivariate GARCH model (Bauwens *et al.*, 2006; Rosenow, 2008), outlier-corrected GARCH model (Park, 2002; Zhang and King, 2005; Ané *et al.*, 2008), Markov chain Monte Carlo (MCMC) sampling techniques to estimate the GARCH model (Gerlach and Tuyl, 2006), other semiparametric or nonparametric specification and estimation such as genetic algorithm, wavelet smoother, kernel density etc. (Franke *et al.*, 2004; Lux and Schornstein, 2005; Renò, 2006; Chen *et al.*, 2008; Feng and McNeil, 2008; Corradi *et al.*, 2009) and combination forecasts from competing approaches (Hu and Tsoukalas, 1999; Dunis and Huang, 2002).

⁵Many studies find that the relative accuracy of various models is also highly sensitive to the measures used to evaluate them (Taylor, 1999; Brooks and Persaud, 2003). Most comparisons are based on the average figure of mean absolute error (MAE) and mean square error (MSE) etc. Diebold and Mariano (1995) and West (1996) show how standard errors for MAE and MSE are derived taking into account serial correlation in the forecast errors for statistical inference. Lehar *et al.* (2002) applies value-at-risk (VaR)-oriented evaluation measures to compare the out-of-sample performance. In addition to the symmetric measures of MAE and MSE, Balaban (2004) also uses asymmetric evaluation criteria such as mean mixed error statistics to compare the forecasting performance, penalizing under/over-predictions of volatility more heavily. Recent research has also suggested that this relative failure of GARCH models arises not from a failure of the model but a failure to specify correctly the true volatility measure against which forecasting performance is measured. It is argued that the standard approach of using *ex post* daily squared returns as the measure of true volatility includes a large noisy component. An alternative measure for true volatility has therefore been suggested based on the cumulative squared returns from intra-day data, also referred to as realized, or integrated volatility (Andersen and Bollerslev, 1998; Andersen *et al.*, 2003; Meddahi, 2003; McMillan and Speight, 2004; Galbraith and Kisinbay, 2005; Ghysels *et al.*, 2006).

⁶In many instances, the researchers find the inclusion of implied volatility or trade volume as an exogenous variable in the framework of the GARCH model to be beneficial (Brooks, 1998; Fleming, 1998; Blair *et al.*, 2001; Koopman *et al.*, 2005; Gospodinov *et al.*, 2006; Becker *et al.*, 2007).

nonlinearity (Franses and Dijk, 2000). The first class of RS volatility model assumes that the regime can be determined by an observable variable, including the nonlinear exponential GARCH (EGARCH) model of Nelson (1991), threshold GJR-GARCH model of Glosten *et al.* (1992) and quadratic GARCH model of Engle *et al.* (1993) and Sentana (1995). The second class of RS model for volatility implements GARCH with a Hamilton (1989) type framework that assumes the regime is the realization of a hidden Markov chain, such as (double) Markov switching GARCH model of Gray (1996), Klaassen (2002) and Chen *et al.* (2008).

Both the linear and nonlinear GARCH model described above are parametric and normally estimated jointly by maximum likelihood estimation (MLE). That is, they make specific assumptions about the functional form of the data generation process and the distribution of error terms that is necessary for MLE. Such parametric models are easy to estimate and readily interpretable, but these advantages may come at a cost. Perhaps nonparametric models are better representations of the underlying data generation process. Instead of specifying a particular functional form and making *a priori* distributional assumption, the nonparametric model will search for the best fit over a large set of alternative functional forms. Thus, in the literature, many nonlinear nonparametric GARCH models are developed and still developing fast, among which the artificial neural network (ANN) is extensively used. This paper focuses on one of the neural network algorithms, the support vector machine (SVM), and investigates its forecasting ability of volatility as compared with the simplest moving average method, standard linear GARCH model, nonlinear EGARCH model and traditional recurrent ANN-based nonlinear GARCH model. The moving average method is chosen as the benchmark because some studies find that it provides more accurate forecasts than GARCH models (Dimson and Marsh, 1990; Tse and Tung, 1992; Figlewski, 1997). Among the number of nonlinear parametric GARCH models the EGARCH model is also the most commonly used (Cao and Tsay, 1992; Cumby *et al.*, 1993; Heynen and Kat, 1994; Chong *et al.*, 1999; Hu and Tsoukalas, 1999; Gokcan, 2000; Balaban, 2004).

In recent years, ANN has been successfully used for forecasting financial time series; for recent work, see Fernandez-Rodriguez *et al.* (2000), Qi and Wu (2003), and Pantelidaki and Bunn (2005). The studies in favor of ANN-based GARCH model as opposed to parametric GARCH model in forecasting conditional volatility include Donaldson and Kamstra (1997), Schittenkopf *et al.* (2000), Taylor (2000), Dunis and Huang (2002), Hamid and Iqbal (2004), Ferland and Lalancette (2006), Tseng *et al.* (2008). However, the traditional ANN algorithm also suffers from its own weaknesses such as the need for many controlling parameters, difficulty in obtaining a global solution and the danger of over-fitting (Tay and Cao, 2001). Thus, SVM that can obtain a unique global solution by solving a quadratic programming is developed by Vapnik and his co-workers (1995, 1997). Naturally, SVM also keeps the advantages of conventional ANN such as the flexibility in approximating any nonlinear function arbitrarily well, without *a priori* assumptions about the properties of the data and without the requirement of large sample size that MLE-based parametric GARCH models have. Unlike traditional ANN implementing the empirical risk minimization (ERM) principle, the most particular principle of SVM is to implement the structural risk minimization (SRM), which seeks to achieve a balance between the training error and generalization error, leading, theoretically, to better forecasting performance than traditional ANN (Gunn, 1998; Haykin, 1999). Recently, SVM has gained popularity in predicting financial variables owing to such attractive features (Cao and Tay, 2001; Härdle *et al.*, 2005, 2007; Chen *et al.*, 2009). Pérez-Cruz *et al.* (2003) also propose an SVM-based GARCH (1, 1) model and shows that it provides better volatility forecasts than the standard GARCH model. However, they use the feedforward SVM procedure, which has the same structure as the autoregressive (AR) process and has poor ability

to model a long-time memory. Inspired by the merit of recurrent ANN (Kuan and Liu, 1995; Dunis and Huang, 2002; Bekiros and Georgoutsos, 2008), in this paper we propose a recurrent SVM procedure which can model the ARMA process and apply it to forecast the conditional variance equation of the GARCH model in real data analysis.

The forecasting accuracy of the recurrent SVM-based GARCH model in one-period-ahead volatility forecasting is compared with the competing models in terms of two evaluation metrics of mean absolute error (MAE) and directional accuracy (DA). The statistical hypothesis of equal forecasting accuracy between pairwise models is also investigated by using the Diebold and Mariano (1995) test, calculated according to the Newey–West procedure (Newey and West, 1987). The Diebold and Mariano (DM) test is one of the most important contributions to the study of out-of-sample forecasting accuracy evaluation over the past two decades, and has been further generalized and extensively used in many studies since then (Corradi and Swanson, 2004; Awartani and Corradi, 2005; Preminger and Franck, 2007; Taylor, 2008; Groen *et al.*, 2009; Wong and Tu, 2009).

This paper is organized as follows. The next section briefly introduces the theory of SVM. The third section specifies the empirical model and forecasting scheme. The fourth section uses the Monte Carlo simulation to evaluate how the models perform under controlled conditions. The fifth section describes the GBP exchange rates and NYSE composite index data and discusses the volatility forecasting performance of all models for the real data. The paper concludes with the sixth section.

SUPPORT VECTOR MACHINE

The support vector machine (SVM) originates from Vapnik's statistical learning theory (Vapnik, 1995, 1997), which has the design of a feedforward network with an input layer, a single hidden layer of nonlinear units and an output layer, and formulates the regression problem as a quadratic programming (QP) problem. SVM estimates a function by nonlinearly mapping the input space into a high-dimensional hidden space and then running the linear regression in the output space. Thus, the linear regression in the output space corresponds to a nonlinear regression in the low-dimensional input space. The theory denotes that if the dimensions of feature space (or hidden space) are high enough, SVM may approximate any nonlinear mapping relations. As the name implies, the design of the SVM hinges upon the extraction of a subset of the training data that serves as support vectors, which represent a stable characteristic of the data.

Given a training dataset (\mathbf{x}_t, y_t) , where input vector $\mathbf{x}_t \in \mathbb{R}^p$ and output scalar $y_t \in \mathbb{R}^1$. Indeed, the desired response y , known as a 'teacher', represents the optimum action to be performed by the SVM. We aim at finding a sample regression function $f(\mathbf{x})$, or denoted by \hat{y} , as below to approximate the latent, unknown decision function $g(\mathbf{x})$:

$$f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b \quad (1)$$

where the superscript T is a transposing operator that should be differentiated from the sample size T of the time series used later in this paper. In equation (1), $\phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \dots, \phi_l(\mathbf{x})]^T$, $\mathbf{w} = [w_1, \dots, w_l]^T$. The $\phi(\mathbf{x})$ is known as the nonlinear transfer function which represents the features of the input space and projects the inputs into the feature space. The dimension of the feature space is l , which is directly related to the capacity of the SVM to approximate a smooth input–output mapping; the higher the dimension of the feature space, the more accurate the approximation will be. Parameter

\mathbf{w} denotes a set of linear weights connecting the feature space to the output space, and b is the threshold.

To get the function $f(\mathbf{x})$, the optimal \mathbf{w}^* and b^* have to be estimated from the data. First, we define a linear ε -insensitive loss function, L_ε , originally proposed by Vapnik (1995):

$$L_\varepsilon(\mathbf{x}, y, f(\mathbf{x})) = \begin{cases} |y - f(\mathbf{x})| - \varepsilon & \text{for } |y - f(\mathbf{x})| \geq \varepsilon \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

This function indicates the fact that it does not penalize errors below ε . The training points within the ε -tube have no loss and do not provide any information for decision. Therefore, these points do not appear in the decision function $f(\mathbf{x})$. Only those data points located on or outside the ε -tube will serve as the support vectors and are finally used to construct the $f(\mathbf{x})$. This property of sparseness algorithm results only from the ε -insensitive loss function and greatly simplifies the computation of SVM. The non-negative slack variables, ξ and ξ' (below or above the ε -tube, or denoted together by $\xi^{(i)}$; see Figure 1) are employed to describe this kind of ε -insensitive loss.

The derivation of SVM follows the principle of structural risk minimization (SRM) that is rooted in the Vapnik–Chervonenkis (VC) dimension theory (Haykin, 1999). Structural risk is the upper boundary of empirical loss, denoted by ε -insensitive loss function, plus the confidence interval (or called margin), which is constructed in equation (3). The primal constrained optimization problem of SVM is obtained below:

$$\min_{\mathbf{w} \in \mathbb{R}^T, \xi^{(i)} \in \mathbb{R}^{2T}, b \in \mathbb{R}} \mathbf{C}(\mathbf{w}, b, \xi_i, \xi'_i) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^T (\xi_i + \xi'_i) \quad (3)$$

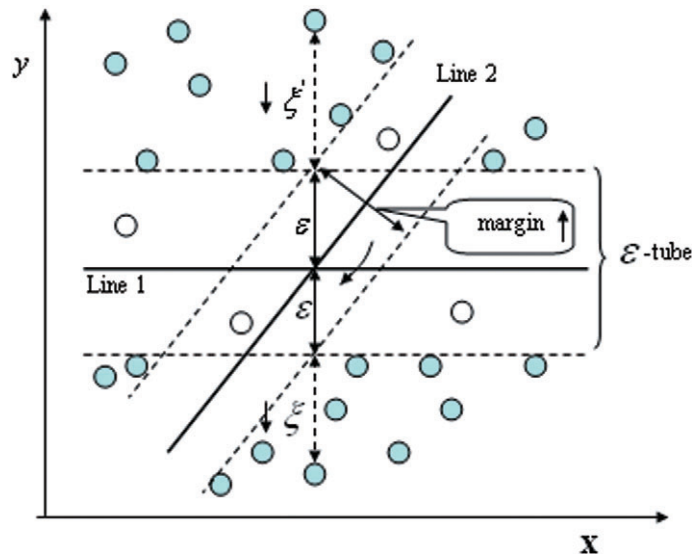


Figure 1. Principle of structural risk minimization (SRM) of SVM

such that

$$\mathbf{w}^T \phi(\mathbf{x}_t) + b - y_t \leq \varepsilon + \xi_t \tag{4}$$

$$y_t - \mathbf{w}^T \phi(\mathbf{x}_t) - b \leq \varepsilon + \xi'_t \tag{5}$$

$$\xi_t \geq 0, \xi'_t \geq 0, t = 1, 2, \dots, T \tag{6}$$

The formulation of the cost function $C(\cdot)$ in equation (3) is in perfect accord with the SRM principle, which is illustrated in Figure 1 (in which the dark circles are data points extracted as support vectors). In equation (3), the first term indicates the Euclidean norm of the weight vector $\mathbf{w}(\|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w})$ and measures the function flatness; to minimize it is equivalent to maximizing the separation margin ($2/\|\mathbf{w}\|$), that is, maximizing the generalization ability. The second term represents the empirical risk loss determined by the ε -insensitive loss function and is similar to the sum of residual squares in the objective function of ANN. Finally, SVM obtains the tradeoff between the two terms; as a result, it not only fits the historical data well but also forecasts the future data excellently. As shown in Figure 1, both regression lines 1 and 2 can classify the data points correctly and then minimize the empirical loss; however, the separation margin of the two lines are different, in which the regression line 1 has the larger margin. It is the special design of minimizing the structural risk that endows SVM with the excellent forecasting ability among all candidates. In addition, the convex quadratic programming and linear restrictions in the above primal problem ensure that SVM can always obtain the global unique optimal solution, which is different from the usual networks that easily get trapped in local minima. The penalty parameter $C > 0$ controls the penalizing extent on the sample points which lie outside ε -tube. Both ε and C , the free parameter of SVM, must be selected by the user.

The corresponding dual problem of the SVM can be derived from the primal problem by using the Karush–Kuhn–Tucker conditions as follows:

$$\min_{\alpha'_i \in \mathbb{R}^{2T}} \frac{1}{2} \sum_{s=1}^T \sum_{t=1}^T (\alpha'_s - \alpha_s)(\alpha'_t - \alpha_t) K(x_s \cdot x_t) + \varepsilon \sum_{t=1}^T (\alpha'_t + \alpha_t) - \sum_{t=1}^T y_t (\alpha'_t - \alpha_t) \tag{7}$$

such that

$$\sum_{t=1}^T (\alpha_t - \alpha'_t) = 0 \tag{8}$$

$$0 \leq \alpha_t, \alpha'_t \leq Cs, t = 1, 2, \dots, T \tag{9}$$

where α_t and α'_t (or $\alpha'_t{}^{(i)}$) are the Lagrange multipliers. The dual problem can be solved more easily than the primal problem (Scholkopf and Smola, 2001; Deng and Tian, 2004). Making use of any solution of α_t and α'_t , the optimal solutions of the primal problem can be calculated in which \mathbf{w}^* is unique and expressed as follows:

$$\mathbf{w}^* = \sum_{t=1}^T (\alpha'_t - \alpha_t) \phi(\mathbf{x}_t) \tag{10}$$

However, b^* is not unique and formulated in terms of different cases. If $i \in \{t | \alpha_t \in (0, C)\}$, then

$$b^* = y_i - \sum_{t=1}^T (\alpha'_t - \alpha_t) K(\mathbf{x}_t \cdot \mathbf{x}_i) + \varepsilon \quad (11)$$

If $j \in \{t | \alpha'_t \in (0, C)\}$, then

$$b^* = y_j - \sum_{t=1}^T (\alpha'_t - \alpha_t) K(\mathbf{x}_t \cdot \mathbf{x}_j) - \varepsilon \quad (12)$$

The cases of both $i, j \in \{t | \alpha_t^{(\cdot)} = 0\}$ and $i, j \in \{t | \alpha_t^{(\cdot)} = C\}$ rarely occur in reality.

Thus the regression decision function $f(\mathbf{x})$ will be computed by using \mathbf{w}^* and b^* in the following forms:

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^{*T} \phi(\mathbf{x}) + b^* \\ &= \sum_{t=1}^T (\alpha'_t - \alpha_t) \phi^T(\mathbf{x}_t) \phi(\mathbf{x}) + b^* \\ &= \sum_{t=1}^T (\alpha'_t - \alpha_t) K(\mathbf{x}_t, \mathbf{x}) + b^* \end{aligned} \quad (13)$$

where $K(\mathbf{x}_t, \mathbf{x}) = \phi^T(\mathbf{x}_t) \phi(\mathbf{x})$ is the inner-product kernel function. In fact, the SVM theory considers only the form of $K(\mathbf{x}_t, \mathbf{x})$ in the feature space without specifying explicitly $\phi(\mathbf{x})$ and without computing all corresponding inner products. Therefore, the kernel function greatly reduces the computational complexity of high-dimensional hidden space and becomes the crucial part of SVM. The function which satisfies the Mercer theorem can be chosen as the SVM kernel. No analytical method is currently available to determine the most suitable kernel for a particular dataset. This paper experiments with three different kernels to investigate the effect of a kernel type in Monte Carlo simulation:

$$\text{Linear: } K(\mathbf{x}_t, \mathbf{x}) = \mathbf{x}_t^T \mathbf{x} \quad (14)$$

$$\text{Polynomial: } K(\mathbf{x}_t, \mathbf{x}) = (\mathbf{x}_t^T \mathbf{x} + 1)^d \quad (15)$$

$$\text{Gaussian: } K(\mathbf{x}_t, \mathbf{x}) = \exp\left(\frac{-\|\mathbf{x} - \mathbf{x}_t\|^2}{2\sigma^2}\right) \quad (16)$$

where d and σ^2 are the parameters for the polynomial and Gaussian kernel. Before implementation of the SVM, the appropriate values of the coefficients ε , C , d and σ^2 must be determined in advance through cross-validation. The sensitivity analysis of the parameters and the kernel type will be illustrated by using the simulated data below ('Monte Carlo Simulation').

EMPIRICAL MODELING

In this study, the forecasts are obtained first by applying the Monte Carlo Simulation, following the suggestions in Andersen and Bollerslev (1998) and Clements and Smith (1999, 2001). The main motivation for conducting a simulation experiment is that, since the true volatility is known, the candidate volatility measures can be compared with certainty. We then fit each of the models to the daily returns on the GBP exchange rate and NYSE stock indexes and forecast their respective volatility. The empirical modeling and forecasting scheme described below are employed for both simulation and real data.

Model specification

In this paper the real data we analyze are the daily financial returns, y_t , converted from the corresponding price or index, I_t , using continuous compounding transformation as

$$y_t = 100 \times (\log I_{t+1} - \log I_t) \quad (17)$$

Empirical findings suggest that GARCH is a more parsimonious model than ARCH, and GARCH (1, 1) specification is sufficient to model the variance changing over long sample periods and has become the most popular structure when capturing financial volatility (Akgiray, 1989; Franses and Dijk, 1996; Brooks, 1998; Gokcan, 2000; Andersson, 2001; Brooks and Persaud, 2003; Poon and Granger, 2003; Gerlach and Tuyl, 2006). As such, throughout the paper, the analysis is restricted to the case of the GARCH (1, 1) process for the second conditional variance function and the AR(1)⁷ process for the conditional mean equation, for the sake of candidate comparison under the same conditions.

Thus the linear standard GARCH (1, 1) model is specified as follows:

$$y_t = c + \phi_1 y_{t-1} + u_t \quad u_t \sim N(0, h_t) \quad (18a)$$

$$h_t = \kappa + \delta_1 h_{t-1} + \alpha_1 u_{t-1}^2 \quad (18b)$$

where c , ϕ_1 , κ , δ_1 and α_1 are constant parameters. Such restrictions on the parameters that κ , δ_1 and α_1 are non-negative and $\delta_1 + \alpha_1 < 1$ prevent negative variances (Bollerslev, 1986).

All odd moments of u_t in the standard GARCH model equal zero, and hence u_t and y_t are symmetric time series. The nonlinear EGARCH (1, 1) model that is able to capture the asymmetry is similar to the linear GARCH model but the h_t process is given by

$$\log(h_t) = \kappa + \delta_1 \log(h_{t-1}) + \alpha_1 \left(\frac{|u_{t-1}|}{\sqrt{h_{t-1}}} - \sqrt{2/\pi} \right) + \beta_1 \frac{u_{t-1}}{\sqrt{h_{t-1}}} \quad (19)$$

where κ , δ_1 , α_1 and β_1 are the constant parameters. The EGARCH model is fundamentally different from the standard GARCH model in that the standardized innovation serves as the forcing variable for the conditional variance. Also, there are no restrictions on the parameters to ensure non-negativity

⁷Franses and Dijk (1996) also denote that the order of autoregression in the first conditional mean equation of the GARCH framework is usually 0 or small. Thus, the order 1 is specified for this study.

of the variances. The coefficient β_1 is introduced to capture the asymmetry. If $\beta_1 = 0$, a positive return shock has the same effect on h_t as the negative return shock of the same amount; if $\beta_1 < 0$, a positive return shock actually reduces h_t ; if $\beta_1 > 0$, then a positive return shock increases h_t . Previous studies have viewed this coefficient as typically negative, indicating that negative return shocks normally generate more volatility than positive return shocks, so generating the so-called leverage effect.

The conditional variance of u_t is given by $h_t = E_{t-1}u_t^2 = \hat{u}_{t-1}^2$. Roughly speaking, in a GARCH process the conditional variances can be modeled by an ARMA type process (Franses and Dijk, 1996). For instance, the ARMA process of the conditional variance of u_t in a linear GARCH model can be expressed as below (Hamilton, 1997; Enders, 2004):

$$u_t^2 = \kappa + (\delta_1 + \alpha_1)u_{t-1}^2 + w_t - \delta_1 w_{t-1} \quad (20)$$

where $w_t \equiv u_t^2 - \hat{u}_{t-1}^2 = u_t^2 - h_t$, which is white noisy error. Inspired by this, the nonparametric recurrent ANN and SVM based nonlinear GARCH (1, 1) model is specified as the following form:

$$y_t = f(y_{t-1}) + u_t \quad (21a)$$

$$u_t^2 = g(u_{t-1}^2, w_{t-1}) + w_t \quad (21b)$$

where $f(\cdot)$ and $g(\cdot)$ are nonlinear nonparametric function forms for conditional mean and variance equations, respectively. Note that equation (21b) is adopted for the analysis of real data because the actual volatility h_t is unobservable, while in the case of simulation the conditional variance equation is just specified as $h_t = f(h_{t-1}, u_{t-1}^2)$ due to h_t being known. Because of the way GARCH (1, 1) class models are constructed, the volatility is known at time $t - 1$. Thus the one-step-ahead forecast of volatility is readily available.

The moving average method uses weighted moving averages of past squared innovations to forecast volatility (Niemira and Klein, 1994). For simulated data, the moving average forecast for the next-day volatility, using the five most recent observations, is expressed as

$$\hat{u}_{t+1}^2 = \frac{1}{5} \sum_{j=t-4}^t u_j^2 \quad (22)$$

For real data, the moving average forecast for the next-day volatility is expressed as (Engle *et al.*, 1993)

$$\hat{u}_{t+1}^2 = \frac{1}{5} \sum_{j=t-4}^t (y_j - \bar{y}_{5,t})^2 \quad (23)$$

where

$$\bar{y}_{5,t} = \frac{1}{5} \sum_{j=t-4}^t y_j$$

The recurrent ANN used in this study is the feedback multilayer perceptrons (MLP) network with the addition of a global feedback connection from the output layer to its input space. We specify

this kind of recurrent back-propagation network with the following architecture: one nonlinear hidden layer with four neurons, each using a tan-sigmoid differentiable transfer function to generate the output, and one linear output layer with one neuron. As a training algorithm, the fast training Levenberg–Marquardt algorithm is chosen. The value of the learning rate parameter used in the training process is set to be 0.05. These specifications and choices are standard in the neural network literature.

Recurrent SVM procedure

As Haykin (1999) said, the standard SVM described above usually appears in the design of a simple network in which an input layer of source nodes projects onto an output layer of computation node, but not vice versa (see Figure 2(a)). This process is known as feedforward SVM and could be easily employed to estimate such AR process as the first conditional mean function (21a), $y_t = f(y_{t-1}) + u_t$, and the second conditional variance function in the situation of simulation, $h_t = f(h_{t-1}, u_{t-1}^2)$. However, because the unobservable error term w_t is introduced into the GARCH model which indeed exhibits the nonlinear ARMA process, how to estimate the conditional volatility model (21b) for real data?

To estimate the nonlinear ARMA model, a feedback process of SVM with unobservable moving average part as inputs, not addressed before our application⁸, has to be described, which distinguishes itself from feedforward SVM in that it has at least one feedback loop (see Figure 2(b)). In this paper, we abuse terminology and refer to this process as ‘recurrent SVM’. The feedback loops involve the use of particular branches composed of *one-delay operator*, z^{-1} , which result in nonlinear dynamical behavior and have a profound impact on the learning capability of SVM. Thus the recurrent SVM will capture more dynamic characteristics of y_t than does feedforward SVM.

To overcome the problem that the series of error term w_t is unavailable, we employ the model residuals as estimates of the errors in an iterative way, which is similar to the way that the linear ARMA model is iteratively estimated by MLE (Box *et al.*, 1994; Hamilton, 1997). Likewise, the

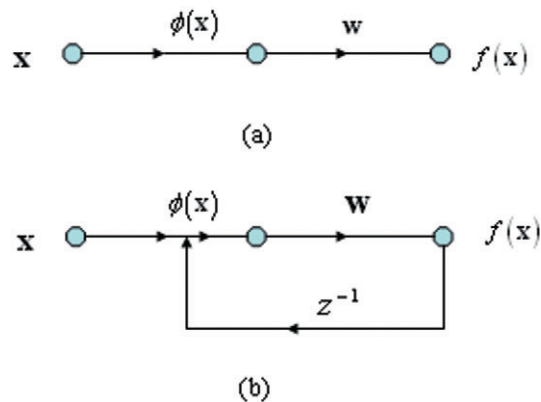


Figure 2. Signal-flow graphs of feedforward and recurrent SVM. (a) Signal-flow graph of a feedforward SVM. (b) Signal-flow graph of a single-loop recurrent SVM

⁸Suykens and Vandewalle (2000) proposed the algorithm of recurrent least squares SVM. The difference between the two recurrent SVM algorithms is their sparseness solutions.

error term is initially set to be its expectation: zero. The empirical procedure of the recurrent SVM executed during the training phase is described as follows. The letter i indicates the iterative epoch and t denotes the period:

- Step 1: Set $i = 1$ and start with all residuals at zero: $w_t^{(1)} = 0$.
- Step 2: Run an SVM procedure to get the decision function $f^{(i)}$ to the points $\{x_t, y_t\} = \{u_{t-1}^2, u_t^2\}$ with all inputs $x_t = \{u_{t-1}^2, w_{t-1}^{(i)}\}$.
- Step 3: Compute the new residuals $w_t^{(i+1)} = u_t^2 - f^{(i)}$.
- Step 4: Terminate the computational process when the stopping criterion is satisfied; otherwise, set $i = i + 1$ and go back to Step 2.

Note that the first iterative epoch is in fact a feedforward SVM process and results in an AR (1) model and that the following epochs provide results of the ARMA (1, 1) model, being estimated by the recurrent SVM.

In general, the procedure cannot be shown to converge, and there are no well-defined criteria for stopping its operation. Rather, some reasonable criteria can be found, although with its own practical drawback, which may be used to terminate the computational process.

To formulate such a criterion, it is logical to think in terms of the properties of the estimated residual series. After sufficiently long iterative steps, the autocorrelation displayed behind the residuals during the first AR epoch should disappear, and the information in the residual behavior has been completely adopted and the final residual series should be white noisy. Accordingly, we may suggest a sensible convergence criterion for the recurrent SVM procedure as follows:

The recurrent SVM procedure is considered to have converged when the corresponding residuals become white noisy, or has no autocorrelation.

To quantify the measurement of white noise, we use the formal hypothesis test, the Ljung–Box–Pierce Q -test, to investigate a departure from randomness based on the ACF of the residuals. Under the null hypothesis of no autocorrelation in residuals, the Q -test statistic is asymptotically distributed as chi-square. In fact, we just check the actual p -values (exact level of significance) of the Q -test of lag 1. It is reasonable to think there is no higher-order autocorrelation if there is no one-order autocorrelation in residuals. Only if the p -values of the Q -test for five consecutive epochs are simultaneously higher than 0.1 is the iterative computational process stopped. To overcome the drawback of this convergence criterion, we use cross-validation to avoid the possible over-fitting problem; see ‘Real data analysis’ below for the iterative process in detail.

Forecasting scheme

To illustrate the forecasting scheme, the SVM-GARCH model is also exemplified. First, estimate the conditional mean equation (21a) by using the feedforward SVM in the full sample period $T(1, 2, \dots, T)$ to obtain residuals, u_1, u_2, \dots, u_T . Then, recursively run the SVM-GARCH (1, 1) model for squared residuals thus obtained to forecast the one-period-ahead volatility. The recursive forecasting scheme is employed with an updating sample window; the estimating and forecasting process is carried out recursively by updating the sample with one observation each time, rerunning the SVM approach and recalculating the model parameters and corresponding forecasts. Here, the SVM approach to estimate the conditional volatility is feedforward for simulation and recurrent, as described in the above subsection, for real data. The first training sample is $u_1^2, u_2^2, \dots, u_{T_1}^2$ ($T_1 < T$). The observations of $T - T_1$ are retained as a forecasting or test sample.

Therefore, we can estimate and forecast the SVM-based conditional volatility equation for $n = T - T_1$ times. We set $n = 60$ for both simulation and real data in this study. Thus, 60 one-period-ahead forecast volatilities, $\hat{u}_{T-59}^2, \hat{u}_{T-58}^2, \dots, \hat{u}_{T-1}^2, \hat{u}_T^2$, will be acquired for out-of-sample forecasting evaluation.

Evaluation measures and pairwise comparison of competing models

We evaluate the forecasting performance using two standard statistical criteria: mean absolute forecast error (MAE) and directional accuracy (DA), expressed as follows (Brooks, 1998; Moosa, 2000):

$$\text{MAE} = \frac{1}{n} \sum_{t=T_1}^{T-1} |u_{t+1}^2 - \hat{u}_{t+1}^2| \quad (24)$$

$$\text{DA}(\%) = \frac{100}{n} \sum_{t=T_1}^{T-1} a_t \quad (25)$$

where

$$a_t = \begin{cases} 1 & (u_{t+1}^2 - u_t^2)(\hat{u}_{t+1}^2 - \hat{u}_t^2) \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

MAE measures the average magnitude of forecasting error which disproportionately weights large forecast errors more gently relative to MSE; and DA measures the correctness of the turning point forecasts, which gives a rough indication of the average direction of the forecast volatility.

The fundamental problem with the evaluation of volatility forecasts of real data is that volatility is unobservable and so actual values with which to compare the forecasts do not exist. Therefore, researchers are necessarily required to make an auxiliary assumption about how the actual *ex post* volatility is calculated. In this paper, we use the square of the return minus its mean value as the surrogate of actual volatility against which MAE and DA can be calculated. This approach is similar to the standard one, squared returns, because the mean of returns is usually close to zero. The proxy of actual volatility in real data is expressed as

$$u_t^2 = (y_t - \bar{y})^2 \quad (26)$$

where y_t is returns and \bar{y} is mean of returns. This proxy has been used in many recent papers, such as Pagan and Schwert (1990), Day and Lewis (1992), Chan *et al.* (1995), West and Cho (1995), Chong *et al.* (1999), Brooks (2001) and Brooks and Persaud (2003).

To test for equal forecasting accuracy of two competing models, we use the two-sided DM test statistic proposed by Diebold and Mariano (1995) for the difference of MAE loss function. The null and alternative hypotheses in this case are

$$H_0: \text{MAE}_1 - \text{MAE}_0 = 0 \text{ versus } H_1: \text{MAE}_1 - \text{MAE}_0 \neq 0$$

where the subscript 0 denotes the benchmark model and 1 the competing model. The DM statistic in a robust form is then based on the following large sample statistic:

$$\text{DM} = \frac{1}{\sqrt{n}} \frac{1}{\sqrt{\hat{S}^2}} \sum_{t=T_1}^{T-1} (|u_{t+1}^2 - \hat{u}_{1,t+1}^2| - |u_{t+1}^2 - \hat{u}_{0,t+1}^2|) \sim N(0, 1) \quad (27)$$

where \hat{S}^2 denotes a heteroscedasticity and autocorrelation consistent (HAC) robust (co)variance matrix which is estimated according to the Newey–West procedure (Newey and West, 1987). We use Andrews' (1991) approximation rule to automatically select the number of lags for the HAC matrix. If n grows at a rate such that as $T \rightarrow \infty$, $n \rightarrow \infty$ and $n/T_1 \rightarrow 0$, then the DM statistic converges in distribution to a standard normal.

MONTE CARLO SIMULATION

Data-generating process

In this section we investigate the forecasting performance of all candidates using artificial simulated data under controlled conditions. To generate the data, we first need to parameterize the GARCH (1, 1) model in equation (18) with the following settings $(c, \phi_1, \kappa, \delta_1, \alpha_1) = (0, 0.5, 0.0005, 0.8, 0.1)$ for medium persistence and a disturbance term u_t distributed first as Gaussian and then as a Student's t with five degrees of freedom (kurtosis = 5). The second distribution tries to model the skewness and excess of kurtosis that usually appears in real financial series. Using the same specified models, two artificial samples of size 500 and 1000 are created under a two-distributions assumption, giving a total of four situations. To limit the computational burden, each situation is replicated only 50 times. Then the multiple simulated y_t and h_t are 500×50 and 1000×50 element matrices for different distribution.

Parameter selection

The use of cross-validation is appealing particularly when we have to design a somewhat complex approach with good generalization as the goal. For example, here we may use cross-validation to determine the values of free parameters of SVM with the best performance. One series of 50 simulated returns and volatility of 1000 size and Student's t distribution, one of the four situations, is exemplified as below. The first training data, that is, the former 940 observations, are used to determine the appropriate values taken by the free parameters. The training data are further randomly partitioned into two disjoint subsets: estimating sample and validating sample (700 and 240 observations, respectively).

As shown above, two free parameters (ε and C) and two kernel coefficients (d and σ^2) have to be selected by users before running the SVM procedure. The motivation for using cross-validation here is to validate the model on a dataset different from the one used for parameter estimation. In this way we may use the training set to assess the performance of various values of parameters, and thereby choose the best one. The sensitivity investigation of SVM (represented by the generalization error, MAE) with respect to four parameters is illustrated in Figures 3 and 4 for conditional mean and variance estimation, respectively.

Figure 3 describes the sensitivity analysis for the conditional mean equation. Parameter C varies from a very small value of 0.0001 to infinity, with ε being fixed at 0.0001 and σ^2 0.4. Clearly, when

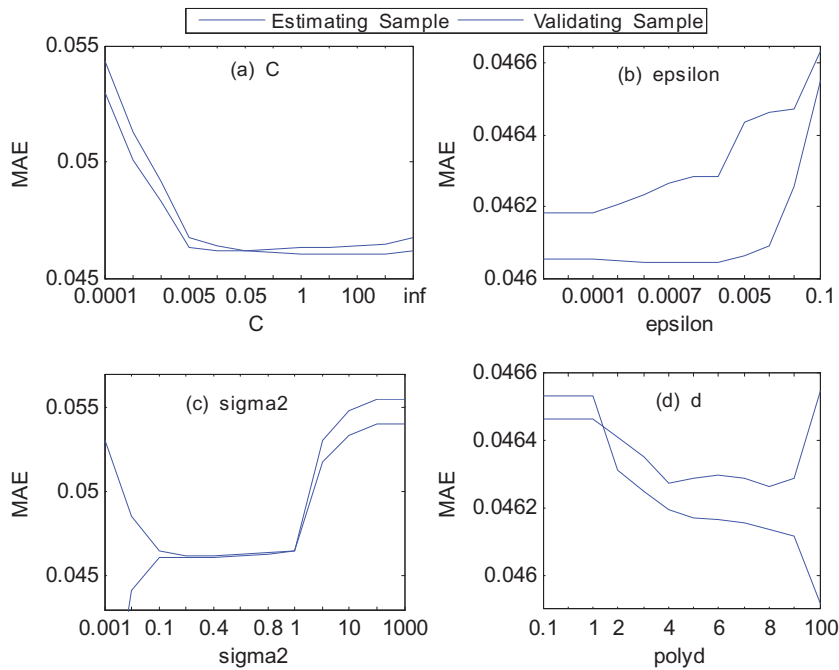


Figure 3. Sensitivity analysis of SVM in conditional mean estimation

$C = 0.05$, MAE of the validation sample obtains the lowest value, 0.046. Parameter ε takes values in the range [0.00001, 0.00005, 0.0001, 0.0003, 0.0005, 0.0007, 0.0009, 0.001, 0.005, 0.01, 0.05, 0.1], with $C = 0.05$ and $\sigma^2 = 0.4$. The values of ε to the left of the point = 0.0001 have no influence on the performance of SVM. Coefficient σ^2 varies from values of 0.001 to 1000, with C being 0.05 and 0.0001. Obviously, the value of $\sigma^2 = 0.4$ leads to the best validation performance. If we set $C = 0.05$ and 0.0001 and the polynomial kernel parameter $d = [0.1, 0.5, 1, 2, 3, 4, 5, 6, 7, 8, 10, 100]$, the validating MAE attains the minima when $d = 8$; after that, over-fitting the training set occurs. Note that the polynomial kernel with $d = 1$ is similar to the linear kernel. Thus, the appropriate parameters of SVM for the conditional mean returns are: $C = 0.05$, $\varepsilon = 0.0001$, $\sigma^2 = 0.4$ and $d = 8$.

Figure 4 describes the parameter selection process for conditional variance series. Similar to the return series, the MAE of both estimating and validating sample decreases as the values of C increase and become stable when C takes a value greater than 10; in contrast to C , as the values of ε increase, both MAE of SVM are considerably more stable before the point of $\varepsilon = 0.0001$ and increase slowly, and sharply after $\varepsilon = 0.001$. The value of $\sigma^2 = 0.01$ results in the best validation performance; namely, its MAE reaches the minimum value, about 0.000065. The values of d taken between 100 and 1000 have not much effect on the performance of SVM but after that range the over-fitting phenomenon becomes serious. Likewise, when one parameter is analyzed, the others are set to be fixed. Therefore, the correct parameters chosen for the conditional variance series are $C = 10$, $\varepsilon = 0.00005$, $\sigma^2 = 0.01$ and $d = 250$, respectively.

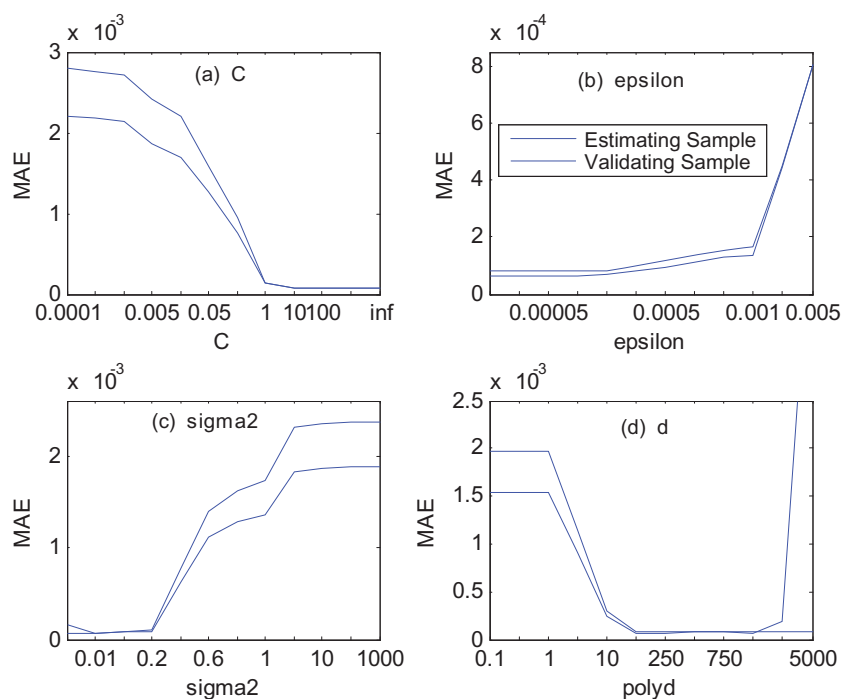


Figure 4. Sensitivity analysis of SVM in conditional variance estimation

Thus far we discuss the sensitivity investigation of parameters by using the simulated data with 1000 observations and t distribution. The parameter selection for the other three random samples is similar to this and not reported here to save space.

EFFECT OF KERNEL TYPE AND FORECASTING EVALUATION

There is still the possibility of over-fitting after training. Therefore, the generalization performance of the competing models is further measured and evaluated on the test set, which is different from the validation subset. For the simulated data, the forecasting sample is the last 60 observations. For each replication, the SVM-based GARCH (1, 1) model and the others are estimated, and the forecasting errors are calculated using the forecasting schemes described above. The results of out-of-sample one-period-ahead volatility forecasting measures for four situations are shown in Table I. The reported results are the mean values of 50 independent replications. Table II presents the p -values of Diebold-Mariano (DM) test for the MAE difference, which are defined as the significance levels at which the null hypothesis under investigation can be rejected. In calculating the DM statistic, the null hypothesis of equal forecasting ability is related to the four benchmark models: moving average, standard GARCH, EGARCH and traditional ANN models. We report the results of the DM test, say DM1, in the third and seventh columns for two simulated series, respectively, under the null hypothesis that the absolute forecast error produced by the moving average method is equal to those obtained using the other models. DM2, DM3 and DM4 are organized in the same manner and show the test results when the benchmark models are respectively the standard GARCH, EGARCH and recurrent ANN models. The DM tests in this study are investigated in a robust form, by simply

Table I. Diebold–Mariano test for the MAE difference on real data

Models	Sample Size = 500				Sample Size = 1000			
	Normality		Student's t		Normality		Student's t	
	MAE	DA	MAE	DA	MAE	DA	MAE	DA
Moving Average	0.0001276	44.07	0.0001747	59.32	0.0001198	54.24	0.0002130	40.68
Standard GARCH	0.0000972	76.27	0.0001765	55.93	0.0000488	79.66	0.0001083	59.32
EGARCH	0.0001312	67.80	0.0002075	64.41	0.0000730	57.63	0.0001864	74.58
ANN-GARCH	0.0001517	72.88	0.0002481	57.63	0.0000904	62.71	0.0001442	67.80
SVMI-GARCH	0.0000960	76.27	0.0001369	71.19	0.0000501	74.58	0.0000715	72.88
SVMp-GARCH	0.0000924	76.27	0.0001371	71.19	0.0000479	71.19	0.0000714	77.97
SVMg-GARCH	0.0000796	86.44	0.0001397	81.36	0.0000456	83.05	0.0000769	98.31

Note: SVMI, SVMp and SVMg represent the SVM with linear, polynomial and Gaussian kernel, respectively, for short.

scaling the numerator by a heteroscedasticity and autocorrelation consistent (HAC) (co)variance matrix calculated according to Newey–West procedures (Newey and West, 1987).

Table I firstly shows the effect of kernel functions on out-of-sample forecasting performance of SVM. The linear kernel behaves better in the sample with 500 sizes and t distribution based on DA measure. The polynomial kernel is the most suitable for forecasting the t -distributed 1000 sample size also based on DA. For all the other six cases, the Gaussian kernel looks promising, however, which is not a general conclusion but only true for the case we are studying. As a whole, three types of kernel-based SVM have a similar volatility forecasting performance and almost behave better than the benchmarks. Since no single kernel function dominates all volatility predictions, practitioners could try any kernel function. In the real data analysis later, for example, we only investigate the performance of the Gaussian kernel-based SVM-GARCH model.

Now, based on Table I, we revert to comparing the volatility forecasting ability among all competing models. In terms of the average ranking of MAE measures, the order of the forecasting ability of the different methods from highest to lowest is displayed in turn as follows: SVMp-GARCH, SVMg-GARCH, SVMI-GARCH⁹, standard GARCH, EGARCH, moving average and ANN-GARCH model. Concretely, in the situation of normal distribution, the standard GARCH model behaves not badly, which is ranked fourth (only inferior to three SVM models) in the 500 sizes and even ranked third (only defeated by Gaussian and polynomial SVM models) in the series of 1000 sizes. Even though the data satisfy the normality assumption that is required for MLE in the standard GARCH model, the SVM-GARCH models still outperform it in forecasting the magnitude of the volatility error. Nonlinear EGARCH and ANN-GARCH models perform worse than the linear GARCH model. In the situation of t distribution, the forecasting performance of the linear GARCH model grows poorer and the difference of MAE values between SVM-GARCH and standard MLE-GARCH models becomes larger than that under normality. Possibly this results from the fact that the normality assumption required for MLE is violated but it is not necessary for the SVM method. Not as expected, the asymmetric EGARCH model is weak in reducing the forecasting error even in the case of skewed distribution.

Based on the DA measures in Table I, on average, the Gaussian SVM-GARCH model ranks highest (for all four situations) in forecasting volatility directions, followed by polynomial and linear

⁹That is, corresponding to SVM-based GARCH models with polynomial, Gaussian and linear kernel function, respectively.

SVM-GARCH models, linear GARCH model, EGARCH model, ANN-GARCH model and moving average, in turn. In the situation of the normal distribution, the standard GARCH model behaves even better than forecasting error magnitude—ranked second for both the series of 500 sizes (only inferior to Gaussian but equal to linear and polynomial SVM models) and 1000 sizes (worse than Gaussian but better than the other two SVM type models). In the case of normality and large sample sizes, particularly favorable for MLE, the standard GARCH model still cannot defeat the Gaussian-based SVM-GARCH model. It is not surprising for EGARCH to behave badly in this case. As for the situation of t distribution, the linear GARCH model is ranked last for the 500 sizes (55.93%) and second last for the 1000 sizes (59.32%); while the asymmetric EGARCH model is good at forecasts of volatility turning points—ranked fourth for short series (only behind the three SVM models) and even third for long series (inferior to Gaussian and polynomial but better than the linear SVM-GARCH model). This time the ANN-GARCH model defeats the linear GARCH model. As for the linear GARCH model and moving average method, in the situation of 500 sizes and t distribution the standard GARCH model performs worse than the moving average, the simplest time series method, in terms of both MAE and DA measures. The conclusions described above are obtained on average based on 50 replications.

Table II displays the p -values of the DM test when the moving average method, standard GARCH, EGARCH and ANN models are compared with each of the other models considered in the study. We denote these tests DM1, DM2, DM3 and DM4, respectively. For instance, DM1 presents the test results for the simple moving average, where a p -value no greater than 0.05 indicates that the moving average method yields a higher forecast error (in terms of absolute error) relative to the competing model at 5% significance level, a p -value no smaller than 0.95 means that the moving average produces a lower forecast error at the 5% level, while a p -value between 0.05 and 0.95 implies that the benchmark and competing model have equivalent forecasting accuracy from the viewpoint of statistics. The same interpretation applies to the p -values reported for DM2–DM4.

Table II. Diebold–Mariano test for the MAE difference on Monte Carlo simulation

Distribution	Models	Sample size = 500				Sample size = 1000			
		DM1	DM2	DM3	DM4	DM1	DM2	DM3	DM4
Normality	Moving average		0.976	0.401	0.070		1.000	0.999	0.875
	Standard GARCH	0.024		0.001	0.000	0.000		0.001	0.000
	EGARCH	0.600	0.999		0.005	0.001	0.999		0.033
	ANN-GARCH	0.930	1.000	0.995		0.125	1.000	0.967	
	SVMl-GARCH	0.018	0.460	0.002	0.000	0.000	0.574	0.002	0.000
	SVMp-GARCH	0.023	0.413	0.004	0.000	0.000	0.420	0.003	0.000
Student's t	SVMg-GARCH	0.002	0.097	0.000	0.000	0.000	0.354	0.000	0.000
	Moving average		0.480	0.036	0.000		1.000	0.822	0.984
	Standard GARCH	0.520		0.054	0.003	0.000		0.000	0.001
	EGARCH	0.964	0.946		0.021	0.178	1.000		0.966
	ANN-GARCH	1.000	0.997	0.979		0.016	0.999	0.034	
	SVMl-GARCH	0.043	0.037	0.002	0.000	0.000	0.019	0.000	0.000
	SVMp-GARCH	0.056	0.043	0.001	0.000	0.000	0.025	0.000	0.000
	SVMg-GARCH	0.070	0.050	0.000	0.000	0.000	0.033	0.000	0.000

Note: DM1, DM2, DM3 and DM4 are the robust Diebold and Mariano (1995) test by using the Newey–West procedures (Newey and West, 1987) when the benchmark models are the moving average, linear GARCH model, EGARCH model and traditional ANN-GARCH model, respectively. For each test we consider the MAE loss functions.

Under the normal distribution, DM1 tests indicate that there is equivalent forecasting ability between moving average and EGARCH for short series, and between moving average and ANN-GARCH for long series. Such models as standard GARCH and the three SVM-GARCH all have higher volatility forecasting accuracy than moving average for both series at least at the 5% significance level. Moving average outperforms the ANN-GARCH model at the 10% level for a series of 500 size and EGARCH outperforms moving average at the 0.1% significance level for long series. According to DM2, three SVM type models have statistically equivalent forecasting ability to standard GARCH model for both series, with only one exception that the Gaussian SVM-GARCH model behaves better than the standard GARCH model at 10% significance level for short series. For both series, the standard GARCH model outperforms EGARCH and ANN-GARCH models at extremely low significance level. The DM3 statistic reveals that, for two series, three SVM-GARCH models perform better than the EGARCH model and EGARCH better than the ANN-GARCH model all at extremely significant levels. Finally, the ANN-GARCH model is found statistically and consistently inferior to the three SVM models for any series based on DM4 tests.

In the case of Student's t distribution, the out-of-sample performance of the standard GARCH model deteriorates. Now, according to DM2, the three SVM-GARCH models forecast volatility significantly better than the standard GARCH model at the 5% level for both series. The standard GARCH model cannot statistically defeat the moving average, either, for short series. However, both EGARCH and ANN-GARCH models are still statistically inferior to the standard GARCH model. In fact, according to DM1, DM3 and DM4, the three SVM-GARCH models all consistently outperform such benchmarks as moving average, EGARCH and ANN-GARCH models in forecasting volatility for any series. In terms of DM1, furthermore, the null hypothesis of equal forecasting accuracy between moving average and EGARCH cannot be rejected for a series of 1000 size rather 500 size. Moving average is significantly better than the ANN-GARCH model for short series, but the case is reversed for long series. In a series of 500 sizes, the ANN-GARCH model is significantly outperformed by the EGARCH model, while for the series of 1000 size the ANN type model statistically defeats the EGARCH model.

In summary, it appears that the three SVM-GARCH models do a better job of forecasting volatility than the moving average, standard GARCH, EGARCH and ANN-GARCH models in terms of MAE measures, which is statistically supported by the DM1, DM3, DM4 tests and DM2 in the case of t distribution. The DM2 test reveals that under the normal distribution the three SVM-GARCH models and standard GARCH model have similar volatility forecasting ability. Based on DA measures, the standard GARCH model too has a better ability in forecasting volatility turning points under normality and large sample sizes, while the asymmetric EGARCH model behaves better under the skewed t distribution. But both linear GARCH and nonlinear EGARCH cannot defeat all SVM-type models, at least the Gaussian-based SVM-GARCH model, in forecasting volatility directions.

REAL DATA ANALYSIS

In this section, we investigate the volatility forecasting performance of all candidates by using real data for two kinds of financial variables: GBP/USD exchange rates and NYSE average index.

Data description

The first dataset consists of the daily nominal bilateral exchange rates of British pounds (GBP) against the US dollar for the period January 5, 2004 to December 31, 2007. The data are obtained

from a database provided by Policy Analysis Computing and Information Facility in Commerce (PACIFIC) at the University of British Columbia, which contains the closing rates for a total of 81 currencies and commodities. The second dataset consists of the daily closing price of the New York Stock Exchange (NYSE) composite stock index for the period January 8, 2004 to December 31, 2007. The data are downloaded directly from the Market Information section of the NYSE web page.

It has been widely accepted that a variety of financial variables including foreign exchange rates and stock prices are integrated of order one. To avoid the issue of possible nonstationarity, both sets of raw real data are transformed into daily returns via equation (17), giving a returns series of 1001 observations and then a residual series is obtained from a fitted conditional mean equation of the GARCH class models. For the squared residuals of 1000 observations, the recursive estimating samples for the conditional volatility function are updated from the former 940 observations through the former 999 and then 60 numbers of one-period-ahead volatility forecasts are obtained, corresponding to an evaluation sample spanned from the 941st through the 1000th data points, that is, out-of-sample period of October 3, 2007 to December 31, 2007 for GBP and October 5, 2007 to December 31, 2007 for NYSE data.

The daily series for the log-levels and the returns of the GBP and NYSE are depicted in Figure 5. This figure shows that the returns series are mean-stationary, and exhibit the typical volatility clustering phenomenon with periods of unusually large volatility followed by periods of relative tranquility. Table III reports the summary of the descriptive statistics for the GBP and NYSE returns. Both series are typically characterized by excessive kurtosis and asymmetry. The Bera and Jarque (1981) tests all strongly reject the normality hypothesis. For GBP series, the Ljung–Box Q(6) statistic

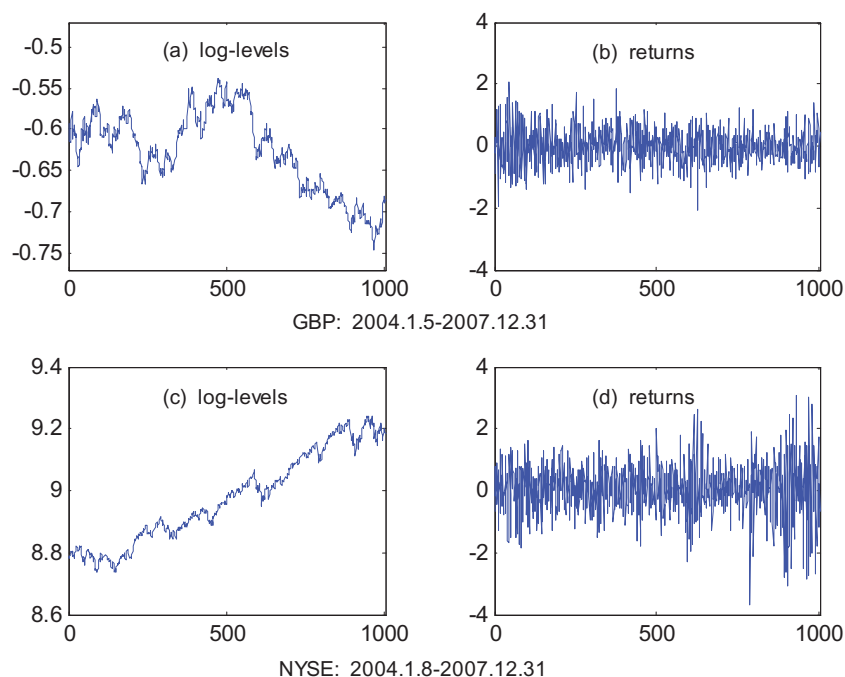


Figure 5. Log levels and returns of GBP exchange rates and NYSE stock index

Table III. Descriptive statistics for daily financial returns

Returns	GBP		NYSE	
	Statistics	<i>p</i> -value	Statistics	<i>p</i> -value
Mean	-0.0092		0.0393	
Variance	0.2827		0.6197	
Skewness	0.1206		-0.3489	
Kurtosis	3.7130		4.9343	
Normality	23.1860	0.00001	174.7200	0.00000
<i>Q</i> (6)	3.0313	0.80490	12.7100	0.04788
<i>Q</i> (6)*	31.6390	0.00002	150.2400	0.00000
ARCH(6)	28.9280	0.00006	101.8400	0.00000

Notes: Normality is the Bera-Jarque (1981) normality test; *Q*(6) is the Ljung-Box *Q* test at 6 order for raw returns; *Q*(6)* is LB *Q* test for squared returns; ARCH(6) is Engle's (1982) LM test for ARCH effect.

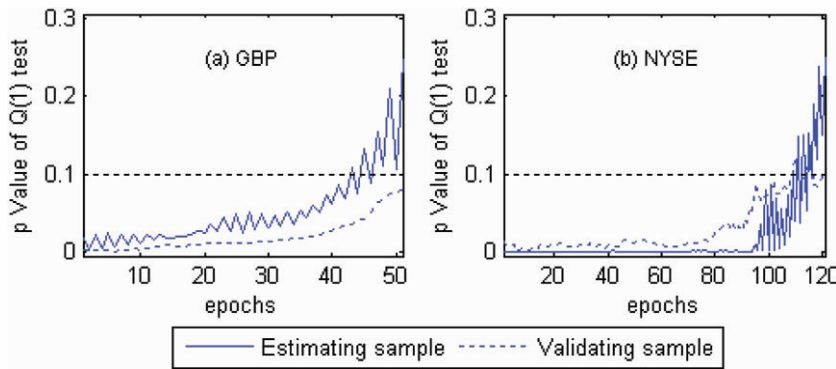


Figure 6. Iterative epochs of recurrent SVR procedure for real data

of raw returns indicates no significant correlation, but the *Q*(6) value of the squared returns reveals that there is significant autocorrelation in the squared returns. The *Q*(6) tests of both raw and squared returns of NYSE are all significant. Engle's (1982) LM tests for ARCH effect show significant evidence in support of GARCH effects (i.e., heteroscedasticity) for both series. Note that the number in parentheses indicates testing at 6 lag order. This examination of daily returns on the GBP and NYSE data reveals that returns can be characterized by heteroscedasticity and time-varying autocorrelation; therefore, we expect the GARCH class models to capture it adequately. Furthermore, as seen from Figure 5 and Table III, it seems that NYSE returns exhibit more variability, skewness, kurtosis and volatility clustering than GBP series such that nonlinear asymmetric EGARCH model should fit it more accurately.

Iterative epochs of recurrent SVM

Because the actual volatility h_t is unobservable for real data analysis, the second conditional variance equation (21b) of the GARCH (1, 1) model should be estimated by using the recurrent SVM procedure, as proposed above. Again, we use cross-validation to determine when the procedure is stopped.

With good forecasting performance as the goal, it is very difficult to figure out when it is best to stop training only in terms of fitting performance. It is possible for the procedure to end up

over-fitting the training data if the training session is not stopped at the right point. We can identify the onset of over-fitting and the stopping point through the use of cross-validation. Figure 6(a) and (b) describes the iterative epochs for volatility prediction of the first training sample of GBP and NYSE, respectively. For the GBP series, the iterative process of recurrent SVM procedure is stopped at the 51st epoch; while, for NYSE, the iterative process is longer and stopped after 121 iterative steps, possibly due to higher kurtosis and more variability and noise behind the NYSE series. Now, we could say, at about the 10% level of significance, the final residuals of equation (21b) obtained from the recurrent SVM procedure have no autocorrelation. In addition, the p -value curves of both estimating and validating samples exhibit a similar pattern (namely, increase with an increasing number of epochs) and point to almost the same stopping point. That is to say, there is no over-fitting phenomenon for the examples illustrated here; the recurrent SVM model does as well on the validating subset as it does on the estimating subset, on which its design is based.

The values taken by the free parameter of SVM and kernel coefficients are also selected according to the sensitivity investigation, similar to that done in Monte Carlo simulation. We do not report the parameter selection process here but present the formal results throughout the real data analysis. For both conditional mean and variance estimation of GBP and NYSE series, fortunately, similar parameter values of feedforward and recurrent SVM procedure could be found as follows: $C = 0.005$, $\varepsilon = 0.05$ and $\sigma^2 = 0.2$. Note that in the analysis of financial returns only the Gaussian kernel is employed for the sake of simplicity due to its best performance among linear, polynomial and Gaussian kernels, as described in Monte Carlo simulation.

Comparing the forecasting ability

The results of out-of-sample volatility forecasting accuracy for each model by using real data are presented in Table IV. Table V reports the p -values of the Diebold–Mariano (DM) test for the difference of MAE loss function in a robust HAC form from Newey–West procedures. In calculating the DM statistic, the null hypothesis of equal forecasting accuracy is related to the four benchmark

Table IV. Measure of volatility forecasting performance for real data

Models	Measures	Moving average	Standard GARCH	EGARCH	ANN-GARCH	SVM-GARCH
GBP	MAE	0.28895	0.24713	0.25719	0.24691	0.23257
	DA	37.29	38.98	49.15	38.98	45.76
NYSE	MAE	1.69610	1.51000	1.44880	1.62980	1.50410
	DA	32.20	42.37	55.93	32.20	57.63

Table V. Diebold–Mariano test for the MAE difference on real data

Models	GBP				NYSE			
	DM1	DM2	DM3	DM4	DM1	DM2	DM3	DM4
Moving average		0.990	0.970	0.981		0.935	0.970	0.813
Standard GARCH	0.010		0.017	0.583	0.065		0.902	0.061
EGARCH	0.030	0.983		0.980	0.030	0.098		0.044
ANN-GARCH	0.019	0.417	0.020		0.187	0.939	0.956	
SVM-GARCH	0.001	0.076	0.000	0.067	0.047	0.054	0.885	0.042

Note: DM1, DM2, DM3 and DM4 are the robust Diebold and Mariano (1995) test by using the Newey–West procedures (Newey and West, 1987) when the benchmark models are the moving average, linear GARCH model, EGARCH model and traditional ANN-GARCH model, respectively. For each test we consider the MAE loss functions.

models: moving average, standard GARCH, EGARCH and ANN models. We specify them as DM1, DM2, DM3 and DM4, respectively. A p -value no greater than 0.05 indicates that the benchmark model yields a higher forecast error (in terms of absolute error) relative to the competing model at the 5% significance level, a p -value no smaller than 0.95 means that benchmark model produces a lower forecast error at 5% level, while a p -value between 0.10 and 0.90 implies that the benchmark and competing models have the equal forecasting accuracy at 10% significance level.

According to MAE measures in Table IV, the SVM-GARCH model is the best one for the GBP series and second for the NYSE series in forecasting the magnitude of volatility error. DM tests in Table V almost statistically favor the SVM-GARCH model as the best model, too, at least at 10% significance level. Even though the MAE metric reveals that the EGARCH model outperforms the SVM-GARCH model for the NYSE series, it is not supported by the DM3 test, which means both models have equal forecasting ability. The better performance of the EGARCH model for NYSE is perhaps due to its ability to capture higher skewness and asymmetry occurring in the SYSE series than in GBP. The standard GARCH model performs modestly in terms of MAE measures, statistically inferior to EGARCH and superior to the ANN-GARCH model for NYSE and significantly better than EGARCH and similar to the ANN-GARCH model for GBP according to DM2 tests. The moving average method is always ranked last in forecasting the magnitude of volatility error, the evidence being significantly supported at least at the 10% level by the DM1 tests in Table V with just one exception, that for NYSE series moving average and ANN-GARCH model have equal forecasting ability. MAE measures and DM3 and DM4 tests denote that the EGARCH model also significantly outperforms the ANN-GARCH model for highly skewed NYSE series but the case is totally reverse for the GBP sample.

Based on DA measures in Table IV, on average, the moving average method is still ranked last, the ANN-GARCH model is ranked second last and the standard GARCH model is ranked at the middle position in forecasting volatility directions. For the GBP series, EGARCH performs best with DA value to be highest 49.15%, followed closely by the SVM-GARCH model; while, for the NYSE model, the best model to forecast volatility turning points is the SVM-GARCH model, with the asymmetric EGARCH model is ranked second, their DA values being 57.63% and 55.93%, respectively.

The empirical evidence of real data also confirms the conclusion obtained in Monte Carlo simulation and favors the theoretical advantage of the SVM-GARCH model. Due to high skewness in financial returns, the asymmetric EGARCH model normally behaves better than the standard GARCH model, particularly in the case of higher skewness or in forecasting volatility turning points. The moving average method always behaves worst and the ANN-GARCH model sometimes good in forecasting one-period-ahead financial volatilities among all candidates.

CONCLUSIONS

In many applications, SVM has shown excellent forecasting performance due to its particular structural design of SRM principle rather than ERM employed by conventional ANN and MLE methods. This inspires us to use it to improve the volatility forecasting ability of the parametric GARCH models. Empirical applications are made for forecasting the simulated data and the real data of daily GBP exchange rates and NYSE stock index.

To avoid the problem that the actual volatility for real data is unobservable, we propose a recurrent SVM procedure with a global feedback loop from the output layer to the input, as opposed to

the feedforward one for simulation, to estimate the conditional volatility equation, that is the ARMA process in nature, of the nonlinear GARCH model. The forecasting performance of the SVM-GARCH model is compared with the moving average, standard GARCH, asymmetric EGARCH and traditional ANN-GARCH models based on two quantitative evaluation measures and robust Diebold–Mariano tests following the Newey–West procedure.

The real data results, together with the simulation evidence, consistently and significantly support the use of the feedforward and recurrent SVM-based GARCH (1, 1) models in forecasting the one-period-ahead volatility error magnitude and direction. The standard GARCH model also performs well in the case of normality and large sample size, while the asymmetric EGARCH model is good at forecasting volatility under the high skewed distribution; but they rarely exceed SVM-GARCH models, at least the Gaussian-type SVM. The recurrent ANN-GARCH model and moving average method behave well only in a few cases. Overall, empirical analysis is in favor of the theoretical advantage of the SVM.

How to choose the appropriate values of free parameters and kernel coefficients and what effect of kernel type in the SVM procedure are investigated by using the sensitivity analysis in Monte Carlo simulation. The iterative process of the proposed recurrent SVM procedure in real data analysis is also examined in detail by the cross-validation method, which is shown to be implemented very easily and could be adopted as another standard SVM construction procedure in other applications.

ACKNOWLEDGEMENTS

The authors acknowledge the editor, Derek Bunn, and the referees for their constructive comments. Thanks also goes to the production editor, Ivry Tan, and my student, Qian Feng, who print and proofread the manuscript. This work is sponsored by Deutsche Forschungsgemeinschaft through the SFB 649 ‘Economic Risk’. Shiyi Chen is also supported by Kyungpook National University Graduate Scholarship for Excellent International Students, Shanghai Leading Academic Discipline Project (No. B101) and State Innovative Institute of Project 985 at Fudan University.

REFERENCES

- Akgiray V. 1989. Conditional heteroskedasticity in time series models of stock returns: evidence and forecasts. *Journal of Business* **62**(1): 55–80.
- Andersen T, Bollerslev T. 1998. Answering the skeptics: yes, standard volatility models do provide accurate forecasts. *International Economic Review* **39**: 885–905.
- Andersen T, Bollerslev T, Diebold F, Labys P. 2003. Modeling and forecasting realized volatility. *Econometrica* **71**: 579–625.
- Andersson J. 2001. On the normal inverse gaussian stochastic volatility model. *Journal of Business and Economic Statistics* **19**: 44–54.
- Andrews D. 1991. Heteroskedasticity and autocorrelation consistent covariance matrix estimation. *Econometrica* **59**: 817–858.
- Ané T, Ureche-Rangau L, Gambet J, Bouverot J. 2008. Robust outlier detection for Asia-Pacific stock index returns. *Journal of International Financial Markets, Institutions and Money* **18**(4): 326–343.
- Awartani B, Corradi V. 2005. Predicting the volatility of the S&P-500 stock index via GARCH models: the role of asymmetries. *International Journal of Forecasting* **21**(1): 167–183.
- Balaban E. 2004. Comparative forecasting performance of symmetric and asymmetric conditional volatility models of an exchange rate. *Economics Letters* **83**(1): 99–105.

- Bauwens L, Sebastien L, Jeroen R. 2006. Multivariate GARCH models: a survey. *Journal of Applied Econometrics* **21**: 79–109.
- Becker R, Clements A, White S. 2007. Does implied volatility provide any information beyond that captured in model-based volatility forecasts? *Journal of Banking and Finance* **31**(8): 2535–2549.
- Becker R, Clements A, McClelland A. 2009. The jump component of S&P 500 volatility and the VIX index. *Journal of Banking and Finance* **33**(6): 1033–1038.
- Bekiros S, Georgoutsos D. 2008. Direction-of-change forecasting using a volatility-based recurrent neural network. *Journal of Forecasting* **27**(5): 407–417.
- Bera A, Jarque C. 1981. An efficient large-sample test for normality of observations and regression residuals. *Australian National University Working Papers in Econometrics*, 40. Canberra.
- Blair B, Poon S,-H, Taylor S. 2001. Forecasting S&P 100 volatility: the incremental information content of implied volatilities and high-frequency index returns. *Journal of Econometrics* **105**: 5–26.
- Bollerslev T. 1986. Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics* **31**: 307–327.
- Bollerslev T, Chou R, Kroner K. 1992. Arch modeling in finance: a review of the theory and empirical evidence. *Journal of Econometrics* **52**: 5–59.
- Box G, Jenkins G, Reinsel G. 1994. *Time Series Analysis: Forecasting and Control*. Prentice Hall: Englewood Cliffs, NJ.
- Brailsford T, Faff R. 1996. An evaluation of volatility forecasting techniques. *Journal of Banking and Finance* **20**: 419–438.
- Brooks C. 1998. Predicting stock index volatility: can market volume help? *Journal of Forecasting* **17**: 59–80.
- Brooks C. 2001. A double-threshold GARCH model for the french franc/deutschmark exchange rate. *Journal of Forecasting* **20**: 135–143.
- Brooks C, Persaud G. 2003. Volatility forecasting for risk management. *Journal of Forecasting* **22**: 1–22.
- Cao C, Tsay R. 1992. Nonlinear time-series analysis of stock volatilities. *Journal of Applied Econometrics* **1**: 165–185.
- Cao L, Tay F. 2001. Financial forecasting using support vector machines. *Neural Computation and Application* **10**: 184–192.
- Chan K, Christie W, Schultz P. 1995. Market structure and the intraday pattern of bid-ask spreads for NASDAQ securities. *Journal of Business* **68**(1): 35–40.
- Chen G, Choi Y, Zhou Y. 2008. Detections of changes in return by a wavelet smoother with conditional heteroscedastic volatility. *Journal of Econometrics* **143**(2): 227–262.
- Chen S, Härdle W, Moro R. 2009. Modeling default risk with support vector machines. *Quantitative Finance* (accepted for publication).
- Chib S, Nardari F, Shephard N. 2002. Markov chain Monte Carlo methods for generalized stochastic volatility models. *Journal of Econometrics* **108**: 281–316.
- Chong C, Ahmad M, Abdullah M. 1999. Performance of GARCH models in forecasting stock market volatility. *Journal of Forecasting* **18**: 333–343.
- Choudhry T, Wu H. 2008. Forecasting ability of GARCH vs kalman filter method: evidence from daily UK time-varying beta. *Journal of Forecasting* **27**(8): 670–689.
- Clements M, Smith J. 1999. A Monte Carlo study of the forecasting performance of empirical SETAR models. *Journal of Applied Econometrics* **14**: 123–141.
- Clements M, Smith J. 2001. Evaluating forecasts from SETAR models of exchange rates. *Journal of International Money and Finance* **20**: 133–148.
- Corradi V, Swanson N. 2004. Some recent developments in predictive accuracy testing with nested models and (generic) nonlinear alternatives. *International Journal of Forecasting* **20**(2): 185–199.
- Corradi V, Distaso W, Swanson N. 2009. Predictive density estimators for daily volatility based on the use of realized measures. *Journal of Econometrics* **150**(2): 119–138.
- Cumby R, Figlewski S, Hasbrouck J. 1993. Forecasting volatility and correlations with EGARCH models. *Journal of Derivatives* Winter: 51–63.
- Day T, Lewis C. 1992. Stock market volatility and the information content of stock index options. *Journal of Econometrics* **52**: 267–287.
- Deng N, Tian Y. 2004. *New Methods in Data Mining: Support Vector Machine*. Science Press: Beijing.
- Diebold F, Mariano R. 1995. Comparing predictive accuracy. *Journal of Business and Economic Statistics* **13**: 253–265.

- Dimson E, Marsh P. 1990. Volatility forecasting without data-snooping. *Journal of Banking and Finance* **44**: 399–421.
- Donaldson R, Kamstra M. 1997. An artificial neural network-GARCH model for international stock return volatility. *Journal of Empirical Finance* **4**: 17–46.
- Dotsis G, Psychoyios D, Skiadopoulos G. 2007. An empirical comparison of continuous-time models of implied volatility indices. *Journal of Banking and Finance* **31**: 3584–3603.
- Dunis C, Huang X. 2002. Forecasting and trading currency volatility: an application of recurrent neural regression and model combination. *Journal of Forecasting* **21**: 317–354.
- Enders W. 2004. *Applied Econometric Time Series* (2nd edn). Wiley: New York.
- Engle R. 1982. Autoregressive conditional heteroskedasticity with estimates of the variance of UK inflation. *Econometrica* **50**: 957–1008.
- Engle R, Hong C-H, Kane A, Noh J. 1993. *Advances in Futures and Options Research*, Vol. 6. JAI Press: Greenwich, CT; 393–415.
- Feng Y, McNeil A. 2008. Modelling of scale change, periodicity and conditional heteroskedasticity in return volatility. *Economic Modelling* **25**(5): 850–867.
- Ferland R, Lalancette S. 2006. Dynamics of realized volatilities and correlations: an empirical study. *Journal of Banking and Finance* **30**(7): 2109–2130.
- Fernandez-Rodríguez F, Gonzalez-Martel C, Sosvilla-Rivero S. 2000. On the profitability of technical trading rules based on artificial neural networks: evidence from the Madrid stock market. *Economics Letters* **69**(1): 89–94.
- Figlewski S. 1997. Forecasting volatility. *Financial Markets, Institutions and Instruments* **6**: 1–88.
- Fleming J. 1998. The quality of market volatility forecasts implied by S&P 100 index option prices. *Journal of Empirical Finance* **5**: 317–345.
- Franke J, Neumann M, Stockis J. 2004. Bootstrapping nonparametric estimators of the volatility function. *Journal of Econometrics* **118**: 189–218.
- Franses P, Dijk DV. 1996. Forecasting stock market volatility using (non-linear) GARCH models. *Journal of Forecasting* **15**(3): 229–235.
- Franses P, Dijk DV. 2000. *Nonlinear Time Series Models in Empirical Finance*. Cambridge University Press: Cambridge, UK.
- Franses P, McAleer M. 2002. Financial volatility: an introduction. *Journal of Applied Econometrics* **17**: 419–424.
- Galbraith J, Kisinbay T. 2005. Content horizons for conditional variance forecasts. *International Journal of Forecasting* **21**: 249–260.
- Gerlach R, Tuyl F. 2006. MCMC methods for comparing stochastic volatility and GARCH models. *International Journal of Forecasting* **22**(1): 91–107.
- Ghysels E, Harvey A, Rebault E. 1996. *Handbook of Statistics: Statistical Methods in Finance*, Vol. 14. Elsevier Science: Amsterdam; 119–191.
- Ghysels E, Santa-Clara P, Valkanov R. 2006. Predicting volatility: how to get most out of returns data sampled at different frequencies. *Journal of Econometrics* **131**: 59–95.
- Glosten L, Jagannathan R, Runkle D. 1992. On the relation between the expected value and the volatility of nominal excess return on stocks. *Journal of Finance* **46**: 1779–1801.
- Gokcan S. 2000. Forecasting volatility of emerging stock markets: linear versus non-linear GARCH models. *Journal of Forecasting* **19**(6): 499–504.
- Gospodinov N, Gavala A, Jiang D. 2006. Forecasting volatility. *Journal of Forecasting* **25**(6): 381–340.
- Gray S. 1996. Modeling the conditional distribution of interest rates as a regime-switching process. *Journal of Financial Economics* **42**: 27–62.
- Groen J, Kapetanios G, Price S. 2009. A real time evaluation of bank of England forecasts of inflation and growth. *International Journal of Forecasting* **25**(1): 74–80.
- Gunn S. 1998. Support vector machines for classification and regression. *Isis-I-98*. Technical report, Image Speech and Intelligent Systems Group, University of Southampton, UK.
- Hamid S, Iqbal Z. 2004. Using neural networks for forecasting volatility of S&P 500 index futures prices. *Journal of Business Research* **57**: 1116–1125.
- Hamilton J. 1989. A new approach to the economic analysis of nonstationary time series and the business cycle. *Econometrica* **57**: 357–384.
- Hamilton J. 1997. *Time Series Analysis*. Princeton University Press: Princeton, NJ.
- Härdle, W, Moro R, Schäfer D. 2005. *Statistical Tools for Finance and Insurance*. Springer: Berlin.

- Härdle W, Moro R, Schäfer D. 2007. *Handbook for Data Visualization*. Springer: Berlin.
- Haykin S. 1999. *Neural Networks: A Comprehensive Foundations* (2nd edn). Prentice Hall: Englewood Cliffs, NJ.
- Heynen R, Kat H. 1994. Volatility prediction: a comparison of stochastic volatility, GARCH(1, 1) and EGARCH(1, 1) models. *Journal of Derivatives* 50–65.
- Hu M, Tsoukalas C. 1999. Combining conditional volatility forecasts using neural networks: an application to the EMS exchange rates. *Journal of International Financial Markets, Institution and Money* 9: 407–422.
- Jorion P. 1995. Predicting volatility in the foreign exchange market. *Journal of Finance* 50: 507–528.
- Jorion P. 1996. *The Microstructure of Foreign Exchange Markets*. Chicago University Press: Chicago, IL.
- Klaassen F. 2002. Improving GARCH volatility forecasts with regime-switching GARCH. *Empirical Economics* 27: 363–394.
- Koopman S, Jungbacker B, Hol E. 2005. Forecasting daily variability of the S&P 100 stock index using historical, realised and implied volatility measurements. *Journal of Empirical Finance* 12: 445–475.
- Kuan C, Liu T. 1995. Forecasting exchange rates using feedforward and recurrent neural networks. *Journal of Applied Econometrics* 10: 347–364.
- Lamoureux C, Lastrapes W. 1993. Forecasting stock-return variance: understanding of stochastic implied volatilities. *Review of Financial Studies* 6: 293–326.
- Lehar A, Scheicher M, Schittenkopf C. 2002. GARCH vs. stochastic volatility: option pricing and risk management. *Journal of Banking and Finance* 26: 323–345.
- Li W, Ling S, McAleer M. 2002. Recent theoretical results for time series models with GARCH errors. *Journal of Economic Surveys* 16: 245–269.
- Lux T, Schornstein S. 2005. Genetic learning as an explanation of stylized facts of foreign exchange markets. *Journal of Mathematical Economics* 41: 169–196.
- Marcucci J. 2005. *Studies in Nonlinear Dynamics and Econometrics*, Vol. 9. Berkeley Electronic Press: Berkeley, CA; 1145.
- McMillan D, Speight A. 2004. Daily volatility forecasts: reassessing the performance of GARCH models. *Journal of Forecasting* 23(6): 449–460.
- McMillan D, Speight A, Gwilym O. 2000. Forecasting UK stock market volatility: a comparative analysis of alternate methods. *Applied Financial Economics* 10: 435–448.
- Meddahi N. 2003. ARMA representations of integrated and realized variances. *Econometrics Journal* 6: 334–355.
- Moosa I. 2000. *Exchange Rate Forecasting: Techniques and Applications*. Macmillan Press: London.
- Neely C. 2009. Forecasting foreign exchange volatility: why is implied volatility biased and inefficient? And does it matter? *Journal of International Financial Markets, Institutions and Money* 19(1): 188–205.
- Nelson D. 1991. Conditional heteroskedasticity in asset returns: a new approach. *Econometrica* 59: 347–370.
- Newey W, West K. 1987. A simple positive, semi-definite, heteroskedasticity and autocorrelation consistent covariance matrix. *Econometrica* 55(3): 703–708.
- Niemira M, Klein P. 1994. *Forecasting Financial and Economic Cycles*. Wiley: New York.
- Pagan A, Schwert G. 1990. Alternative models for conditional stock volatility. *Journal of Econometrics* 45: 267–290.
- Pantelidaki S, Bunn D. 2005. Development of a multifunctional sales response model with the diagnostic aid of artificial neural networks. *Journal of Forecasting* 24: 505–521.
- Park B. 2002. An outlier robust GARCH model and forecasting volatility of exchange rate returns. *Journal of Forecasting* 21(5): 381–393.
- Pérez-Cruz F, Afonso-Rodríguez J, Giner J. 2003. Estimating GARCH models using SVM. *Quantitative Finance* 3: 163–172.
- Pong S, Shackleton M, Taylor S, Xu X. 2004. Forecasting currency volatility: a comparison of implied volatilities and AR(FI) MA models. *Journal of Banking and Finance* 28: 2541–2563.
- Poon S-H, Granger C. 2003. Forecasting volatility in financial markets: a review. *Journal of Economic Literature* 41: 478–539.
- Preminger A, Franck R. 2007. Forecasting exchange rates: a robust regression approach. *International Journal of Forecasting* 23(1): 71–84.
- Qi M, Wu Y. 2003. Nonlinear prediction of exchange rates with monetary fundamentals. *Journal of Empirical Finance* 10: 623–640.

- Renò R. 2006. Nonparametric estimation of stochastic volatility models. *Economics Letters* **90**(3): 390–395.
- Rosenow B. 2008. Determining the optimal dimensionality of multivariate volatility models with tools from random matrix theory. *Journal of Economic Dynamics and Control* **32**(1): 279–302.
- Schittenkopf C, Dorffner G, Dockner E. 2000. Forecasting time-dependent conditional densities: a semi-nonparametric neural network approach. *Journal of Forecasting* **19**: 355–374.
- Scholkopf B, Smola A. 2001. *Learning with Kernels*. MIT Press: Cambridge, MA.
- Sentana E. 1995. Quadratic ARCH models. *Review of Economic Studies* **62**: 639–661.
- Suykens J, Vandewalle J. 2000. Recurrent least squares support vector machines. *IEEE Transactions on Circuits and Systems I* **47**(7): 1109–1114.
- Tay F, Cao L. 2001. Application of support vector machines in financial time series forecasting. *Omega* **29**: 309–317.
- Taylor J. 1999. Evaluating volatility and interval forecasts. *Journal of Forecasting* **18**: 111–128.
- Taylor J. 2000. A quantile regression neural network approach to estimating the conditional density of multiperiod returns. *Journal of Forecasting* **19**: 299–311.
- Taylor N. 2008. Can idiosyncratic volatility help forecast stock market volatility? *International Journal of Forecasting* **24**(3): 462–479.
- Taylor S. 1986. *Modelling Financial Time Series*. Wiley: Chichester.
- Tse Y, Tung S. 1992. Forecasting volatility in the singapore stock market. *Asia Pacific Journal of Management* **9**: 1–13.
- Tseng C, Cheng S, Wang Y, Peng J. 2008. Artificial neural network model of the hybrid EGARCH volatility of the taiwan stock index option prices. *Physica A: Statistical Mechanics and its Applications* **387**(13): 3192–3200.
- Vapnik V. 1995. *The Nature of Statistical Learning Theory*. Springer: New York.
- Vapnik V. 1997. *Statistical Learning Theory*. Wiley: New York.
- West K. 1996. Asymptotic inference about predictive ability. *Econometrica* **64**: 1067–1084.
- West K, Cho D. 1995. The predictive ability of several models of exchange rate volatility. *Journal of Econometrics* **69**: 367–391.
- Wong W, Tu A. 2009. Market imperfections and the information content of implied and realized volatility. *Pacific Basin Finance Journal* **17**(1): 58–79.
- Zhang X, King M. 2005. Influence diagnostics in generalized autoregressive conditional heteroscedasticity processes. *Journal of Business and Economic Statistics* **23**: 118–129.

Authors' biographies:

Shiyi Chen started teaching at the School of Economics of Fudan University in China as an Assistant Professor after receiving his PhD degree in econometrics from the School of Economics and Trade at Kyungpook National University in the Republic of Korea in February 2006. From November 2008, Shiyi Chen became an Associate Professor of Econometrics. His research interests are time series forecasting, nonparametric econometrics, and energy and emission economics. One of his articles, Modeling Default Risk with Support Vector Machines, co-authored with Wolfgang K. Härdle and Rouslan A. Moro, was accepted by the *Journal of Quantitative Finance* in January 2009.

Wolfgang K. Härdle gained his Dr rer. nat. in mathematics at Universität Heidelberg in 1982 and his Habilitation at Universität Bonn in 1988. He is currently Chair Professor of Statistics at the Department of Economics and Business Administration, Humboldt-Universität zu Berlin. He is also director of CASE (Center for Applied Statistics and Economics) and of the Collaborative Research Center 'Economic Risk'. His research focuses on dimension reduction techniques, computational statistics and quantitative finance. He has published 34 books and more than 200 papers in leading statistical, econometrics and finance journals. He is one of the 'Highly Cited Scientists' according to the Institute of Scientific Information.

Kiho Jeong received a PhD degree in econometrics from the University of Wisconsin at Madison in 1991. After working for two years at the Korea Energy Economic Institute as an economist, he joined the School of Economics and Trade at Kyungpook National University in 1994 as an assistant professor, where he is now a full professor. His research interests are forecasting energy/financial markets, modelling climate change effects and nonparametric kernel methods.

Authors' addresses:

Shiyi Chen, School of Economics, Fudan University, 600 Guoquan Road, Shanghai 200433, China.

Wolfgang K. Härdle, Center for Applied Statistics and Economics, Humboldt University, Berlin, Germany.

Kiho Jeong, School of Economics and Trade, Kyungpook National University, Daegu, Republic of Korea.