

# Tree Structured GARCH Models

Francesco Aurino and Peter Bühlmann\*  
ETH Zürich, Switzerland

Revised Version  
October 2000

## Abstract

We propose a new GARCH model with tree-structured multiple thresholds for volatility estimation in financial time series. The approach relies on the idea of a binary tree where every terminal node parameterizes a (local) GARCH model for a partition cell of the predictor space. Fitting of such trees is constructed within the likelihood framework for non-Gaussian observations: it is very different from the well-known CART procedure for regression which is based on residual sum of squares. Our strategy includes the classical GARCH model as a special case and allows to increase model-complexity in a systematic and flexible way. We derive a consistency result and conclude with simulations and real data analysis that the new method has better predictive potential in comparison with other approaches.

*Keywords:* Conditional variance; Financial time series; GARCH model; Maximum likelihood; Threshold model; Tree model; Volatility.

## 1 Introduction

We propose a new method for estimating volatility in stationary financial time series. The real data examples of interest are daily log-returns  $X_t = \log(P_t/P_{t-1})$ , where  $P_t$  denotes the price of an asset at day  $t$ . Our modeling technique is parametric and potentially high-dimensional: it relies on estimating thresholds by using the idea of binary tree construction for partitioning a predictor space.

As a starting point, consider a nonparametric GARCH(1,1) model,

$$\begin{aligned} X_t &= \sigma_t Z_t \quad (t \in \mathbb{Z}), \\ \sigma_t^2 &= f(X_{t-1}, \sigma_{t-1}^2), \quad f: \mathbb{R} \times \mathbb{R}^+ \rightarrow \mathbb{R}^+, \end{aligned} \tag{1.1}$$

where  $(Z_t)_{t \in \mathbb{Z}}$  is a sequence of i.i.d. innovation variables with  $\mathbb{E}[Z_t] = 0$ ,  $\text{Var}(Z_t) = 1$  and  $Z_t$  independent from  $\{X_s; s < t\}$ . The so-called volatility  $\sigma_t$  is then defined as

$$\sigma_t^2 = \text{Var}(X_t | \mathcal{F}_{t-1}) = f(X_{t-1}, \sigma_{t-1}^2),$$

where  $\mathcal{F}_{t-1}$  denotes the  $\sigma$ -algebra (the information) of the variables  $\{X_s; s \leq t-1\}$ . The restriction to model the squared volatility  $\sigma_t^2$  as a function of the previous values  $X_{t-1}$  and  $\sigma_{t-1}^2$  only is natural in finance. Note that it still generates a dependence of  $\sigma_t^2$  from *all previous* observations  $\{X_s; s < t\}$  due to the recursive definition with  $\sigma_{t-1}^2$ : this is the important mathematical

---

\*Corresponding author. Address: Seminar für Statistik, LEO D12, ETH Zentrum, CH-8092 Zürich, Switzerland. E-mail: buhlmann@stat.math.ethz.ch.

difference between ARCH and GARCH models. Also, the implicit assumption in (1.1) that  $\mathbb{E}[X_t|\mathcal{F}_{t-1}] \equiv 0$  is a reasonable approximation for many financial time series: the substantial modeling effort goes into the dominant volatility, although for real data, we usually subtract first a linear AR(1) estimate for the conditional mean. The simplest but often used example for (1.1) is the classical GARCH(1,1) model (Bollerslev, 1986),

$$f(x, \sigma^2) = \alpha_0 + \alpha_1 x^2 + \beta \sigma^2, \quad \alpha_0, \alpha_1, \beta > 0. \quad (1.2)$$

Generally, the unknown function  $f(\cdot, \cdot)$  in (1.1) may be nonlinear and even not smooth; for example, an asymmetry in financial trading with positive and negative return values from the past implies an asymmetric or even discontinuous behavior for  $f(x, y)$  around  $x = 0$ . Estimation of  $f(\cdot, \cdot)$  in general is very difficult due to the non-observable volatility in the second argument. An iterative nonparametric estimation procedure has been proposed in Bühlmann and McNeil (1999). It has the usual advantage about flexibility for  $f(\cdot, \cdot)$ , although quite a few theoretical issues are not rigorously settled yet. Its disadvantages are mainly: poor performance at edges, including the high values of volatility which are of particular interest in practice; lack of ability to deal with non-Gaussian observations; sensitivity to the choice of smoothing parameters. Our approach here is more in the spirit of a *sieve approximation* with parametric models for the nonparametric function  $f(\cdot, \cdot)$ . The approximation builds on the following principles:

- (1) It includes the classical GARCH(1,1) model as a simple special case.
- (2) It uses a binary tree type selection strategy to estimate thresholds (splits) for building up an approximating multiple threshold GARCH model. The binary tree construction, where every terminal node represents a (local) three-dimensional GARCH model, is based on the likelihood in model (1.1).

Item (1) has an important link to practice: there is a relatively strong believe that the classical GARCH(1,1) model is appropriate, despite its simplicity with only three parameters. Our tree structured nested modeling strategy allows to verify such a hypothesis by using known selection techniques for nested models: as we will see, there is potential to improve upon the classical GARCH(1,1) for real data and we will quantify such gains in terms of prediction accuracy for volatility, rather than testing about structural properties of  $f(\cdot, \cdot)$ .

The likelihood driven tree method mentioned in item (2) marks an essential difference to CART (Breiman et al., 1984). Underlying an approximate normality assumption for observations, CART uses residual sum of squares. For financial data, the normality assumption for observations and the corresponding techniques are not appropriate and can result in very poor performance. Our approach resembles more the general tree fitting with the deviance criterion used by Clark and Pregibon (1993). Another difference to CART (or more general versions driven by deviance) is that our tree structured scheme for conditional variance estimation employs a three-dimensional (local) GARCH model in every terminal node from the binary tree; CART uses only one location parameter per node. Finally, our tree GARCH procedure is modeling the function  $f(\cdot, \cdot)$  in (1.1) and hence the *infinite past* in terms of the observations; whereas CART (or versions thereof) in autoregressive modeling deals with a  $p$ -dimensional predictor space ( $p < \infty$ ) from finitely many lagged observations. Our methodology is also markedly different from autoregressive threshold models (SETAR) for conditional expectations, cf. Tong (1990), since we focus on the conditional variance of very non-Gaussian observations, and because we allow for non-Markovian models.

Extending the GARCH(1,1) model in (1.2) in the direction of adding potentially high parametric complexity hasn't been considered yet. Other versions of GARCH(1,1) with three or four

parameters and a GARCH model with one or two thresholds at *fixed* locations (Rabemananjara and Zakoian, 1993) have been proposed. But there seems to be no systematic flexible route how to build up a class of models from the classical GARCH(1,1) in (1.2) to a potentially high dimensional approximation of the general model in (1.1). The paper here deals mainly with this latter task: we describe the methodology in Section 2, present a consistency result in Section 3 and demonstrate the power and use of the new procedure on simulated and real data in Section 4.

## 2 Tree structured GARCH estimation

We describe here our methodology for approximating  $f(\cdot, \cdot)$  in (1.1) by a piecewise linear function. The novel part thereby is the estimation of thresholds in  $\mathbb{R} \times \mathbb{R}^+$ , the starting and end points of piecewise approximating functions. The working model is

$$\begin{aligned} X_t &= \phi X_{t-1} + \sigma_t(\theta) Z_t \quad (t \in \mathbb{Z}), \\ \sigma_t^2(\theta) &= f_\theta(X_{t-1}, \sigma_{t-1}^2(\theta)), \end{aligned} \quad (2.1)$$

with  $\sigma_t(\theta) Z_t$  as in model (1.1), but the functional form  $f(\cdot, \cdot) = f_\theta(\cdot, \cdot)$  now parameterized by a threshold function, see formula (2.2) below. We also add here a linear autoregressive term for estimating a conditional mean (being of minor importance for many financial time series). The function  $f_\theta(\cdot, \cdot)$  is parameterized as binary tree structured GARCH(1,1). It involves a partition

$$\mathcal{P} = \{\mathcal{R}_1, \dots, \mathcal{R}_k\}, \quad \cup_{j=1}^k \mathcal{R}_j = \mathbb{R} \times \mathbb{R}^+, \quad \mathcal{R}_i \cap \mathcal{R}_j = \emptyset \quad (i \neq j)$$

for the predictor space. For every partition cell  $\mathcal{R}_j$ , we employ a GARCH(1,1) model: the parametric form of the function then depends on  $\mathcal{P}$ ,

$$f_\theta(x, \sigma^2) = f_\theta^{\mathcal{P}}(x, \sigma^2) = \sum_{j=1}^k (\alpha_{0,j} + \alpha_{1,j} x^2 + \beta_j \sigma^2) I_{[(x, \sigma^2) \in \mathcal{R}_j]}, \quad (2.2)$$

where  $\theta$  denotes the parameter set  $\{\alpha_{0,j}, \alpha_{1,j}, \beta_j; j = 1, \dots, k\}$ . For  $k = 1$ , we have the classical GARCH(1,1) model from (1.2). As we will discuss in section 2.1, the partition  $\mathcal{P} = \{\mathcal{R}_1, \dots, \mathcal{R}_k\}$  is constructed from a binary tree: every terminal node represents a rectangular partition cell  $\mathcal{R}_j$  whose edges are determined by thresholds.

FIGURE 2.1 ABOUT HERE.

Figure 2.1 represents an example of a binary tree partition of the predictor  $\{(x, \sigma^2); x \in \mathbb{R}, \sigma^2 \in \mathbb{R}^+\}$ : the partition cells, corresponding to terminal nodes in the tree, are  $\mathcal{R}_1 = \{(x, \sigma^2); x \leq d_1\}$ ,  $\mathcal{R}_2 = \{(x, \sigma^2); x > d_1 \text{ and } \sigma^2 \leq d_2\}$ ,  $\mathcal{R}_3 = \{(x, \sigma^2); x > d_1 \text{ and } \sigma^2 > d_2\}$ .

The negative log-likelihood in the working model (2.1) is

$$-\ell(\phi, \theta; X_2^n) = - \sum_{t=2}^n \log \left( \sigma_t^{-1}(\theta) p_Z \left( \frac{X_t - \phi X_{t-1}}{\sigma_t(\theta)} \right) \right), \quad (2.3)$$

where  $p_Z(\cdot)$  denotes the density of the innovation  $Z_t$ . The log-likelihood is always considered conditional on  $X_1$  and some reasonable starting value  $\sigma_1^2(\theta)$ , e.g.  $\sigma_1^2(\theta) = \text{Var}(X_1)$ .

The strategy of our flexible tree structured GARCH estimation is as follows.

- (1) The minimizing criterion is always the negative log-likelihood in (2.3) with innovation density  $p_Z(\cdot)$ , either specified (e.g. standard normal) or of parametric form such as scaled Student's  $t$  with unknown degrees of freedom. The parametric form of  $\sigma_t(\theta)$  in (2.3) is always of the form of a threshold model in (2.2). As we will discuss later in Section 3, the model and innovation density in (2.1) do not need to be true for good approximating properties.
- (2) Optimization with threshold functions in item (1) becomes an estimation and model selection problem. The former is done by maximum likelihood. For the latter, an exhaustive search is computationally prohibitive and we propose a tree structured partial search: within a data-determined tree structure, the optimal model is estimated using the AIC criterion.

## 2.1 Forward entering of thresholds: growing the binary tree

Forward entering of threshold variables, using a binary tree construction, induces a partition for  $\mathbb{R} \times \mathbb{R}^+$  as follows. A first threshold  $d_1 \in \mathbb{R}$  or  $\mathbb{R}^+$  together with a component index  $\iota_1 \in \{1, 2\}$  partitions

$$\mathbb{R} \times \mathbb{R}^+ = \mathcal{R}_{left} \cup \mathcal{R}_{right},$$

where  $\mathcal{R}_{left} = \{(x, \sigma^2) \in \mathbb{R} \times \mathbb{R}^+; (x, \sigma^2)_{\iota_1} \leq d_1\}$  and  $\mathcal{R}_{right}$  analogously but with the relation ' $>$ ' instead. Then, one of the partition cells  $\mathcal{R}_{left}$  or  $\mathcal{R}_{right}$  is again partitioned with a threshold  $d_2$  and a component index  $\iota_2$  in the same fashion. We iterate this procedure: for the  $m$ th iteration step, we specify a pair  $(d_m, \iota_m)$  and an existing partition cell for further refinement by splitting it into two cells. The refinement of an existing partition is always constructed according to the following rule:

$$\begin{aligned} \mathcal{P}^{(old)} = \cup_j \mathcal{R}_j \text{ an existing partition} &\rightarrow \text{pick an element } \mathcal{R}_{j^*} \in \mathcal{P}^{(old)} \\ &\rightarrow \text{split } \mathcal{R}_{j^*} = \mathcal{R}_{j^*,left} \cup \mathcal{R}_{j^*,right} \text{ according to a pair } (d, \iota) \in \mathbb{R} \times \{1, 2\} \\ &\rightarrow \mathcal{P}^{(new)} = \cup_{j \neq j^*} \mathcal{R}_j \cup (\mathcal{R}_{j^*,left} \cup \mathcal{R}_{j^*,right}), \end{aligned} \tag{2.4}$$

with  $(d, \iota)$  describing a threshold and component index,  $\mathcal{R}_{j^*,left} = \{(x, \sigma^2) \in \mathcal{R}_{j^*} \subset \mathbb{R} \times \mathbb{R}^+; (x, \sigma^2)_\iota \leq d\}$  and analogously for  $\mathcal{R}_{j^*,right}$  with the relation ' $>$ '. The whole procedure then produces a partition  $\mathcal{P} = \{\mathcal{R}_1, \dots, \mathcal{R}_k\}$  which can be described as a binary tree whose terminal nodes represent the partition cells, see also Figure 2.1. This is conceptually as in CART (Breiman et al., 1984). But as discussed below, the data-driven construction of a tree-structured partition and estimation of parameters are very different.

Our algorithm below constructs a binary tree, corresponding to a partition of  $\mathbb{R} \times \mathbb{R}^+$ , by optimizing reduction of a conditional negative log-likelihood.

*Step 1.* Compute the negative log-likelihood from (2.3) without partitioning (i.e. in the partition  $\mathcal{P}_{opt}^{(0)} = \mathbb{R} \times \mathbb{R}^+$ ) with

$$f_{\theta^{(0)}}^{\mathcal{P}_{opt}^{(0)}}(x, \sigma^2) = \alpha_0 + \alpha_1 x^2 + \beta \sigma^2, \quad \theta^{(0)} = (\alpha_0, \alpha_1, \beta),$$

and derive from it the maximum likelihood estimates  $\hat{\phi}^{(0)}, \hat{\theta}^{(0)}$  using a quasi-Newton method, cf. Nocedal and Wright (1999). Set  $m = 0$ .

*Step 2.* Increment  $m$  by one. Search for the best refined partition  $\mathcal{P}_{opt}^{(m)}$  by binary splitting of a cell from  $\mathcal{P}_{opt}^{(m-1)}$  as described in (2.4). The details are as follows:

(I) Given  $\mathcal{P}_{opt}^{(m-1)} = \{\mathcal{R}_1, \dots, \mathcal{R}_m\}$ , consider a new partition  $\mathcal{P}^{(m)}$ , where one partition cell  $\mathcal{R}_{j^*} \in \mathcal{P}^{(m-1)}$  is split into  $\mathcal{R}_{j^*} = \mathcal{R}_{j^*,left} \cup \mathcal{R}_{j^*,right}$  as in (2.4). The function associated with  $\mathcal{P}^{(m)}$  is

$$\begin{aligned} f_{(\theta^{(m-1)\setminus*}, \theta^*)}^{\mathcal{P}^{(m)}}(x, \sigma^2) &= \sum_{j \neq j^*} (\alpha_{0,j} + \alpha_{1,j}x^2 + \beta_j\sigma^2) I_{[(x, \sigma^2) \in \mathcal{R}_j]} \\ &+ \sum_{i \in \{j_{left}^*, j_{right}^*\}} (\alpha_{0,i} + \alpha_{1,i}x^2 + \beta_i\sigma^2) I_{[(x, \sigma^2) \in \mathcal{R}_i]}, \end{aligned} \quad (2.5)$$

where

$$\begin{aligned} \theta^{(m-1)\setminus*} &= \{\alpha_{0,j}, \alpha_{1,j}, \beta_j; j = 1, \dots, m, j \neq j^*\} \in (\mathbb{R}^+)^{3(m-1)}, \\ \theta^* &= \{\alpha_{0,i}, \alpha_{1,i}, \beta_i; i \in \{j_{left}^*, j_{right}^*\}\} \in (\mathbb{R}^+)^6. \end{aligned}$$

(II) Compute the minimal negative conditional log-likelihood in the refined partition  $\mathcal{P}^{(m)}$ , holding the parameter vector  $\hat{\phi}^{(m-1)}, \hat{\theta}^{(m-1)\setminus*}$  fixed,

$$\min_{\theta^*} \left( -\ell^{\mathcal{P}^{(m)}}(\hat{\phi}^{(0)}, (\hat{\theta}^{(m-1)\setminus*}, \theta^*); X_2^n) \right), \quad (2.6)$$

by numerical minimization over  $\theta^*$  using a quasi-Newton method. Here,  $-\ell^{\mathcal{P}^{(m)}}$  is as in (2.3) with the function  $f^{\mathcal{P}^{(m)}}(\cdot, \cdot)$  from (2.5). For this numerical minimization over  $\theta^*$ , use as starting values for the parameters  $\theta^*$  in both new cells  $\mathcal{R}_{j^*,left}, \mathcal{R}_{j^*,right}$  the components of  $\hat{\theta}^{(m-1)}$  corresponding to cell  $\mathcal{R}_{j^*}$ .

(III) Optimize (2.6) by varying  $\mathcal{P}^{(m)}$  in (I) and recomputing (II). Denote the optimal refined partition by  $\mathcal{P}_{opt}^{(m)}$ .

*Step 3.* Compute the maximum likelihood in partition  $\mathcal{P}_{opt}^{(m)}$ . Minimize with a quasi-Newton method the negative log-likelihood in (2.3) with  $f^{\mathcal{P}_{opt}^{(m)}}(\cdot, \cdot)$  from (2.5) to obtain  $(\hat{\phi}^{(m)}, \hat{\theta}^{(m)})$ . Thereby, use the starting values  $\hat{\phi}^{(m-1)}, \hat{\theta}^{(m-1)\setminus*}$  and the minimizer  $\hat{\theta}^*$  in (2.6) which is computed in Step 2.

*Step 4.* Repeat Steps 2 and 3 until  $m = M$ . This yields a partition  $\mathcal{P}_{opt}^{(M)}$  corresponding to a large binary tree equipped with parameter estimates  $(\hat{\phi}^{(M)}, \hat{\theta}^{(M)})$ .

*Remark 2.1.* The value  $M$ , corresponding to  $M + 1$  partition cells or terminal nodes in the binary tree, is pre-specified in advance such that the binary tree is sufficiently large. With financial return data, choosing  $M$  around 6 is often appropriate.

*Remark 2.2.* The search for the splitting value in Step 2 is done on a grid: we propose grid-points being empirical  $\alpha$ -quantiles of the data with  $\alpha = i/\text{mesh}$ ,  $i = 1, \dots, \text{mesh} - 1$ . We typically choose  $\text{mesh} = 8$  or  $16$ .

*Remark 2.3.* The conditional log-likelihood in (2.6) yields a substantial computational shortcut compared to a full likelihood approach. For every given partition  $\mathcal{P}^{(m)}$ , the numerical non-linear minimization in (2.6) involves only the 6-dimensional parameter  $\theta^*$ . Since our algorithm searches over many candidate partitions  $\mathcal{P}^{(m)}$  in every iteration step  $m$ , a relatively fast non-linear minimization is important. Finding the best split is thus determined by maximal reduction of the negative conditional log-likelihood.

*Remark 2.4.* The parameter estimates in Step 3 are computed from the full likelihood. For  $(\hat{\phi}^{(m)}, \hat{\theta}^{(m)})$  we take advantage of the fact that the starting values specified in Step 3 are very reasonable for obtaining a reliable and fast maximum likelihood estimate in a possibly high-dimensional parameter space.

## 2.2 Pruning the tree

The binary tree, or the partition  $\mathcal{P}_{opt}^{(M)}$ , constructed in Section 2.1 is too large, or too fine, respectively. We correct by pruning: we search for a best subtree (w.r.t. AIC in (2.7) below) which is most often computationally feasible since  $M$  around 6 seems usually large enough for financial time series. Denote by  $\tau$  the set of all binary subtrees from  $\mathcal{P}_{opt}^{(M)}$ : its elements are denoted by  $\mathcal{P}_i$ . Note that  $\tau$  is generally larger than the set  $\{\mathcal{P}_{opt}^{(0)}, \mathcal{P}_{opt}^{(1)}, \dots, \mathcal{P}_{opt}^{(M)}\}$ .

For every  $\mathcal{P}_i$ , we compute the maximum likelihood estimates  $(\hat{\phi}^{\mathcal{P}_i}, \hat{\theta}^{\mathcal{P}_i})$  with a quasi-Newton method, according to (2.3) with  $f^{\mathcal{P}_i}(\cdot, \cdot)$  of the form (2.2). Note that reasonable starting values are again at hand by going backwards from  $(\hat{\phi}^{(M)}, \hat{\theta}^{(M)})$  in a stage-wise manner. We then consider the penalized negative log-likelihood

$$\text{AIC}(\mathcal{P}_i) = -2\ell(\hat{\phi}^{\mathcal{P}_i}, \hat{\theta}^{\mathcal{P}_i}; X_2^n) + 2(\dim(\hat{\theta}^{\mathcal{P}_i}) + 1) \quad (2.7)$$

as a measure for predictive performance, namely the AIC statistic. The additional contribution 1 in the penalty term arises whenever the conditional expectation parameter  $\phi$  in (2.1) is estimated. Choose the binary tree, or the partition  $\hat{\mathcal{P}}$ , minimizing (2.7). The final tree structured GARCH model is thus given by (2.1) with  $\hat{\phi}^{\hat{\mathcal{P}}}$ , and  $f_{\hat{\theta}^{\hat{\mathcal{P}}}}^{\hat{\mathcal{P}}}(\cdot, \cdot)$  as in (2.2) based on partition  $\hat{\mathcal{P}}$ .

*Remark 2.5.* The AIC-statistic in (2.7) can be replaced by any other sensible model selection criterion. We have experimented with two other versions but found that overall performance with AIC is satisfactory.

## 3 Consistency

We give here some supporting asymptotics for a threshold model specified by (2.1) with Gaussian innovations and (2.2). Extensions to non-Gaussian innovations with suitably nice densities will be analogous. We exclude the effect of model selection and assume that the model structure is fixed, i.e. the structure of a binary tree partition  $\mathcal{P}$  in (2.2) is specified. The unknown parameters of the model are then

$$\xi = (\phi, \theta, d)$$

for the autoregressive, GARCH (within a partition cell) and threshold parameters, respectively. The volatility can thus be written as

$$\begin{aligned} \sigma_t^2(\theta, d) &= f_{\theta}^{\mathcal{P}_d}(X_{t-1}, \sigma_{t-1}^2(\theta, d)), \\ f_{\theta}^{\mathcal{P}_d} &\text{ as in (2.2) but with } \mathcal{P} = \mathcal{P}_d \text{ parameterized with unknown thresholds } d. \end{aligned} \quad (3.1)$$

Since we view the tree GARCH model only as a suitable approximation for the data generating process, we argue here that our procedure yields consistent estimates for the best parameters  $\xi_0$ , projected on the model specified by (2.1) with Gaussian innovations and (3.1),

$$\xi_0 = \operatorname{argmax}_{\xi} h(\xi), \quad h(\xi) = \mathbf{E}\left[\log\left(\sigma_t(\theta, d)^{-1} \varphi\left(\frac{X_t - \phi X_{t-1}}{\sigma_t(\theta, d)}\right)\right)\right]. \quad (3.2)$$

We assume that  $\xi_0$  is unique, e.g.  $h(\cdot)$  is strictly concave. If the model is true,  $\xi_0$  will be the true parameter. Although  $\xi_0$  is unique, the maximum likelihood estimator  $\hat{\xi}$  is not: given data, the log-likelihood is constant between observations with respect to the threshold parameters  $d$ . For theoretical purposes, we define  $\hat{d}$  as the smallest values representing a maximum likelihood estimator.

**Theorem 1.** *Assume that the data generating process  $(X_t)_{t \in \mathbb{Z}}$  is stationary. Under the regularity conditions (A1)–(A5) given in the Appendix, the maximum likelihood estimator in the model specified by (2.1) with Gaussian innovations and (3.1) satisfies,*

$$\hat{\xi} \rightarrow \xi_0 \text{ in probability, as } n \rightarrow \infty.$$

A proof is given in the Appendix. Note that Theorem 1 describes consistency in mis-specified models. The true model isn't necessarily of GARCH-type as in (1.1); or the correct model is of the form (1.1) but with non-normal innovations  $(Z_t)_t$ . The interpretation of the consistency result is then as follows: the model with parameter  $\xi_0$  is the closest element in the class of threshold models, described by (2.1) and (3.1), to the true underlying stationary process with respect to the Kullback-Leibler divergence, see for example Shibata (1997). Particularly, if the fitted model is true, the true parameters are consistently estimated.

Asymptotic normality with convergence rate  $n^{-1/2}$  for the parameters  $\phi, \theta$  usually holds due to partial differentiability with respect to  $\phi, \theta$  of the log-likelihood function; this is outlined in Pollard (1984, Ch. VII). Under the realistic assumption that the data generating model is not exactly of threshold-type form, the estimated thresholds  $\hat{d}$  will generally have a different limiting distribution and, provided that the true model has a smooth distribution, its convergence rate is slower than  $n^{-1/2}$ , due to a non-continuous likelihood function with respect to the split point parameters. A general relevant theory for this phenomenon is given by Kim and Pollard (1990); see also Hansen (2000).

## 4 Numerical results

We consider here the performance of tree structured GARCH models for simulated and real data. We compare the results with the GARCH(1,1) estimate in (1.2) and a nonparametric generalized additive model for log-transformed squared data which is described in the Appendix. We always report with the use of  $M = 5$  in Step 4 from Section 2.1 and with grid search using  $\text{mesh} = 8$  as described in Remark 2.2 (except in Table 3.2 where  $\text{mesh} = 16$  is used in addition): these specifications lead to good tree structured model fits, despite their somewhat simple nature.

### 4.1 Simulations

The models that we use for simulating data are as in (1.1) with various  $f(\cdot, \cdot)$ . One is a threshold model

$$f(x, \sigma^2) = \begin{cases} 0.1 + 0.5x^2, & \text{if } x \leq d_1 = 0, \\ 0.2 + 0.2x^2 + 0.75\sigma^2, & \text{if } x > d_1 = 0 \text{ and } \sigma^2 \leq d_2 = 0.5, \\ 0.8 + 0.5\sigma^2, & \text{if } x > d_1 = 0 \text{ and } \sigma^2 > d_2 = 0.5, \end{cases} \quad (4.1)$$

The parameters are chosen to mimic time series of real log-returns. Also, the first threshold  $d_1 = 0$  allows for a natural asymmetry in finance. Another model is a classical GARCH(1,1)

$$f(x, \sigma^2) = 0.05 + 0.1x^2 + 0.85\sigma^2. \quad (4.2)$$

A third model is neither GARCH nor threshold GARCH

$$f(x, \sigma^2) = (0.1 + 0.2|x| + 0.9x^2) \cdot (0.8 \exp(-1.5|x||\sigma|)) + (0.4x^2 + 0.5\sigma^2)^{3/4} \quad (4.3)$$

The distribution of innovations is either standard normal  $Z_t \sim \mathcal{N}(0, 1)$  or scaled  $t_6$  so that  $\sqrt{6/4}Z_t \sim t_6$  has again variance one. We always take sample size  $n = 1000$ : for real daily data,

this would correspond to about four years which is a reasonable window where stationarity is expected to be approximately true.

Estimation is always based using the knowledge that  $\mu_t = \mathbf{E}[X_t | \mathcal{F}_{t-1}] \equiv 0$ , i.e.  $\phi = 0$  in (2.1). Figures 4.1 and 4.2 display some results from the tree structured GARCH model, in comparison with the classical fit of a GARCH(1,1) in (1.2) and with an estimated generalized additive model (GAM) described in the Appendix.

FIGURE 4.1 ABOUT HERE.

We observe the following in Figure 4.1. The tree structured GARCH fit using  $\mathcal{N}(0, 1)$  innovations overestimates the number of thresholds: the first two thresholds are approximately correct. An improvement is given by the tree structured GARCH fit with scaled  $t_\nu$ -distributed innovations: the thresholds and also the maximum likelihood estimated degrees of freedom  $\hat{\nu} = 5.12$  for the true  $\nu = 6$  are very satisfactory. The classical GARCH(1,1) model with scaled  $t_\nu$ -distributed innovations yields  $\hat{\nu} = 4.37$  and of course, it doesn't exhibit any thresholds in the volatility surface. There is no surprise that the tree structured GARCH with scaled  $t_\nu$ -distributed innovations is best, since the true model is of this form. But the tree GARCH estimate with  $\mathcal{N}(0, 1)$  misspecified innovations fits, as a quasi-maximum-likelihood, still reasonably good.

FIGURE 4.2 ABOUT HERE

Figure 4.2 exploits a desirable important feature: the tree structured GARCH with scaled  $t_\nu$ -distributed innovations performs very well in regions of high conditional variance. This is not the case with the classical GARCH(1,1) fit; and the GAM estimate described in the Appendix is very poor in regions of high volatility. Note the different scales for the various procedures in Figure 4.2.

For quantifying the goodness of fit, we consider various measures:

$$\text{IS-L}_p = \sum_{t=1}^n |\sigma_t^2 - \hat{\sigma}_t^2|^p, \quad p = 1, 2, \quad (\text{in-sample loss})$$

the AIC statistic from (2.7),

$$\text{OS-L}_p = \sum_{t=1}^n |\sigma_t^2 - \hat{\sigma}_t^2(Y_1^{t-1})|^p, \quad p = 1, 2, \quad Y_1^n \text{ a new test set (out-sample loss),}$$

where for OS-L,  $\hat{\sigma}_t^2(Y_1^{t-1})$  is using the estimated model from the data  $X_1^n$  but evaluates it on new test data  $Y_1^n$  being another independent realization of the data. Both, the out-sample OS-L- and AIC-statistic are measures for predictive performance. The IS- and even more the OS-L-statistics are interesting measures for our simulations, but we can't calculate them for real data examples. Detailed results are reported in Tables 4.1–4.2, where we denote by data 1–7 different independent realizations from model (2.1):

- data 1–3 with (4.1) and  $\mathcal{N}(0, 1)$  innovations,
- data 4 with (4.1) and scaled innovations  $\sqrt{6/4}Z_t \sim t_6$ ,
- data 5 with (4.2) and  $\mathcal{N}(0, 1)$  innovations,
- data 6 and 7 with (4.3) and  $\mathcal{N}(0, 1)$  innovations.

TABLES 4.1–4.2 ABOUT HERE.

As an out-sample statistic, we view OS-L<sub>2</sub> as the most important measure for simulations. It gives more weight to large deviations than the OS-L<sub>1</sub> criterion: therefore, it is more appropriate when studying the estimate in regions of high volatility.



The tree structured GARCH procedure consistently outperforms the classical GARCH(1,1): for data 5 (GARCH(1,1) realization), it doesn't fit any thresholds and coincides with the GARCH(1,1) fit. Our new tree GARCH procedure compares favorably with the GAM estimate: it is much better for data 5 (GARCH(1,1) realization) and data 1–4 (threshold GARCH realizations), but slightly worse for data 6 and 7 (realizations with (4.3)). The classical GARCH(1,1) can exhibit huge variation in out-sample accuracy OS-L, for example a poor performance in data 2. This may be due to a very flat negative log-likelihood at the observed data: see Zumbach (1999) who also proposes a remedy. Interestingly, the tree structured GARCH model exhibits here much more stability.

## 4.2 Two real data examples

We consider two financial instruments with 1000 daily negative log-returns  $X_t = -100 \log(P_t/P_{t-1})$  (in percentages): from the German DAX index between January 18, 1994 and November 17, 1997; and from the BMW stock price between September 23, 1992 and July 23, 1996. We consider the tree structured GARCH model, again in comparison with the GARCH(1,1) from (1.2) (both with the additional model term  $\phi X_{t-1}$  from (2.1) for  $\mathbf{E}[X_t|\mathcal{F}_{t-1}]$ ), and with the GAM model described in the Appendix.

FIGURE 4.3 ABOUT HERE.

Figure 4.3 shows the result for the DAX index from a tree GARCH fit with  $\mathcal{N}(0, 1)$ -distributed innovations. Three thresholds are fitted in the volatility surface. Graphical diagnostics for the residuals is satisfactory, with a tendency for heavier tails than standard normal.

For these real-data examples, we measure goodness of fit with the AIC statistic from (2.7) and with in- and out-sample losses for predicting centered second moments,

$$\begin{aligned} \text{IS-PL}_2 &= \sum_{t=1}^n \left( \hat{\sigma}_t^2 - (X_t - \hat{\mu}_t)^2 \right)^2 \quad (\text{in-sample prediction loss}), \\ \text{OS-PL}_2 &= \sum_{t=1}^n \left( \hat{\sigma}_t^2(Y_1^{t-1}) - (Y_t - \hat{\mu}_t(Y_1^{t-1}))^2 \right)^2, \quad Y_1^n \text{ a new test set (out-sample prediction loss)} \end{aligned}$$

with  $\hat{\mu}_t = \hat{\phi} X_{t-1}$  and  $\hat{\mu}_t(Y_1^{t-1})$ ,  $\hat{\sigma}_t^2(Y_1^{t-1})$  using the estimated model from the data  $X_1^n$  but evaluating at test data  $Y_1^n$  with  $n = 1000$ .

TABLE 4.3 ABOUT HERE.

We view OS-PL<sub>2</sub> as the most important measure, followed by AIC: as an in-sample loss, IS-PL<sub>2</sub> is not very relevant. The tree structured model improves upon the classical GARCH(1,1) (both with  $\mathcal{N}(0, 1)$ -distributed innovations) for the DAX index; for the BMW data, the two procedures have about the same performance. This is consistent with a common belief that classical GARCH(1,1) is better for individual prices than indices. The GAM estimate has the poorest performance in both data sets: for the BMW series, its behavior is extremely poor.

## 4.3 Summarizing numerical results

The tree structured GARCH procedure often outperforms the estimate from the GARCH(1,1) model in (1.2) and the nonparametric GAM fit described in the Appendix: such better performance is with respect to many goodness of fit and graphical criteria. More specifically:

- (1) For real data, the tree structured GARCH procedure was best for log-returns from an index (DAX) and equally good as the classical GARCH(1,1) fit for an individual share (BMW). The GAM estimate was poorest on real data.
- (2) For simulated data, the tree structured GARCH is better than classical GARCH(1,1). In one case, the GAM fit was found slightly better than tree structured GARCH; but overall, GAM was poorer and it may be very unstable, resulting in extremely low performance.
- (3) The tree structured GARCH estimate may be much better in the interesting regions where the true volatility is high, see Figure 4.2. The nonparametric GAM fit can be extremely poor in regions of high volatility (this problem doesn't disappear when trying other smoothing parameters).
- (4) The tree structured fitting procedure is sometimes slightly improved by assuming scaled  $t_\nu$ -distributed innovations  $Z_t$  for the likelihood in (2.3), provided that the underlying innovations are heavier tailed. This weakly evident in Figure 4.3 for real data. See also Table 4.1.
- (5) The AIC-statistic is an indicator for ranking out-sample performance; and pruning with AIC in (2.7) works reasonably well.

Items (1), (2) and (3) indicate a strong advantage of parametric, likelihood based methods over nonparametric least squares smoothing techniques such as the GAM specification used here or multiplicative nonparametric ARCH models in Hafner (1998) or Yang et al. (1999). As pointed out in (4), the likelihood approach can be easily modified to heavier tailed innovations inducing then even more heavy tails for the observations in the model. If performance is judged with a criterion putting emphasis on accurate prediction in high volatility regions, our tree based GARCH model seems clearly best among all alternative methods considered here.

#### 4.4 How appropriate is GARCH(1,1) for daily returns?

The GARCH(1,1) model in (1.2) is very popular for analyzing daily log-returns of financial assets: it is often argued that it performs well despite that it has only three parameters describing a very low-dimensional model for sample size in the range of 1000. We quantify here the possible gains by using the flexible tree structured GARCH model. In virtually all examples of daily log-return stock data, the new tree GARCH procedure often improves upon the classical GARCH(1,1) and was never found to be significantly worse: the difference in performance has typically quantitative magnitude as reported in Sections 4.1- 4.3, with a remarkable gain for regions of high volatility. With real data, the first split has always been found in the  $x$ -axis around zero: this is compatible with the interpretation that there is an asymmetric behavior depending on the sign of the previous log-return.

## 5 Concluding remarks

We propose a tree structured GARCH model which is more flexible and accurate for prediction of volatility in financial time series than classical GARCH(1,1). The modeling strategy includes the classical GARCH(1,1) as a special case (no thresholds) and allows to increase complexity in a *systematic* way. On finite data, the new method compares favorably with a nonparametric technique based on additive models: especially in the interesting regions where the true volatility is moderate or large.

As supporting asymptotics, we present a consistency result about estimation of a best tree structured GARCH model for approximating a general stationary process, see Section 3. More refined statistical inference is difficult due to the non-continuous nature of thresholds or trees: limiting distributions of estimated thresholds are typically non-Gaussian, see the discussion following Theorem 1. If the primary goal is construction of better volatility forecasts, which in turn can be used for dynamic risk management (cf. McNeil and Frey, 2000), we choose the route to optimize an information or complexity criterion rather than the somewhat inappropriate structural tool of testing for a model structure. As a simple solution, we use the AIC criterion: we have gained evidence on numerical examples, that it can be used as a reasonable guideline.

Our univariate tree structured GARCH procedure has a straightforward application in multivariate models where the conditional variance of an individual series is modeled as a function of the individual lagged values and individual lagged volatility. The multivariate cross-dependence is then modeled with cross-dependent innovations. For example, individual tree GARCH models for volatilities lead to an attractive version of the multivariate, constant conditional correlation model (Bollerslev, 1990). Another straightforward extension of our methodology is tree structured GARCH( $p, q$ ) modeling with  $p > 1$  or  $q > 1$ . As already mentioned in Section 1, this may be of minor importance since the general model in (1.1) is in vogue and believed to capture the most important aspects of the underlying mechanism.

**Acknowledgments:** Comments by two referees helped to improve the presentation of our results.

## References

- [1] Bollerslev, T. (1986). Generalized autoregressive conditional heteroskedasticity. *J. of Econometrics* **31**, 307–327.
- [2] Bollerslev, T. (1990). Modelling the coherence in short-run nominal exchange rates: a multivariate generalized ARCH model. *The Review of Economics and Statistics* **72**, 498–505.
- [3] Breiman, L., Friedman, J.H., Olshen, R.A. & Stone, C.J. (1984). *Classification and Regression Trees*. Wadsworth, Belmont (CA).
- [4] Bühlmann, P. and McNeil, A.J. (1999). Nonparametric GARCH models. Preprint, ETH Zürich.
- [5] Clark, L.A. and Pregibon, D. (1993). Tree-based models. In *Statistical Models in S*, pp.377–419, auths. J.M. Chambers and T.J. Hastie. Chapman & Hall, London.
- [6] Hafner, C.M. (1998). Estimating high-frequency foreign exchange rate volatility with non-parametric ARCH models. *J. of Statistical Planning and Inference* **68**, 247–269.
- [7] Hansen, B.E. (2000). Splitting and threshold estimation. *Econometrica* **68**, 575–603.
- [8] Hastie, T.J. and Tibshirani, R.J. (1990). *Generalized Additive Models*. Chapman & Hall, London.
- [9] Kim, J. and Pollard, D. (1990). Cube root asymptotics. *Annals of Statistics* **18**, 191–219.
- [10] McNeil, A.J. and Frey, R. (2000). Estimation of tail-related risk measures for heteroscedastic financial time series: an extreme value approach. To appear in *J. of Empirical Finance*.

- [11] Nocedal, J. and Wright, S.J. (1999). *Numerical Optimization*. Springer, New York.
- [12] Pollard, D. (1984). *Convergence of Stochastic Processes*. Springer, New York.
- [13] Rabemananjara, R. and Zakoian, J.M. (1993). Threshold ARCH models and asymmetries in volatility. *J. of Applied Econometrics* **8**, 31–49.
- [14] Shibata, R. (1997). Bootstrap estimate of Kullback-Leibler information for model selection. *Statistica Sinica* **7**, 375–394.
- [15] Tong, H. (1990). *Non-linear Time Series. A Dynamical System Approach*. Oxford University Press.
- [16] Yang, L., Härdle, W. and Nielsen, J.P. (1999). Nonparametric autoregression with multiplicative volatility and additive mean. *J. of Time Series Analysis* **20**, 579–604.
- [17] Yu, B. (1994). Rates of convergence for empirical processes of stationary mixing sequences. *Ann. Probab.* **22**, 94–116.
- [18] Zumbach, G. (1999). The pitfalls in fitting GARCH(1,1) processes. Preprint. Olsen & Associates, Zürich.

## Appendix

### Assumptions and proof of Theorem 1

We first give and discuss a set of regularity conditions for Theorem 1.

- (A1) The data generating process  $(X_t)_{t \in \mathbb{Z}}$  is stationary,  $\beta$ -mixing (cf. Yu, 1994) with  $\beta(k) \leq Ck^{-\delta}$  for some  $C, \delta > 0$ . Moreover,  $\mathbb{E}|X_t|^2 < \infty$ .
- (A2) Consider the parameter set  $\mathcal{L} = \{(\theta, d); \sup_t \sigma_t(\theta, d) = O_P(1)\}$ , i.e. the set where the volatility doesn't 'explode'. Denote by  $\Xi = \{\xi; (\theta, d) \in \mathcal{L}, |\phi| < 1\}$ . Assume that the best projected parameter vector  $\xi_0$  is an interior point of  $\Xi$ .
- (A3) The function  $f_\theta^{\mathcal{P}^d}(\cdot, \cdot)$  in (3.1) and  $(X_t)_{t \in \mathbb{Z}}$  are such that  $h(\cdot)$  in (3.2) has a unique maximizer and is continuous for  $\xi \in \Xi$ .
- (A4) Define the truncated squared volatility as

$$\begin{aligned} (\sigma_t^{\text{trunc}(p)}(\theta, d))^2 &= f_{X_{t-1}} \circ f_{X_{t-2}} \circ \dots \circ f_{X_{t-p}}(X_{t-p}^2), \\ f_x(\tau^2) &= f_\theta^{\mathcal{P}^d}(x, \tau^2), \end{aligned}$$

i.e. an approximation with  $p$  (instead of infinitely many) lagged observed values for the volatility. Assume that the function  $f_\theta^{\mathcal{P}^d}(\cdot, \cdot)$  in (3.1) satisfies:

$$\sup_{t, (\theta, d) \in \mathcal{L}} |\sigma_t(\theta, d) - \sigma_t^{\text{trunc}(p)}(\theta, d)| = O_P(\gamma(p)), \text{ for some } \gamma(p) \rightarrow 0 \text{ (} p \rightarrow \infty \text{)}.$$

- (A5) The function  $f_\theta^{\mathcal{P}^d}(\cdot, \cdot)$  is such that

$$\inf_{x, \sigma^2, (\theta, d) \in \mathcal{L}} f_\theta^{\mathcal{P}^d}(x, \sigma^2) > 0.$$

Assumption (A1) indicates that we do not require the data being generated from a GARCH-type model as in (1.1) or (2.1) and (2.2). (A2) is related to stationarity of the model. Continuity of  $h(\cdot)$  in (A3) is implied by smoothness of the distribution of  $(X_t)_{t \in \mathbb{Z}}$ . (A4) requires that the non-Markovian volatility process can be approximated by the truncated Markovian model; this holds for example, if all  $|\beta_j| < 1$  in (2.2) and  $(X_t)_{t \in \mathbb{Z}}$  having a smooth distribution. The assumption is related to the Shannon-McMillan-Breiman Theorem. (A5) implies strict positivity of the volatility function (being reasonable).

*Proof of Theorem 1.* For  $n$  large enough, the maximizer of the likelihood is in  $\Xi$  (otherwise, the values  $\sigma_t(\hat{\theta}, \hat{d})$  would explode implying a small likelihood).

Denote the log-likelihood (conditioned on the first observation and using any  $\sigma_1 \neq 0$ ) by

$$\ell_n(\xi) = n^{-1} \sum_{t=2}^n \log \left( \sigma_t^{-1} \varphi \left( \frac{X_t - \phi X_{t-1}}{\sigma_t} \right) \right),$$

abbreviating  $\sigma_t(\theta, d)$  by  $\sigma_t$ . Analogously, define the truncated log-likelihood as

$$\ell_n^{trunc(p)}(\xi) = n^{-1} \sum_{t=p+1}^n \log \left( (\sigma_t^{trunc(p)})^{-1} \varphi \left( \frac{X_t - \phi X_{t-1}}{\sigma_t^{trunc(p)}} \right) \right).$$

We will show that

$$\sup_{\xi \in \Xi} |\ell_n(\xi) - \ell_n^{trunc(p)}(\xi)| = O_P(\gamma(p)) \quad (5.1)$$

with  $\gamma(p)$  as in assumption (A4). Similarly, denoting by

$$h^{trunc(p)}(\xi) = \mathbf{E} \left[ \log \left( \frac{1}{\sigma_t^{trunc(p)}(\theta, d)} \varphi \left( \frac{X_t - \phi X_{t-1}}{\sigma_t^{trunc(p)}(\theta, d)} \right) \right) \right],$$

we get for the population version,

$$\sup_{\xi \in \Xi} |h(\xi) - h^{trunc(p)}(\xi)| = O(\gamma(p)). \quad (5.2)$$

We give the proofs for (5.1) and (5.2) at the end. Now, the idea is to work with the truncated log-likelihood which is of the form

$$\begin{aligned} \ell_n^{trunc(p)}(\xi) &= n^{-1} \sum_{t=p+1}^n g(X_t, \dots, X_{t-p+1}; \xi), \\ g(X_t, \dots, X_{t-p+1}; \xi) &= \log \left( (\sigma_t^{trunc(p)})^{-1} \varphi \left( \frac{X_t - \phi X_{t-1}}{\sigma_t^{trunc(p)}} \right) \right). \end{aligned}$$

By the mixing property of  $(X_t)_{t \in \mathbb{Z}}$ ,

$$\sup_{\xi \in \Xi} |\ell_n^{trunc(p)}(\xi) - h^{trunc(p)}(\xi)| = o_P(1), \quad (5.3)$$

for every  $p$ , see Yu (1994). More details proving (5.3) are given at the end. By (5.1)–(5.3) it follows that

$$\sup_{\xi \in \Xi} |\ell_n(\xi) - h(\xi)| = o_P(1). \quad (5.4)$$

Now, use a sandwich argument:

$$\sup_{\xi \in \Xi} |\ell_n(\xi) - h(\xi)| + h(\xi_0) \geq |\ell_n(\hat{\xi}) - h(\hat{\xi})| + h(\hat{\xi}) \geq \ell_n(\hat{\xi}) \geq \ell_n(\xi_0) = (\ell_n(\xi_0) - h(\xi_0)) + h(\xi_0).$$

We thereby have used that  $\xi_0$  and  $\hat{\xi}$  are the maximizers of  $h(\cdot)$  and  $\ell_n(\cdot)$ , respectively. Due to (5.4), the left and right hand side are asymptotically equal to  $h(\xi_0)$ , implying that  $h(\hat{\xi}) = h(\xi_0) + o_P(1)$ . Since  $h(\cdot)$  is continuous and  $\xi_0$  its unique maximizer, it follows that  $\hat{\xi} = \xi_0 + o_P(1)$ , i.e. consistency.

Proof of (5.1) and (5.2): using a first order Taylor expansion,

$$\begin{aligned} \log(\sigma^{-1}\varphi(u/\sigma)) &= \log(\tau^{-1}\varphi(u/\tau)) + h(u, \bar{\sigma})(\sigma - \tau), \quad |\sigma - \bar{\sigma}| < |\sigma - \tau|, \\ h(u, \bar{\sigma}) &= -\bar{\sigma}^{-1} + \bar{\sigma}^{-4}u^2\varphi(u/\bar{\sigma}). \end{aligned}$$

Use this with  $u = X_t - \phi X_{t-1}$ ,  $\sigma = \sigma_t$ ,  $\tau = \sigma_t^{trunc(p)}$ . Assumption (A5), implying that  $\sigma_t, \sigma_t^{trunc(p)}$  are bounded away from zero, yields that  $\Delta_t(\xi) = h(X_t - \phi X_{t-1}, \bar{\sigma}_t(\theta, d))$  is uniformly bounded over  $t$  and  $\xi$ . Thus,

$$\begin{aligned} \ell_n(\xi) &= \ell_n^{trunc(p)}(\xi) + n^{-1} \sum_{t=p+1}^n \Delta_t(\xi)(\sigma_t - \sigma_t^{trunc(p)}), \\ \sup_{t, \xi \in \Xi} |\Delta_t(\xi)| &= O(1). \end{aligned}$$

Now use (A4): since  $\gamma(p)$  is uniform with respect to  $t$  and  $(\theta, d)$ , we obtain (5.1). Formula (5.2) follows similarly.

Proof of (5.3): apply Theorem 3.4 in Yu (1994). Consider the permissible class  $\mathcal{G} = \{g(\cdot; \xi); \xi \in \Xi\}$ . An envelope function  $G(\cdot)$  can be constructed straightforwardly using (A5), and integrability of  $G(\cdot)$  is implied by the moment assumption in (A1). Furthermore, the class  $\mathcal{G}$  consists of functions which are compositions of a smooth, with a piecewise (involving indicators) quadratic function (from GARCH model in a partition cell): a good bound for the metric entropy condition then follows.  $\square$

## Estimation with generalized additive model (GAM)

A nonparametric estimate for  $\sigma_t^2$  can be derived as follows.

1. Estimate the conditional mean  $\mu_t = \mathbf{E}[X_t | \mathcal{F}_{t-1}]$  by an AR(1) model,

$$\hat{\mu}_t = \hat{\phi} X_{t-1}$$

with parametric estimate  $\hat{\phi}$  obtained from least squares fitting.

2. Compute  $Y_t = \log((X_t - \hat{\mu}_t)^2)$ ,  $t = 2, \dots, n$ .

3. In model (2.1), we have  $Y_t \approx \beta + \log(\sigma_t^2) + (\log(Z_t^2) - \beta)$  with  $\beta = \mathbf{E}[\log(Z_t^2)]$ . Denote by  $\gamma_t = \beta + \log(\sigma_t^2)$ . Fit a GAM model with the transformed data  $Y_2^n$ ,

$$\hat{\gamma}_t = \hat{h}_1(X_{t-1}) + \hat{h}_2(X_{t-2}) + \dots + \hat{h}_k(X_{t-k}), \quad t = k+1, \dots, n,$$

with nonparametric estimates  $\hat{h}_i(\cdot)$  obtained from a least squares backfitting algorithm with response variables  $Y_t$ , cf. Hastie and Tibshirani (1990). The optimal value of  $k$  is chosen by minimizing the AIC statistic for Gaussian additive modeling of  $Y_t$ .

4. Back-transform  $\delta_t = \exp(\hat{\gamma}_t) \approx \exp(\beta)\sigma_t^2 = \frac{1}{c}\sigma_t^2$  and build  $R_t^2 = (X_t - \hat{\mu}_t)^2/\delta_t \approx c Z_t^2$ . Thus, set

$$\hat{c} = (n)^{-1} \sum_{t=1}^n R_t^2.$$

5. Then, set

$$\hat{\sigma}_t^2 = \hat{c}\delta_t.$$

6. Iterate steps 1.-5. Thereby use weighted estimation in step 1. with weights  $w_t = \frac{1}{\hat{\sigma}_t}$ , where  $\hat{\sigma}_t^2$  is the estimate from the previous iteration step. Stop iterating by monitoring convergence of  $\hat{\sigma}_t^2$  and  $\hat{\mu}_t$ .

A related technique is given in Yang et al. (1999): they don't use the log-transform but work with the squared observations and dependent, but uncorrelated innovations.

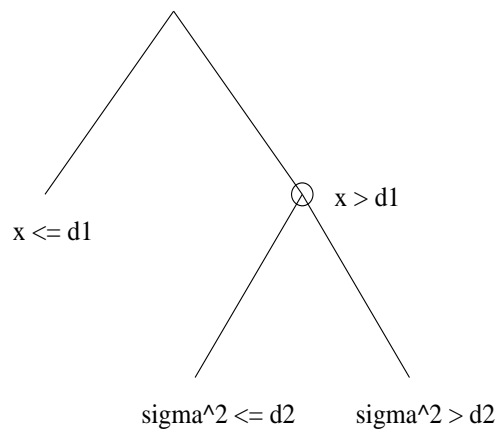


Figure 2.1: Binary tree partition  $\mathcal{P} = \{\mathcal{R}_1, \mathcal{R}_2, \mathcal{R}_3\}$  of predictor space  $\{(x, \sigma^2); x \in \mathbb{R}, \sigma^2 \in \mathbb{R}^+\}$ .



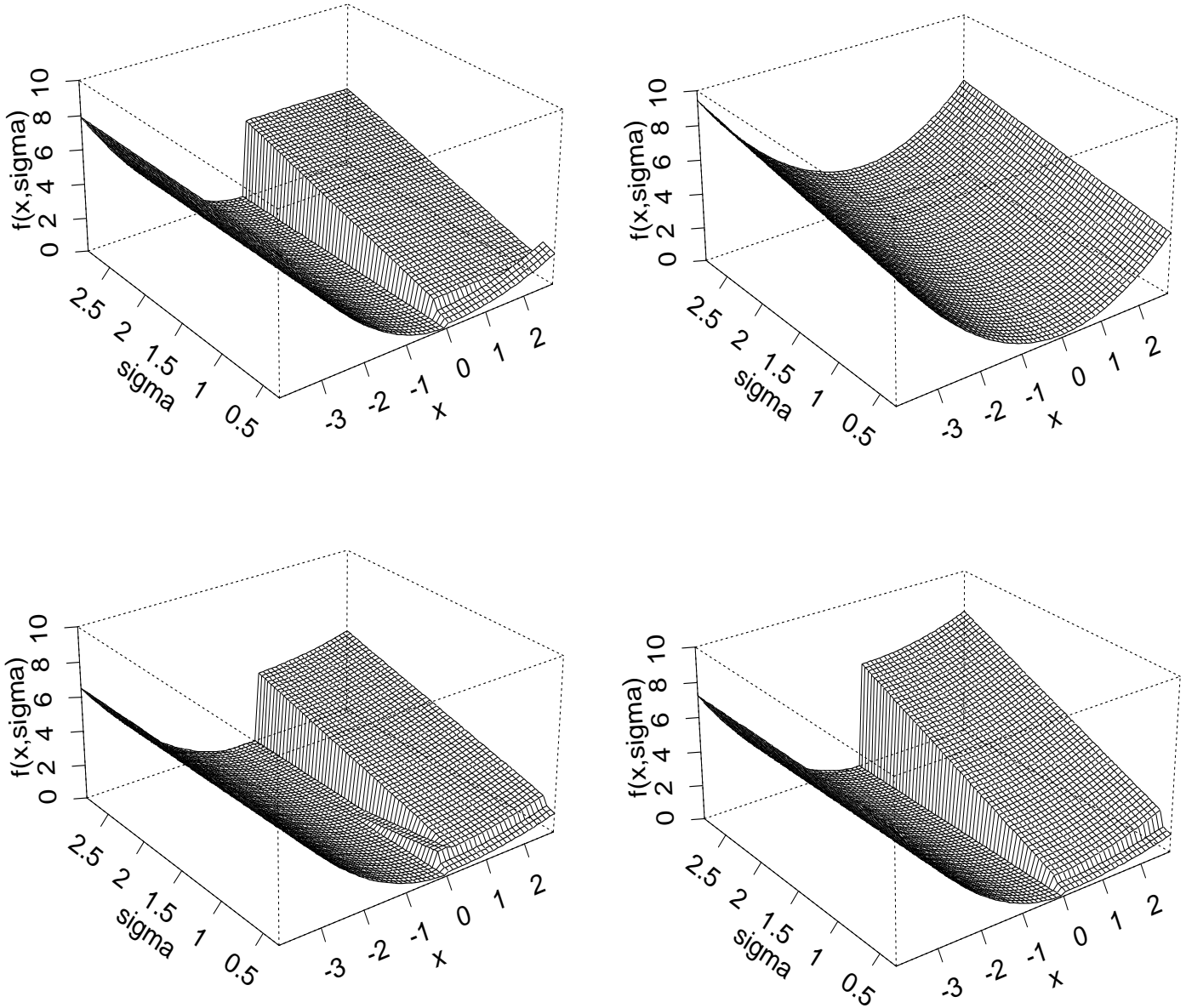


Figure 4.1: Top left: true conditional variance  $f(x, \sigma)$  given by (4.1), plotted against  $x$  and  $\sigma$ . This is used in model (2.1) with scaled  $t_6$ -distributed innovations for a data realization (data 4 from Table 4.1) which is at the basis for the other panels. Top right: estimated conditional variance from classical GARCH(1,1) model with scaled  $t_\nu$ -distributed innovations ( $\nu$  unknown). Bottom left: estimated conditional variance from tree GARCH model with standard normal innovations. Bottom right: estimated conditional variance from tree GARCH model with scaled  $t_\nu$ -distributed innovations ( $\nu$  unknown).

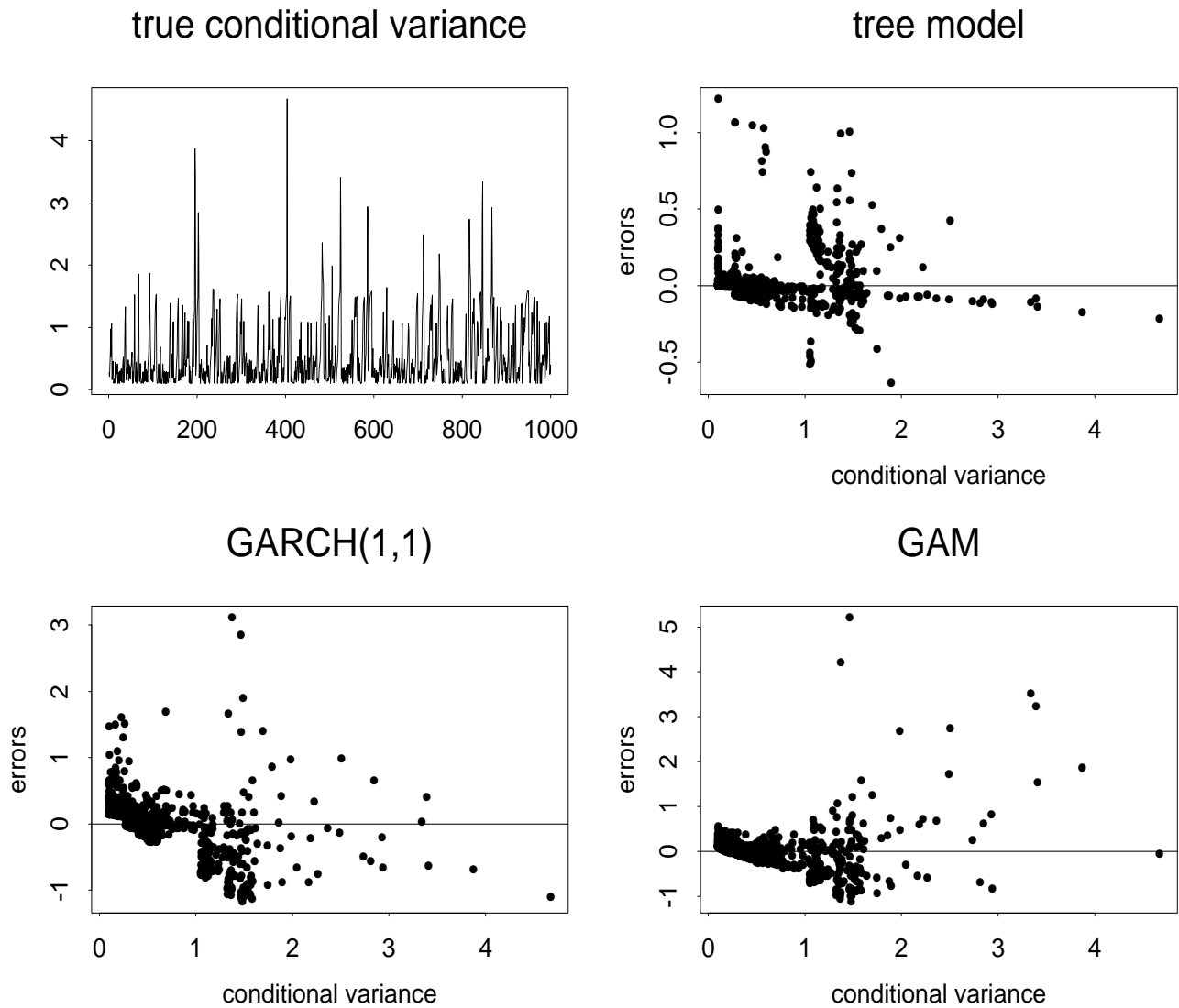
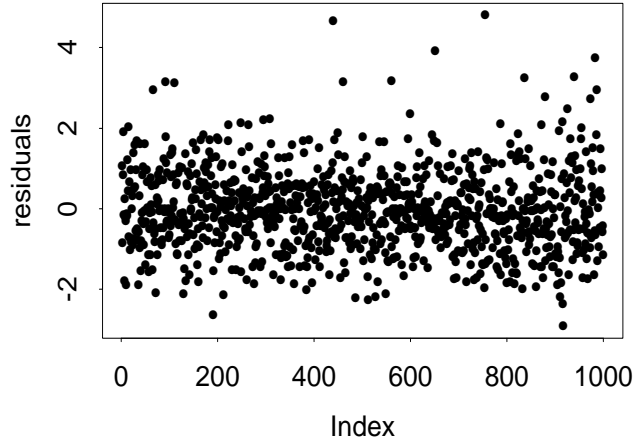
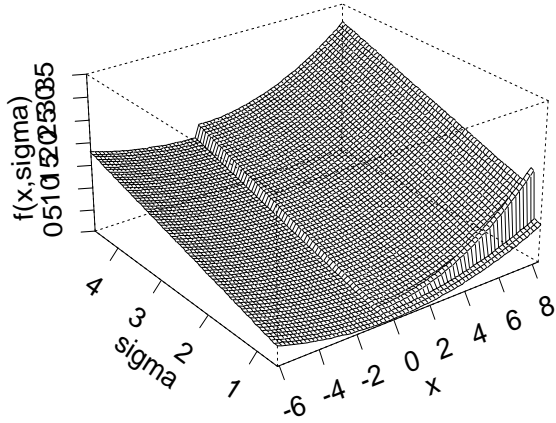


Figure 4.2: Top left: realization of volatility from model (2.1) with (4.1) and  $\mathcal{N}(0,1)$ -distributed innovations (data 2 from Table 4.1). Top right: errors  $\hat{E}_t = \hat{\sigma}_t^2 - \sigma_t^2$  from tree GARCH fit with scaled  $t_\nu$ -distributed innovations ( $\nu$  unknown) versus true conditional variance  $\sigma_t^2$ . Bottom left: errors  $\hat{E}_t$  from classical GARCH(1,1) fit with scaled  $t_\nu$ -distributed innovations ( $\nu$  unknown) versus  $\sigma_t^2$ . Bottom right: errors  $\hat{E}_t$  from GAM estimate (as described in the Appendix) versus  $\sigma_t^2$ .



Series : abs(residuals)

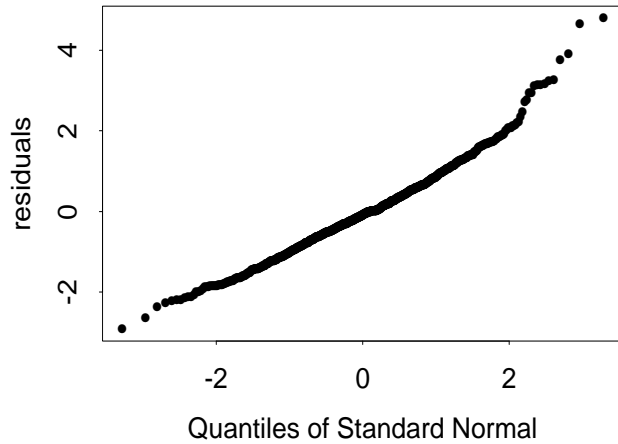
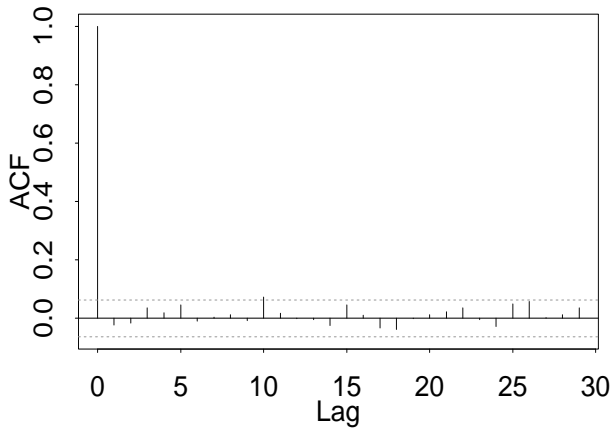


Figure 4.3: Results for negative log-returns of the DAX index from tree structured GARCH model with  $\mathcal{N}(0, 1)$ -distributed innovations. Top left: estimated function  $\hat{f}(x, \sigma)$  for the conditional variance, plotted against  $x$  and  $\sigma$ . Top right: residuals  $\hat{Z}_t = (X_t - \hat{\mu}_t) / \hat{\sigma}_t$  from the tree GARCH fit versus time  $t$ . Bottom left and right: autocorrelation function of the absolute residuals  $|\hat{Z}_t|$  and normal-plot for the residuals  $\hat{Z}_t$ , respectively.

Tree structured GARCH						
	Thresholds	AIC	L <sub>1</sub>	L <sub>2</sub>	OS-L <sub>1</sub>	OS-L <sub>2</sub>
data 1	$\hat{d}_1 = -0.044736$ in $x$ $\hat{d}_2 = 0.773178$ in $\sigma^2$ $\hat{d}_3 = 0.556751$ in $\sigma^2$	1842.201	102.8911	43.47956	110.6775 120.9755	44.39251 61.36875
data 2	$\hat{d}_1 = -0.016448$ in $x$ $\hat{d}_2 = 0.468750$ in $\sigma^2$ $\hat{d}_3 = 0.724830$ in $x$	1862.673	74.45144	25.22954	70.78948 79.80958	28.18894 43.02765
data 3	$\hat{d}_1 = -0.008402$ in $x$ $\hat{d}_2 = 0.188899$ in $x$ $\hat{d}_3 = 0.981620$ in $\sigma^2$	1885.969	101.6080	47.08707	86.50218 96.44229	32.48511 31.12207
average	–	1863.614	92.98351	38.59872	94.19942	40.09751
data 4	$\hat{d}_1 = -0.008463$ in $x$ $\hat{d}_2 = 0.385467$ in $\sigma^2$ $\hat{d}_3 = 0.313737$ in $x$	1743.976	114.0162	40.50913	131.1046 122.4830	65.67391 55.86674
	$\hat{d}_1 = -0.008463$ in $x$ $\hat{d}_2 = 0.406501$ in $\sigma^2$	1709.942	99.52060	50.98907	122.9608 101.5268	117.3457 70.69741
data 5	no thresholds	2935.83	85.12536	20.48578	67.41800 82.59603	9.245963 16.50112
data 6	$\hat{d}_1 = -0.871437$ in $x$ $\hat{d}_2 = 0.910214$ in $x$ $\hat{d}_3 = 1.452422$ in $\sigma^2$ $\hat{d}_4 = 1.706365$ in $\sigma^2$	3350.364	327.9303	301.0929	330.0598 354.0650	321.552 297.113
data 7	$\hat{d}_1 = -0.488943$ in $x$ $\hat{d}_2 = 0.845371$ in $x$	3371.991	252.4075	117.1972	265.7738 273.0419	128.3586 129.7451

Table 4.1: Estimated thresholds  $\hat{d}_i$  and goodness of fit measures for seven simulations (data 1–7). The likelihood for estimation is based on standard normal innovations (data 1–3, data 5–7); for data 4, we use standard normal innovation likelihood (upper part) and scaled  $t_\nu$  innovation likelihood with  $\nu$  unknown (lower part). Out-sample performances OS-L are evaluated at two independent test-sets.

	GARCH(1,1)				
	AIC	L <sub>1</sub>	L <sub>2</sub>	OS-L <sub>1</sub>	OS-L <sub>2</sub>
data 1	2009.407	225.6830	116.2419	248.5142 265.0486	150.1245 175.4757
data 2	2078.053	261.5570	159.4396	233.4011 273.3468	119.8525 178.9309
data 3	2073.761	284.3699	227.1664	243.6371 265.5155	140.5448 181.0429
average	2053.7403	257.2033	167.6160	254.9106	157.6619
data 4	1897.215	257.8995	154.1101	284.7026 280.5269	281.6552 2763654
	1818.409	265.6812	165.3186	287.8279 290.5641	348.4059 350.2534
data 5	2935.83	85.12536	20.48578	67.41800 82.59603	9.245963 16.50112
data 6	3375.200	236.4022	98.20690	476.3733 495.4368	515.8148 566.0920
data 7	3393.054	265.5798	136.6222	270.6098 267.3454	137.0235 138.8106

	GAM			
	L <sub>1</sub>	L <sub>2</sub>	OS-L <sub>1</sub>	OS-L <sub>2</sub>
data 1	200.0058	138.74708	225.2556 227.3836	152.1955 206.5867
data 2	208.8412	171.2537	197.0049 245.7308	168.2254 496.2520
data 3	250.1494	377.67642	231.0466 282.6617	348.6248 549.1308
average	219.6655	229.2257	234.8472	320.1692
data 4	219.9875	427.9997	274.9725 281.0249	459.5426 869.4288
data 5	549.5868	4352.951	416.3903 451.2216	987.4203 1028.240
data 6	315.8304	388.2322	259.9256 347.8295	152.1466 325.6214
data 7	227.7931	124.5836	229.5549 231.4537	132.9865 136.8912

Table 4.2: Goodness of fit measures for the same seven simulations from Table 4.1 using the GARCH(1,1) in (1.2) with  $\mathcal{N}(0, 1)$ -distributed innovations and the GAM model described in the Appendix.

Tree structured GARCH				
	Thresholds	AIC	IS-PL <sub>2</sub>	OS-PL <sub>2</sub>
DAX	$\hat{d}_1 = -0.354193$ in $x$ $\hat{d}_2 = 1.235410$ in $\sigma^2$ $\hat{d}_3 = 1.657356$ in $\sigma^2$	2776.238	8555.143	20001.30
	$\hat{d}_1 = -0.354193$ in $x$ $\hat{d}_2 = -0.889087$ in $x$ $\hat{d}_3 = 1.657356$ in $\sigma^2$ $\hat{d}_4 = -0.508199$ in $x$	2764.430	8507.640	20535.59
BMW	$\hat{d}_1 = -0.321663$ in $x$ $\hat{d}_2 = 1.110003$ in $\sigma^2$	3155.012	12059.22	15111.60

	GARCH(1,1)			GAM	
	AIC	IS-PL <sub>2</sub>	OS-PL <sub>2</sub>	IS-PL <sub>2</sub>	OS-PL <sub>2</sub>
DAX	2785.297	9190.588	20387.93	12588.24	24787.75
BMW	3165.068	12063.92	15110.95	10964.83	50789.17

Table 4.3: Estimated thresholds  $\hat{d}_i$  and goodness of fit measures for negative log-returns from the DAX index and the BMW stock price. The tree GARCH model (with  $\mathcal{N}(0, 1)$ -distributed innovations) is fitted with mesh = 8 (DAX, upper part; BMW) and mesh = 16 (DAX, lower part).