

## Efficient estimation in semiparametric GARCH models

Authors	Drost,F.C.; Klaassen,C.A.J.
Published in	Journal of Econometrics
Publication Date	1997
Link	<a href="https://research.tilburguniversity.edu/en/publications/c7de3f1c-c456-433e-a1c6-28ac16de51eb">https://research.tilburguniversity.edu/en/publications/c7de3f1c-c456-433e-a1c6-28ac16de51eb</a>
Citation	Drost, F C & Klaassen, C A J 1997, 'Efficient estimation in semiparametric GARCH models', Journal of Econometrics, vol. 81, no. 1, pp. 193-221.
Download Date	2026-05-18 22:05:23
Rights	<p>General rights</p> <p>Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.</p> <ul style="list-style-type: none"> <li>- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.</li> <li>- You may not further distribute the material or use it for any profit-making activity or commercial gain</li> <li>- You may freely distribute the URL identifying the publication in the public portal"</li> </ul> <p>Take down policy</p> <p>If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.</p>



ELSEVIER

Journal of Econometrics 81 (1997) 193–221

**JOURNAL OF  
Econometrics**

## Efficient estimation in semiparametric GARCH models

Feike C. Drost<sup>a,\*</sup>, Chris A.J. Klaassen<sup>b</sup>

<sup>a</sup> *Department of Econometrics, Tilburg University, P.O. Box 90153, 5000 LE Tilburg, Netherlands*

<sup>b</sup> *Department of Mathematics, University of Amsterdam, Plantage Muidergracht 24, 1018 TV Amsterdam, Netherlands*

---

### Abstract

It is well-known that financial data sets exhibit conditional heteroskedasticity. GARCH-type models are often used to model this phenomenon. Since the distribution of the rescaled innovations is generally far from a normal distribution, a semiparametric approach is advisable. Several publications observed that adaptive estimation of the Euclidean parameters is not possible in the usual parametrization when the distribution of the rescaled innovations is the unknown nuisance parameter. However, there exists a reparametrization such that the efficient score functions in the parametric model of the autoregression parameters are orthogonal to the tangent space generated by the nuisance parameter, thus suggesting that adaptive estimation of the autoregression parameters is possible. Indeed, we construct adaptive and hence efficient estimators in a general GARCH in mean-type context including integrated GARCH models.

Our analysis is based on a general LAN theorem for time-series models, published elsewhere. In contrast to recent literature about ARCH models we do not need any moment condition. © 1997 Elsevier Science S.A.

**Key words:** LAN in time-series; Semiparametrics; Adaptivity; (integrated) GARCH (in mean)

**JEL classification:** C22; C14

---

\* Corresponding author.

The first author is research fellow of the Royal Netherlands Academy of Arts and Sciences (K.N.A.W.). Part of this research was done by the second author at the Euler International Mathematical Institute, St. Petersburg. The authors gratefully acknowledge the helpful comments of Bas Werker, the editor Helmut Lütkepohl, and three anonymous referees.

## 1. Introduction

It is a well-established empirical fact in financial economics that time series like exchange rates and stock prices exhibit conditional heteroskedasticity. Big shocks are clustered together. The original paper of Engle (1982) proposes the ARCH model to incorporate conditional heteroskedasticity in econometric modeling of financial data sets. Bollerslev (1986) introduces the GARCH model as a generalization of ARCH. This facilitates a parsimonious parametrization which is particularly useful when shocks are important for a longer period (the idea corresponds to the generalization of AR to ARMA models). Several variations and extensions have been proposed in the literature. Nelson (1991) proposes the exponential GARCH model to capture the fact that the stock market is smoother in upward directions than in the opposite case (because of the leverage effect). Gouriéroux and Monfort (1992) suggest a nonparametric approach. They do not restrict attention to conditional variances that depend only upon past squared observations, but they try to estimate the functional form of the conditional heteroskedastic variance from the data. Another important extension is the GARCH-M-type model (cf. Engle et al., 1987). According to the Capital Asset Pricing Model one expects higher returns due to risk premia if the asset is more risky. To model this phenomenon the conditional variance is also included in the mean equation. Lots of applications have shown the strength of the GARCH type of modeling. In this paper we do not refer to original application papers but we want to draw attention to the monograph of Diebold (1988) and the survey paper of Bollerslev et al. (1992).

Despite the success of the GARCH history there are several topics that require attention. In this paper we consider the distributional assumptions on the rescaled innovations. The original formulations of GARCH-type models assume that these residuals are standard normal. Diebold (1988), however, shows that this assumption is often violated in empirical examples. Typically, the innovations have fat-tailed distributions and they are also non-symmetric in several applications. Drost and Werker (1996) provide an explanation for high kurtosis if the observations arise from a GARCH data-generating process in continuous time (see also Drost and Nijman, 1993; Nelson, 1990a). Diebold (1988) suggests that the errors will be 'more normal' if the process is more and more aggregated. Despite the observed non-normality of the error structure, Weiss (1986) and Lee and Hansen (1994) have shown that quasi-maximum likelihood estimation (QMLE), based upon the false assumption of normality, yields  $\sqrt{n}$ -consistent estimators: see also Lumsdaine (1989). However, the efficiency loss may be considerable. Therefore, several authors try to avoid efficiency loss, allowing the error structure to belong to some flexible parametric family of distributions. Student  $t$ -distributions are very popular, cf., e.g., Baillie and Bollerslev (1989). As a drawback of the introduction of such parametric models of the innovation distribution, we mention that the results of Weiss (1986) and Lee and Hansen (1994) do not carry over

to general error distributions. While QMLE based upon the normal distribution yields  $\sqrt{n}$ -consistent estimators, QMLE based upon other distributions (e.g., the student distributions) generally even fails to be consistent if the true distribution is different.

In the approaches mentioned above, the stochastic error structure is still described by some finite-dimensional statistical model. To avoid the introduction of a wrong parametric family of innovation distributions leading to inconsistent estimators and to be more flexible, a semiparametric approach is to be preferred. We want to estimate the conditional heteroskedastic character of the GARCH process but we do not want to restrict the class of error distributions too much. Apart from some regularity conditions we will assume the distribution of the innovations to be completely unknown. In passing we will consider the case of symmetrically distributed innovations. At first sight, these types of estimation problems seem to be much harder than the corresponding parametric ones and one would expect that optimal semiparametric estimators are less precise asymptotically than optimal parametric estimators. For lots of interesting econometric models this presumption turns out to be too pessimistic. Adaptive estimation is often possible. Adaptive estimation is just a special instance of semiparametric efficient estimation. Just as in parametric models, in semiparametric models an efficient estimator is an asymptotically normal estimator with minimal variance. If this minimal variance is the same as when the error distribution is known, one calls the efficient estimator adaptive since it adapts, so to say, to the underlying error distribution. Typically, an estimator based on a (wrongly) specified error distribution is not efficient in the semiparametric sense. For i.i.d. observations a lot of adaptive and semiparametric results are available [cf., e.g., Bickel et al. (1993) (BKRW, 1993 from now on) and the survey papers of Robinson (1988) and Newey (1990)]. Rigorous results are sparse in a time-series context. ARMA models are considered in detail by Kreiss (1987a,b). Some results for GARCH are obtained in Engle and González-Rivera (1991) and Steigerwald (1992) (see also Pötscher, 1995; Steigerwald, 1995). Linton (1993) discusses the semiparametric properties of ARCH models in more detail. However, these papers impose rather high moment conditions. The parameter estimates obtained in empirical work generally fail these moment conditions and, therefore, the scope for application seems to be limited.

In Drost et al. (1997) (henceforth DKW, 1997) a general LAN theorem for time-series models is presented together with conditions guaranteeing the existence of efficient estimators. We will apply these results to GARCH type models, including, e.g., I-GARCH and GARCH-M, thus avoiding severe moment conditions. We do not need the existence of moments neither of the rescaled innovations in the GARCH model (admitting, for example, Cauchy errors) nor of the observations (as is clear from the inclusion of integrated GARCH models). Since we only assume the existence of a stationary solution of the GARCH equations, our approach captures the models commonly used. A general LAN

theorem for time-series models is also contained in Theorem 13, Section 4, of Jeganathan (1995). Based hereon is his Theorem 17, Section 4, which yields adaptive estimators for ARMA-type location models. However, this result is not directly applicable to GARCH scale models and, moreover, it heavily leans on symmetry of the innovations.

To keep notation simple we restrict attention to the popular and most commonly used GARCH(1,1)-type models. This preserves the essential difficulty of GARCH (with respect to ARCH) since both the AR and the MA part are present in the conditional variance equation. All past observations show up (at an exponentially decaying rate). The statement of our theorem with respect to GARCH(1,1) is easily generalized to the general case of GARCH( $p, q$ ).

The first semiparametric results in a GARCH context were only partially successful. Engle and González-Rivera (1991) state “Monte Carlo evidence suggest that this semiparametric method (*i.e. the discrete maximum penalized likelihood estimation technique of Tapia and Thompson, 1978*) can improve the efficiency of the parameter estimates up to 50% over QMLE, but it does not seem to capture the total potential gain in efficiency. In this sense we say that the estimator is not adaptive in the class of densities with mean 0 and variance 1; that is, the estimator is not fully efficient, and it does not achieve the Cramér–Rao lower bound. The information matrix is not block-diagonal between the parameters of interest (the ones in the mean and in the variance equation) and the nuisance parameters (the knots of the density). If we choose the parametric form of the model with a conditional parametric density defined by a shape parameter, this one being part of the parameters to estimate, we can show easily that the expectation of the cross-partial derivatives of the log-likelihood function respects the parameter of interest and the shape parameter is different from 0. In other words, the estimation of the shape parameter affects the efficiency of the estimates of the parameters of interest” (pp. 355–356). These statements imply that the finite-dimensional parameter describing the GARCH model (with the standardized error distribution as nuisance parameter) is not adaptively estimable. This is not surprising since the classical GARCH formulation contains a scale parameter, and in most models the variance is not adaptively estimable. Therefore, the scale parameter is often included into the (infinite-dimensional) nuisance parameter. For the GARCH model this procedure does not work: the scores w.r.t. the remaining autoregression parameters are still not orthogonal to the tangent space. Hence, complete adaptive estimation of the conditional heteroskedastic character is not possible in GARCH models. This explains the efficiency loss observed by Engle and González-Rivera (1991). However, calculation of the scores w.r.t. the parameters of the GARCH model shows that there are several orthogonality relations between the score space and the tangent space generated by the unknown shape. Linton (1993) and Drost et al. (1994) (henceforth DKW, 1994) obtain along different lines a reparametrization of the ARCH and GARCH model, respectively, such that the autoregression parameters are adaptively estimable and

the location-scale parameters generate the most difficult one-dimensional sub-problems. So, knowledge of the shape of the error distribution does not help to construct better estimators of the conditional heteroskedastic character of the GARCH process. This resembles the regression model with unknown location  $\mu \in \mathbb{R}$ , regression parameter  $\beta \in \mathbb{R}^k$  and completely unknown error distribution, where the regression parameter  $\beta$  is adaptively estimable if the location parameter is included into the nuisance (see Bickel, 1982).

The paper is organized along the following lines. In Section 2 we state the LAN theorem for a large set of GARCH(1,1)-type models, including all stationary classical GARCH models such as, e.g., I-GARCH and GARCH-M. This LAN property is derived for the parametric model with the shape of the innovations known, and it implies the Convolution Theorem of Hájek (1970) which we will state next. This Convolution Theorem yields a bound on the asymptotic performance of estimators in the parametric model and is valid, a fortiori, for the semiparametric model as well. Section 3 is devoted to the construction of an estimator of the autoregression parameters on the assumption that the shape of the innovations is unknown, i.e. within the semiparametric model. This estimator happens to attain the bound from the parametric Convolution Theorem and therefore is asymptotically efficient in the parametric model and hence in the semiparametric model, since it does not use knowledge about the shape of the innovations. Such an estimator, which attains the parametric bound in a semiparametric model, is called adaptive. The proofs of these results are based on DKW(1997) and most of them are given in the appendix.

A small simulation study is presented in Section 4. It turns out that the suggested optimal estimator performs as expected: the estimator performs better than QMLE and the difference with MLE (if the error distribution is known) becomes negligible when the sample size is growing large. The empirical illustration in this section shows that the efficiency loss by using QMLE may be considerable. Some conclusions are drawn in Section 5.

## 2. LAN and Convolution Theorem

We consider a generalization of the reparametrized GARCH( $p, q$ ) model as given in Linton (1993), with  $p=0$ , and motivated by adaptation arguments in DKW(1994). For notational simplicity, we take  $p=q=1$ . In this manner the essential difficulty of an infinite number of lags is retained. To obtain the corresponding results for the general case (with  $p, q \in \mathbb{N}$  fixed) a careful replacement of coefficients by vectors suffices.

Let  $\mu \in \mathbb{R}$ ,  $\sigma > 0$ ,  $\alpha > 0$ , and  $\beta > 0$  be parameters and let  $\{\varepsilon_t: t \in \mathbb{Z}\}$  be an i.i.d. sequence of innovation errors with location zero, scale one, and density  $g$ . Put  $\xi_t = \mu + \sigma\varepsilon_t$  and note that  $\xi_t$  is a random variable with location  $\mu$ , scale  $\sigma$ , and density  $\sigma^{-1}g(\{\cdot - \mu\}/\sigma)$ . We introduce the following convention: random

variables, like  $\varepsilon$  and  $\xi$ , denote a typical element of the corresponding sequences  $\{\varepsilon_t: t \in \mathbb{Z}\}$  and  $\{\xi_t: t \in \mathbb{Z}\}$ .

Consider the model with observations

$$y_t = h_t^{1/2} \xi_t = \mu h_t^{1/2} + \sigma h_t^{1/2} \varepsilon_t, \quad (2.1)$$

where the unobservable heteroskedasticity factors  $h_t$  depend on the past via

$$h_t = 1 + \beta h_{t-1} + \alpha y_{t-1}^2 = 1 + h_{t-1}(\beta + \alpha \xi_{t-1}^2). \quad (2.2)$$

Observe that the Euclidean parameter  $\theta = (\alpha, \beta, \mu, \sigma)'$  is identifiable. Throughout we assume that Eq. (2.2) admits a stationary solution  $\{h_t: t \in \mathbb{Z}\}$ . A necessary and sufficient condition is given by (Nelson, 1990b, Theorem 2):

*Assumption A.*

$$E \ln\{\beta + \alpha \xi^2\} < 0. \quad (2.3)$$

Our semiparametric analysis treats the density  $g$  as an infinite-dimensional nuisance parameter and includes all strictly stationary GARCH models of type (2.1), (2.2). These equations contain, e.g., the classical Engle (1982)–Bollerslev (1986) GARCH model with a different parametrization and with finite second moments ( $\beta + \alpha \sigma^2 < 1$  and  $\mu = 0$ ), and the I-GARCH model of Engle and Bollerslev (1986) ( $\beta + \alpha \sigma^2 = 1$  and  $\mu = 0$ ). Furthermore, our model resembles the GARCH-M model of Engle et al. (1987). In the mean equation (2.1) we have included the conditional standard deviation of  $y_t$  while Engle et al. (1987) include a kind of conditional variance. More precisely stated, their model is given by  $z_t = \delta h_t + y_t$  and  $\mu = 0$ , i.e.,  $z_t = \delta h_t + \sigma h_t^{1/2} \varepsilon_t$ . Inserting  $\mu = 0$  in (2.1) or  $\mu = \delta = 0$  in the GARCH-M model yields the classical GARCH model. Generally, risk aversion is stronger pronounced in the original GARCH-M model than in our formulation.

Suppose that we observe  $y_1, \dots, y_n$ , and some starting value  $h_{01}$  initializing (2.2). It is not needed that  $h_{01}$  arises from the stationary solution of (2.2). We are considering estimation of  $\theta$ , based on  $h_{01}, y_1, \dots, y_n$ , in the presence of the infinite-dimensional nuisance parameter  $g$ . However, in this section we will fix the nuisance parameter  $g$  and in the resulting parametric model we will derive a bound on the asymptotic performance of regular estimators of  $\theta$ , a so-called Convolution Theorem. To that end we choose local submodels and we will study estimation of  $\theta$  locally asymptotically. The above-mentioned Convolution Theorem holds once the log-likelihood ratios of the observed random variables are locally asymptotically normal (LAN).

Observe that the model with the autoregression parameters  $\alpha$  and  $\beta$  fixed too, corresponds to the location-scale model for i.i.d. random variables since the information provided by the observations  $h_{01}, y_1, \dots, y_n$  is equal to the information

contained in the i.i.d. random variables  $\xi_1, \dots, \xi_n$ . Consequently, the location-scale model is a parametric submodel of our time-series model and it makes sense to assume that this submodel is regular, i.e. (see Hájek and Šidák, 1967):

*Assumption B.* The distribution of  $\varepsilon$  possesses an absolutely continuous Lebesgue density  $g$  with derivative  $g'$  and finite Fisher information for location

$$I_l(g) = \int \{g'/g\}^2 g(\varepsilon) d\varepsilon \tag{2.4}$$

and for scale

$$I_s(g) = \int \{1 + \varepsilon g'/g(\varepsilon)\}^2 g(\varepsilon) d\varepsilon. \tag{2.5}$$

Moreover, the random variable  $\varepsilon$  has location zero and scale one.

To be able to derive an asymptotic lower bound we have to rely on semi-parametric methods as presented in, e.g., BKRW(1993) and DKW(1994, 1997). So we fix  $\theta$  at  $\theta_0 = (\alpha_0, \beta_0, \mu_0, \sigma_0)'$  and choose local parametrizations  $\theta_n = (\alpha_n, \beta_n, \mu_n, \sigma_n)'$  and  $\tilde{\theta}_n = (\tilde{\alpha}_n, \tilde{\beta}_n, \tilde{\mu}_n, \tilde{\sigma}_n)'$  such that  $|\theta_n - \theta_0| = O(n^{-1/2})$ ,  $|\tilde{\theta}_n - \theta_0| = O(n^{-1/2})$ , and even

$$\lambda_n = \sqrt{n}(\tilde{\theta}_n - \theta_n) \rightarrow \lambda, \quad \text{as } n \rightarrow \infty. \tag{2.6}$$

In the remainder expectations, convergences, etc., are implicitly taken under  $\theta_n$  and  $g$  (unless otherwise indicated).

To obtain a uniform LAN theorem we consider the log-likelihood ratio  $A_n$  of  $h_{01}, y_1, \dots, y_n$  for  $\tilde{\theta}_n$  with respect to  $\theta_n$  under  $\theta_n$  (and  $g$  fixed). Observe that the residuals and the conditional variances up to time  $t$  can be recursively calculated from  $\theta$  and the observations  $h_{01}, y_1, \dots, y_t$ : with  $h_1(\theta) = h_{01}$ , obtain for  $t = 1, 2, \dots$

$$\xi_t(\theta) = y_t/h_t^{1/2}(\theta), \tag{2.7}$$

$$\varepsilon_t(\theta) = \{\xi_t(\theta) - \mu\}/\sigma, \tag{2.8}$$

$$h_{t+1}(\theta) = 1 + \beta h_t(\theta) + \alpha y_t^2. \tag{2.9}$$

Conditionally on  $h_{01}$  the density of  $y_1, \dots, y_n$  under  $\theta_n$  is

$$\begin{aligned} \prod_{t=1}^n \sigma_n^{-1} h_{nt}^{-1/2} g(\sigma_n^{-1} \{h_{nt}^{-1/2} y_t - \mu_n\}) &= \prod_{t=1}^n \sigma_n^{-1} h_{nt}^{-1/2} g(\{\xi_{nt} - \mu_n\}/\sigma_n) \\ &= \prod_{t=1}^n \sigma_n^{-1} h_{nt}^{-1/2} g(\varepsilon_{nt}), \end{aligned}$$

where  $h_{nt} = h_t(\theta_n)$ ,  $\xi_{nt} = \xi_t(\theta_n)$ , and  $\varepsilon_{nt} = \varepsilon_t(\theta_n)$ .

To enhance the interpretation of this formula and to stress the link between the present time-series model and the i.i.d. location-scale model we introduce the notation  $\tilde{h}_{nt} = h_t(\tilde{\theta}_n)$ ,

$$\begin{aligned} \ell\{\mu, \sigma\}(x) &= \log g(\{x - \mu\}/\sigma) - \log \sigma, \\ \begin{pmatrix} M_{nt} \\ S_{nt} \end{pmatrix} &= n^{1/2} \sigma_n^{-1} h_{nt}^{-1/2} \begin{pmatrix} \tilde{\mu}_n \tilde{h}_{nt}^{1/2} - \mu_n h_{nt}^{1/2} \\ \tilde{\sigma}_n \tilde{h}_{nt}^{1/2} - \sigma_n h_{nt}^{1/2} \end{pmatrix}, \end{aligned} \tag{2.10}$$

and  $\tilde{\varepsilon}_{nt} = \varepsilon_t(\tilde{\theta}_n)$ . With  $A_n^s$  the log-likelihood ratio for  $h_{01}$ , the log-likelihood ratio  $A_n$  may be written as

$$\begin{aligned} A_n &= \log \left\{ \prod_{t=1}^n \tilde{\sigma}_n^{-1} \tilde{h}_{nt}^{-1/2} g(\tilde{\varepsilon}_{nt}) \middle/ \prod_{t=1}^n \sigma_n^{-1} h_{nt}^{-1/2} g(\varepsilon_{nt}) \right\} + A_n^s \\ &= \sum_{t=1}^n \{ \ell\{(\mu_n, \sigma_n) + \sigma_n n^{-1/2} (M_{nt}, S_{nt})\}(\tilde{\varepsilon}_{nt}) - \ell\{\mu_n, \sigma_n\}(\tilde{\varepsilon}_{nt}) \} + A_n^s \\ &= \sum_{t=1}^n \{ \ell\{(0, 1) + n^{-1/2} (M_{nt}, S_{nt})\}(\varepsilon_{nt}) - \ell\{0, 1\}(\varepsilon_{nt}) \} + A_n^s. \end{aligned} \tag{2.11}$$

This expression resembles the log-likelihood ratio statistic for the i.i.d. location-scale model but here the deviations  $M_{nt}$  and  $S_{nt}$  are random. In the i.i.d. case the LAN theorem is obtained with deterministic sequences. We will apply the results of DKW(1997) which allow for such random sequences.

To get rid of the starting condition in the log-likelihood ratio statistic we will use the following regularity condition [compare assumption (A.3) of Kreiss (1987a) and Assumption A of DKW(1997)].

*Assumption C.* The density  $\bar{g}_\theta$  of the initial value  $h_{01}$  satisfies, under  $\theta_n$ ,

$$A_n^s = \log\{\bar{g}_{\tilde{\theta}_n}/\bar{g}_{\theta_n}(h_{01})\} \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty. \tag{2.12}$$

To make an appropriate expansion of  $A_n$  it will be handy to introduce the notation  $\dot{\ell}_{nt}$  for the four-dimensional conditional score at time  $t$ . To be more precise, denote the two-dimensional vector derivative of the conditional variance by

$$H_t(\theta) = \frac{\partial}{\partial(\alpha, \beta)} h_t(\theta) = \beta H_{t-1}(\theta) + \begin{pmatrix} y_{t-1}^2 \\ h_{t-1}(\theta) \end{pmatrix}, \tag{2.13}$$

with  $H_1(\theta) = 0_2$ . Define the  $(4 \times 2)$ -derivative matrix  $W_t(\theta)$  [motivated by differentiation of  $(M_{nt}, S_{nt})$  with respect to  $\tilde{\theta}_n$  at  $\theta_n$ ] by

$$W_t(\theta) = \sigma^{-1} \begin{pmatrix} \frac{1}{2} h_t^{-1}(\theta) H_t(\theta)(\mu, \sigma) \\ I_2 \end{pmatrix}, \tag{2.14}$$

denote the location-scale score by (with  $l' = g'/g$ )

$$\psi_t(\theta) = - \left( \begin{array}{c} \ell'(\varepsilon_t(\theta)) \\ 1 + \varepsilon_t(\theta)\ell'(\varepsilon_t(\theta)) \end{array} \right), \tag{2.15}$$

and put

$$\dot{\ell}_t(\theta) = W_t(\theta)\psi_t(\theta).$$

Then, the conditional score at time  $t$  may be denoted by  $\dot{\ell}_{nt} = \dot{\ell}_t(\theta_n)$ . Observe that  $\dot{\ell}$  is just the heuristic score. An expansion of (2.11) shows that the log-likelihood ratio  $A_n$  may be alternatively written as

$$A_n = \lambda' n^{-1/2} \sum_{t=1}^n \dot{\ell}_{nt} - \frac{1}{2} n^{-1} \sum_{t=1}^n \{\lambda' \dot{\ell}_{nt}\}^2 + R_n. \tag{2.16}$$

The LAN result for the parametric version of model (2.1), (2.2) is stated in the following theorem. The proof is deferred to Appendix A.

**Theorem 2.1 (LAN).** *Suppose that Assumptions A–C are satisfied. Then the local log-likelihood ratio statistic  $A_n$ , as defined by (2.11) and (2.16), is asymptotically normal. More precisely, under  $\theta_n$ ,*

$$R_n \xrightarrow{P} 0, \quad A_n \xrightarrow{D} N\left(-\frac{1}{2} \lambda' I(\theta_0) \lambda, \lambda' I(\theta_0) \lambda\right) \quad \text{as } n \rightarrow \infty, \tag{2.17}$$

where  $I(\theta_0)$  is the probability limit of the averaged score products  $\dot{\ell}_{nt}' \dot{\ell}_{nt}$ .

We are now in a position to apply the Convolution Theorem of Hájek (1970); cf. Theorem 2.3.1 of BKRW (1993, p. 24).

**Theorem 2.2 (Convolution Theorem).** *Under the assumptions of the LAN Theorem 2.1, let  $\{T_n : n \in \mathbb{N}\}$  be a regular sequence of estimators of  $q(\theta)$ , where  $q : \mathbb{R}^4 \rightarrow \mathbb{R}^k$  is differentiable with total differential matrix  $\dot{q}$ . As usual, regularity at  $\theta = \theta_0$  means that there exists a random  $k$ -vector  $Z$  such that for all sequences  $\{\theta_n : n \in \mathbb{N}\}$ , with  $n^{1/2}(\theta_n - \theta_0) = O(1)$ ,*

$$n^{1/2} \{T_n - q(\theta_n)\} \xrightarrow{D} Z \quad \text{as } n \rightarrow \infty, \tag{2.18}$$

where the convergence is under  $\theta_n$ . Let  $\tilde{\ell} = \dot{q}(\theta_0)I(\theta_0)^{-1}\dot{\ell}(\theta_0)$  be the efficient influence function, then, under  $\theta_0$ ,

$$\left( \begin{array}{c} n^{1/2} \{T_n - q(\theta_0) - n^{-1} \sum_{i=1}^n \tilde{\ell}_i\} \\ n^{-1/2} \sum_{i=1}^n \tilde{\ell}_i \end{array} \right) \xrightarrow{D} \begin{pmatrix} \Delta_0 \\ Z_0 \end{pmatrix} \quad \text{as } n \rightarrow \infty, \tag{2.19}$$

where  $\Delta_0$  and  $Z_0$  are independent and  $Z_0$  is  $N(0, \overset{\circ}{q}(\theta_0)I(\theta_0)^{-1}\overset{\circ}{q}(\theta_0)')$ . Moreover,  $\{T_n : n \in \mathbb{N}\}$  is efficient if  $\{T_n : n \in \mathbb{N}\}$  is asymptotically linear in the efficient influence function, i.e. if  $\Delta_0 = 0$  (a.s.).

As a conclusion from the Convolution Theorem we obtain that a regular estimator  $\hat{\theta}_n$  of  $\theta$  satisfies, under  $\theta_0$ ,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{D} \Delta_0 + Z_0,$$

i.e. the limit distribution of  $\hat{\theta}_n$  is the convolution of the random vector  $\Delta_0$  and a Gaussian random vector with mean zero and variance the inverse of the information matrix  $I(\theta_0)$ . Since  $\Delta_0$  adds noise to the Gaussian vector  $Z_0$ , it is clear that  $\Delta_0 = 0$  would be preferred. This motivates the usual terminology (as lower bound, etc.) because  $\Delta_0 = 0$  is attainable in the majority of situations.

In the remainder of this paragraph we simplify exposition by supposing that the scores given above are stationary such that we may restrict attention to just one specific element; compare DKW(1994). In this way it is easier to comprehend the specific adaptiveness features in the GARCH model. These results are derived along the lines of Sections 2.4 and 3.4 of BKRW (1993). This expository simplification will be suppressed again in the next section when deriving a (semiparametric) efficient estimator. This optimal estimator satisfies the properties obtained in (I)–(IV) below.

In a stationary setting the Fisher information matrix defined in the LAN Theorem 2.1 simplifies to

$$I(\theta_0) = E\ell'\ell' = E W \psi \psi' W' = E W I_s(g) W',$$

where  $I_s(g)$  is the information matrix in the location-scale model,

$$I_s(g) = E \psi \psi' = \begin{pmatrix} E(\ell')^2 & E\epsilon(\ell')^2 \\ E\epsilon(\ell')^2 & E(1 + \epsilon\ell')^2 \end{pmatrix}.$$

If the location parameter  $\mu$  is known to be zero, as in the classical GARCH case, this formula simplifies even further to

$$I(\theta_0) = I_s(g) E W_s W_s', \tag{2.20}$$

where  $W_s$  is the three-dimensional subvector of  $W$  concerning the relevant derivatives with respect to the scale parameter  $\sigma$  and where  $I_s(g) = E(1 + \epsilon\ell')^2$  is the information for scale in the i.i.d. scale model.

(I) If  $g$  is known and if we want to estimate the autoregression parameter  $v = (\alpha, \beta)'$  in the presence of the nuisance parameter  $\eta = (\mu, \sigma)'$  then we see that the efficient influence function, as defined in the Convolution Theorem 2.2, equals

$$\tilde{\ell} = (I_2, 0_{2 \times 2}) [E\ell'\ell']^{-1} \ell' = (I_2, 0_{2 \times 2}) I(\theta_0)^{-1} \ell'.$$

As in Proposition 2.4.1.A and formula (2.4.3) of BKRW (1993, pp. 28,30) we may write

$$\tilde{\ell} = [E\ell_1^* \ell_1^{*'}]^{-1} \ell_1^*, \tag{2.21}$$

where the so-called efficient score function  $\ell_1^*$  of  $v$  is obtained by the componentwise projection of  $\ell_1$ , the first two elements of  $\ell$ , onto the orthocomplement of  $[\ell_2]$ , the linear span of the last two components of  $\ell$ . Here the inner product is the covariance and the orthocomplement is taken in the linear space spanned by all components of  $\ell$ . It is easy to verify that

$$\ell_1^* = \frac{1}{2} \sigma^{-1} \{ (H/h) - E(H/h) \} (\mu, \sigma) \psi \tag{2.22}$$

and that  $\ell_1^*$  is orthogonal to  $\ell_2$  indeed, since  $H/h = H_t/h_t$  depends on the past only and is independent of the present innovation  $\varepsilon_t$ .

(II) If  $g$  is unknown and if we want to estimate  $v$  in the presence of the nuisance parameters  $\eta$  and  $g$  then we obtain the same efficient influence function. To see this note that the components of  $\ell_1^*$  as given in (2.22) are orthogonal to every element of  $L_2^0(\varepsilon)$  by the independence of present ( $\xi$  and  $\varepsilon$ ) and past ( $h$  and  $H$ ). By (3.4.2) and Corollary 3.4.1.A of BKRW (1993, pp. 70,72) we obtain

$$I(P_0 | v, \mathcal{Q}) \geq E\ell_1^* \ell_1^{*'} \tag{2.23}$$

for all regular parametric submodels  $\mathcal{Q}$  of our semiparametric model  $\mathcal{P}$ , i.e. the information at  $P_0$  in estimation of  $v$  within the parametric submodel  $\mathcal{Q}$  equals at least the information at  $P_0$  in estimation of  $v$  within the parametric model, studied in (I), with  $g$  known. In other words, as far as estimation of  $v$  is concerned, no parametric model  $\mathcal{Q}$  is asymptotically more difficult to first order (contains less information) than the model from (I). Consequently, the semiparametric model  $\mathcal{P}$  itself is asymptotically to first order as difficult as the parametric model with  $g$  known, i.e. the information matrix with respect to  $v$  evaluated at  $P_0$  for the semiparametric model  $\mathcal{P}$  equals the lower bound in the parametric model with  $g$  known (case (I)),

$$I(P_0 | v, \mathcal{P}) = E\ell_1^* \ell_1^{*'}.$$

Once more, the efficient influence function is given by (2.21). Apparently, introduction of the nuisance parameter  $g$  in the presence of the Euclidean nuisance parameter  $\eta$  does not change the efficient influence function for  $v$ . Hence, estimation of  $v$  is asymptotically as hard not knowing  $g$  as knowing  $g$ . One usually calls this *adaptivity*. Observe, however, that the presence of the nuisance parameter  $\eta$  is important to derive this result. If  $\eta$  is known adaptive estimation of  $v$  is not possible! The same conclusion applies if  $\eta$  is included into the ‘big’ infinite dimensional nuisance parameter  $g$ . So, the nuisance parameter  $\eta$  is

treated in another way than the nuisance parameter  $g$ . Since location-scale parameters are almost always present in econometric models a different treatment is not unreasonable and the usage of the protected notion ‘adaptivity’ is legalized. However, with the comments above in mind, a more appropriate way of saying this is to call the parameter  $\nu$   $\eta$ -adaptive, explicitly referring to the remaining nuisance parameters present in the model. [Of course, a similar remark applies to, e.g., the non-symmetric regression model as discussed in Bickel (1982), where the regression parameter  $\beta$  is not fully adaptively estimable. In fact  $\beta$  is  $\mu$ -adaptive.]

(III) Estimation of the remaining parameter  $\eta$  is completely analogous to the location-scale problem for i.i.d. variables. Obtain the well-known lower bound for  $\eta$  in the semiparametric location-scale model. It suffices to construct a sequence of estimators  $\{\hat{\eta}_n, n \in \mathbb{N}\}$  for  $\eta$  attaining this bound. Let  $\hat{\theta}_n$  be some initial  $\sqrt{n}$ -consistent estimator of  $\theta$ , calculate  $\hat{h}_{nt} = h_t(\hat{\theta}_n)$  by plugging in  $\hat{\theta}_n$  into (2.9) and obtain the residuals  $\hat{\xi}_{nt} = \xi_t(\hat{\theta}_n) = y_t / \hat{h}_{nt}^{1/2}$ , similarly. If one proceeds as if the  $\hat{\xi}_{nt}$  are i.i.d. observations from some location-scale model, one obtains a semiparametric efficient estimator for  $\eta$  in our model (as is easily verified from the Convolution Theorem 2.2 by choosing an appropriate function  $q$ ). To be more explicit, we assume that  $g$  has finite second moment and we define the location and scale parameters by standardizing  $g$  via the equations  $E_g \varepsilon = 0, E_g \varepsilon^2 = 1$ . Then the square root of the sample variance is optimal for  $\sigma$  both in the symmetric and non-symmetric case. The sample mean is optimal for  $\mu$  if no symmetry is assumed and under the assumption of symmetry one has to use an efficient estimator for the symmetric location-problem (cf. Example 7.8.1 of BKRW, 1993, p. 400). If one wants to avoid moment conditions on  $\varepsilon$  one may define the location-scale parameter in another way, see the discussion of the M-estimator in Section 3.

(IV) Finally, when estimating the whole Euclidean parameter  $\theta$ , the efficient score is simply obtained from (II) and (III). Following the arguments leading to (2.23) in (II) this score function yields a lower bound indeed. Optimality of this bound follows from (III) by choosing the most difficult direction from the location-scale problem.

Obvious substitutions in Theorems 2.1 and 2.2 show that the conclusions above are also valid for the classical GARCH model with  $\mu = E_g \xi = 0$ . An optimal estimator of  $\sigma$  in the non-symmetric case is given then by the square root of (cf. Example 3.2.3 of BKRW, 1993, pp. 53–55)

$$n^{-1} \sum_{t=1}^n \hat{\xi}_{nt}^2 - n^{-1} \frac{\sum_{t=1}^n \hat{\xi}_{nt}^3}{\sum_{t=1}^n \hat{\xi}_{nt}^2} \sum_{t=1}^n \hat{\xi}_{nt}.$$

In the symmetric case the limiting behavior of this estimator and the square root of the sample variance are the same.

### 3. Adaptive Estimators

In classical parametric models the maximum likelihood estimator is asymptotically efficient, typically. In semiparametric models such an estimation principle yielding efficient estimators does not exist. However, there exist methods to upgrade  $\sqrt{n}$ -consistent estimators to efficient ones by a Newton–Raphson technique, provided it is possible to estimate the relevant score or influence functions sufficiently accurately. In Klaassen (1987) such a method based on ‘sample splitting’ is described for i.i.d. models. Schick (1986) uses both ‘sample splitting’ and Le Cam’s ‘discretization’, again in i.i.d. models. See, e.g., Section 7.8 of BKRW (1993) for details. Schick’s (1986) method has been adapted to time-series models in Theorem 3.1 of DKW(1997). We assume the existence of such a preliminary,  $\sqrt{n}$ -consistent estimator.

*Assumption D.* There exists a  $\sqrt{n}$ -consistent estimator  $\hat{\theta}_n$  of  $\theta_n$  (under  $\theta_n$  and  $g$ ).

For our GARCH model a natural candidate for such an initial estimator is the MLE based on the assumption of normality of the innovations  $\varepsilon_t$ . One often calls this estimator the quasi-MLE. Probably, this QMLE is  $\sqrt{n}$ -consistent under every density  $g$  with  $E_g \varepsilon^4 < \infty$ ; this has been shown by Weiss (1986) for ARCH models and under restrictions by Lee and Hansen (1994) for GARCH models, which are slightly different from ours; see also Lumsdaine (1989). The additional moment condition on  $\varepsilon$  is needed there since a quadratic term appears in the score function of the scale parameter. To avoid moment conditions altogether, one could use, e.g., another preliminary M-estimator, instead. Let  $\chi: \mathbb{R} \rightarrow \mathbb{R}^2$  be a sufficiently smooth bounded function with monotonicity properties. As an example we mention  $\chi = (\chi_1, \chi_2)'$  with

$$\chi_1(x) = \frac{2}{1 + \exp\{-x\}} - 1, \quad x \in \mathbb{R},$$

the location score function for the logistic distribution and

$$\chi_2(x) = \int_0^x 2y \frac{\exp\{-y\}}{(1 + \exp\{-y\})^2} dy - 1, \quad x \in \mathbb{R}.$$

The M-estimator will solve the equations [cf. (2.7)–(2.9) and (2.13), (2.14)]

$$\sum_{i=1}^n W_i(\theta) \chi(\varepsilon_i(\theta)) = 0. \tag{3.1}$$

Use of this M-estimator implies that one standardizes  $g$  at location 0 and scale 1 by the equation  $E_g \chi(\varepsilon) = 0$ ; in the normal case with QMLE this yields  $\mu$  as expectation and  $\sigma$  as standard deviation.

To prove that estimation via (3.1) shows validity of Assumption D we have to prove existence of this M-estimator and its  $\sqrt{n}$ -consistency. It should be possible to show the existence along the lines of Scholz (1971) by studying the 4 by 4 pseudo-information matrix  $EW\chi\chi'W'$ ; see also Huber (1981, pp. 138–139). Here we will not attempt to do this, since the situation is much more complicated than the location-scale problem studied in the literature. At the cost of some generality we suppose here that  $\sqrt{n}$ -consistent estimators  $\hat{\alpha}_n$  and  $\hat{\beta}_n$  are given. The  $\sqrt{n}$ -consistency of  $\hat{\alpha}_n$  and  $\hat{\beta}_n$  together with the contiguity obtained from the LAN Theorem 2.1 implies that we may treat the parameters  $\alpha$  and  $\beta$  as given. So, we are in fact in the i.i.d. location-scale model and the M-estimators for  $\mu$  and  $\sigma$  solving the latter two equations in (3.1) are  $\sqrt{n}$ -consistent; see Huber (1981) and Bickel (1982). We conjecture that the proof of the more general M-estimator solving (3.1) can be given along similar lines.

Here we will focus on efficient and hence adaptive estimation of the autoregression parameters  $\alpha$  and  $\beta$  (cf. (I)–(IV) of Section 2); alternatively, in view of (2.14), note that the score  $\dot{\ell}_{nt}$  satisfies the form discussed in Example 3.1 of DKW(1997). In the appendix we verify the conditions of Theorem 3.1 in DKW(1997), this yields the following theorem.

**Theorem 3.1.** *Under Assumptions A–D adaptive estimators of  $\alpha$  and  $\beta$  do exist.*

To describe our adaptive estimator more accurately, let  $\hat{\theta}_n = (\hat{\alpha}_n, \hat{\beta}_n, \hat{\mu}_n, \hat{\sigma}_n)'$  be a  $\sqrt{n}$ -consistent estimator of  $\theta$  and compute  $W_t(\hat{\theta}_n)$  via (2.13) and (2.14). Let  $\hat{\varepsilon}_{n1}, \dots, \hat{\varepsilon}_{nn}$  be the residuals computed from  $h_1, y_1, \dots, y_n$  and  $\hat{\theta}_n$  using (2.8). Via a kernel estimate based on  $\hat{\varepsilon}_{n1}, \dots, \hat{\varepsilon}_{nn}$  with the logistic kernel, say  $k(\cdot)$ , and bandwidth  $b_n$  we estimate  $g(\cdot)$  by

$$\hat{g}_n(\cdot) = \frac{1}{n} \sum_{t=1}^n \frac{1}{b_n} k\left(\frac{\cdot - \hat{\varepsilon}_{nt}}{b_n}\right)$$

and subsequently  $\psi(\cdot)$  by  $\hat{\psi}_n(\cdot)$ ; here  $b_n \rightarrow 0$  and  $nb_n^4 \rightarrow \infty$ . Now our estimator may be written as

$$\begin{aligned} & (\hat{\alpha}_n, \hat{\beta}_n)' + (I_2, 0_{2 \times 2}) \left( \frac{1}{n} \sum_{t=1}^n W_t(\hat{\theta}_n) \hat{\psi}_n(\hat{\varepsilon}_{nt}) \hat{\psi}_n(\hat{\varepsilon}_{nt})' W_t(\hat{\theta}_n)' \right)^{-1} \\ & \times \frac{1}{n} \sum_{t=1}^n \left\{ W_t(\hat{\theta}_n) - \frac{1}{n} \sum_{s=1}^n W_s(\hat{\theta}_n) \right\} \hat{\psi}_n(\hat{\varepsilon}_{nt}). \end{aligned} \tag{3.2}$$

With  $\hat{\theta}_n$  the QMLE this is the estimator used in the simulations of Section 4. To prove that such estimators are adaptive we need the following two technical modifications.

**Discretization:**  $\hat{\theta}_n$  is discretized by changing its value in  $(0, \infty) \times (0, \infty) \times \mathbb{R} \times (0, \infty)$  into (one of) the nearest point(s) in the grid  $(c/\sqrt{n})(\mathbb{N} \times \mathbb{N} \times \mathbb{Z} \times \mathbb{N})$ .

This technical trick enables one to consider  $\hat{\theta}_n$  to be non-random, and therefore independent of  $\hat{\varepsilon}_{nt}$ ,  $y_t$ , and  $h_t$ .

*Sample splitting:* The set of residuals  $\hat{\varepsilon}_{n1}, \dots, \hat{\varepsilon}_{nn}$  is split into two samples, which may be viewed as independent now. For  $\hat{\varepsilon}_{nt}$  in the first sample, the second sample is used to estimate  $\psi(\cdot)$  by  $\hat{\psi}_{n2}(\cdot)$  and  $\hat{\psi}_n(\hat{\varepsilon}_{nt})$  in (3.2) is replaced by  $\hat{\psi}_{n2}(\hat{\varepsilon}_{nt})$ . Similarly for  $\hat{\varepsilon}_{nt}$  in the second sample, the first sample is used to estimate  $\psi(\cdot)$ . In this way, again some independence is introduced artificially to make the proof work.

This approach has been adopted in DKW(1997). It should be emphasized that both tricks are merely introduced as a technical device to make proofs work. Other approaches have also been studied in the literature. Klaassen (1987) has shown that discretization may be avoided at the cost of an extra sample splitting. Schick (1986) and Koul and Schick (1995) show that sample splitting may be avoided at the cost of some extra conditions.

#### 4. Simulations and an empirical example

To enhance the interpretation and validity of the theoretical results of the previous sections we present a small simulation experiment. Furthermore, a case study concerning some exchange rate series is given.

We simulated several GARCH(1,1) series of length  $n = 1000$ , parameters  $(\alpha, \beta, \sigma) = (0.3, 0.6, 1)$ ,  $(0.1, 0.8, 1)$ , and  $(0.05, 0.9, 1)$  (the parameter  $\mu$  is set to zero and is not estimated to allow for a better comparison with previous simulation studies), and eight different innovation distributions: normal, a balanced mixture of two standard normals with means 2 and  $-2$ , respectively, double exponential, Student's distributions with  $\nu = 5, 7$ , and 9 degrees of freedom, and (skew) chi-squared distributions with  $\nu = 6$  and 12 degrees of freedom. These densities are rescaled such that they have the required zero mean and unit variance.

It is the purpose of the simulations to evaluate the moderate sample properties of the autoregression parameters  $\alpha$  and  $\beta$  which are adaptively estimable, in principle. For each series we estimated these parameters with MLE, QMLE, and a one-step semiparametric procedure. For the latter estimation method we made two estimates: one under general assumptions on the innovation distribution and one under the extra assumption of symmetry. The theoretical results imply that there should be no difference between these two semiparametric methods if the true underlying density is symmetric indeed but small sample properties may differ. In the semiparametric part we used standardized logistic kernels with a bandwidth of  $h = 0.5$ . Reasonable changes of the bandwidth, say  $0.25 \leq h \leq 0.75$ , or another kernel like the normal one do not alter the conclusions below.

In the first part of the simulation experiment we compared the ML estimator with the semiparametric ones (with the MLE as initial starting value). Asymptotically both semiparametric estimators should behave as well as the MLE but one

may expect that the small sample properties of the semiparametric estimators are worse due to the inherent problems of choosing the bandwidth. These results are not reported here but they are comparable to those given in Table 1, from which MLE can be compared to the semiparametric procedure with the less efficient QMLE starting value.

Of course, ML estimation is not feasible in practice since the underlying distribution is not known. Therefore, we used the QMLE as starting point. Since  $\mu$  vanishes for the situation chosen here and  $\alpha\sigma^2 + \beta < 1$ , Theorems 2 and 3 of Lee and Hansen (1994) are applicable and the QMLE is  $\sqrt{n}$ -consistent. This estimator has been improved by the one-step Newton method. For convenience we also report the behavior of the unfeasible MLE in Table 1. The mean values of the estimates in 2500 replications are given together with their sample standard deviations.

To calculate the efficiency of the QMLE, observe that the asymptotic variance of the QMLE is equal to the well-known variance formula  $A^{-1}BA^{-1}$ , where  $A$  is the expectation under  $(\alpha, \beta, \sigma, g)$  of the second derivative of the pseudo log-likelihood (with a wrongly specified normal density) and  $B$  the expectation of the squared first derivative. With  $W_s$  as defined just below (2.20), straightforward calculations show

$$A = 2EW_s W_s',$$

$$B = (\kappa - 1)EW_s W_s',$$

where  $\kappa = \int \varepsilon^4 g(\varepsilon) d\varepsilon$ . Except for the normal distribution, the matrices  $A^{-1}$  and  $B^{-1}$  are generally not equal. Since the asymptotic variance of the QMLE is equal to the lower bound up to a constant, the asymptotic efficiency of each component of the QMLE is given by

$$\frac{4}{(\kappa - 1)I_s(g)} = \frac{4}{\int (\varepsilon^2 - 1)^2 g(\varepsilon) d\varepsilon \int (1 + \varepsilon \ell'(\varepsilon))^2 g(\varepsilon) d\varepsilon} \leq 1.$$

The latter inequality follows from Cauchy–Schwarz applied to the following identity:

$$-2 = E(\varepsilon^2 - 1)(1 + \varepsilon \ell'(\varepsilon)) = \int (\varepsilon^2 - 1)(1 + \varepsilon \ell'(\varepsilon))g(\varepsilon) d\varepsilon.$$

Since the lower bound for  $\alpha$  and  $\beta$  does not change in the semiparametric setting, this expression also entails the loss in the semiparametric model and shows the (potential) gain of the semiparametric estimator (3.2).

Except for the mixture distribution we can exactly calculate the efficiency of QMLE with respect to MLE. For the standardized double exponential the relative efficiency is  $\frac{4}{5}$ , for standardized Student distributions with  $\nu$  degrees of freedom it is  $1 - 12/\nu(\nu - 1)$ , and for standardized chi-squared distributions with  $\nu$  degrees of freedom it is  $(\nu - 4)/(\nu + 6)$ . For these heavy-tailed distributions the efficiency

Table 1

Comparison of MLE, QMLE, and two semiparametric one-step estimators in the GARCH(1,1) model with eight different standardized innovation distributions. Number of observations  $n = 1000$ , true parameters  $(\alpha, \beta) = (0.3, 0.6)$ ,  $(0.1, 0.8)$ , and  $(0.05, 0.9)$ , respectively. The sample means of 2500 independent replications and their sample standard deviations are given

		$\alpha = 0.300$				$\beta = 0.600$				$\alpha = 0.100$				$\beta = 0.800$				$\alpha = 0.050$				$\beta = 0.900$			
		$\hat{\alpha}$	$\hat{\beta}$	$\hat{\sigma}_\alpha$	$\hat{\sigma}_\beta$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\sigma}_\alpha$	$\hat{\sigma}_\beta$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\sigma}_\alpha$	$\hat{\sigma}_\beta$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\sigma}_\alpha$	$\hat{\sigma}_\beta$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\sigma}_\alpha$	$\hat{\sigma}_\beta$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\sigma}_\alpha$	$\hat{\sigma}_\beta$
N	ML=QML	0.298	0.593	0.071	0.056	0.099	0.786	0.035	0.073	0.047	0.891	0.022	0.051	0.298	0.593	0.072	0.057	0.098	0.786	0.036	0.074	0.047	0.892	0.022	0.050
	1-step	0.298	0.593	0.072	0.056	0.099	0.786	0.036	0.073	0.047	0.891	0.022	0.051	0.298	0.593	0.072	0.056	0.099	0.786	0.036	0.073	0.047	0.891	0.022	0.051
	1-step(sym)	0.298	0.593	0.072	0.056	0.099	0.786	0.036	0.073	0.047	0.891	0.022	0.051	0.298	0.593	0.072	0.056	0.099	0.786	0.036	0.073	0.047	0.891	0.022	0.051
DE	ML	0.299	0.592	0.080	0.070	0.099	0.782	0.038	0.083	0.048	0.885	0.023	0.061	0.303	0.588	0.089	0.079	0.100	0.776	0.043	0.094	0.048	0.880	0.026	0.073
	QML	0.299	0.592	0.080	0.070	0.099	0.782	0.038	0.083	0.048	0.885	0.023	0.061	0.294	0.593	0.085	0.074	0.097	0.784	0.040	0.087	0.046	0.886	0.024	0.067
	1-step	0.299	0.592	0.080	0.070	0.099	0.782	0.038	0.083	0.048	0.885	0.023	0.061	0.295	0.592	0.083	0.073	0.097	0.783	0.039	0.086	0.046	0.885	0.024	0.065
	1-step(sym)	0.299	0.592	0.080	0.070	0.099	0.782	0.038	0.083	0.048	0.885	0.023	0.061	0.295	0.592	0.083	0.073	0.097	0.783	0.039	0.086	0.046	0.885	0.024	0.065
NM	ML	0.295	0.595	0.058	0.041	0.098	0.790	0.029	0.054	0.047	0.898	0.018	0.030	0.295	0.595	0.059	0.042	0.097	0.790	0.030	0.054	0.046	0.897	0.018	0.032
	QML	0.295	0.595	0.059	0.042	0.097	0.790	0.030	0.054	0.046	0.897	0.018	0.032	0.295	0.595	0.060	0.043	0.098	0.793	0.030	0.056	0.047	0.901	0.018	0.032
	1-step	0.295	0.595	0.059	0.042	0.097	0.790	0.030	0.054	0.046	0.897	0.018	0.032	0.295	0.595	0.060	0.043	0.098	0.793	0.030	0.056	0.047	0.901	0.018	0.032
	1-step(sym)	0.295	0.595	0.059	0.042	0.099	0.793	0.030	0.056	0.047	0.901	0.018	0.032	0.295	0.595	0.059	0.042	0.099	0.793	0.030	0.056	0.047	0.901	0.018	0.032
$t_5$	ML	0.295	0.592	0.076	0.067	0.100	0.787	0.036	0.071	0.048	0.888	0.021	0.054	0.296	0.586	0.098	0.086	0.101	0.777	0.047	0.101	0.048	0.879	0.027	0.083
	QML	0.296	0.586	0.098	0.086	0.101	0.777	0.047	0.101	0.048	0.879	0.027	0.083	0.284	0.594	0.080	0.071	0.094	0.791	0.037	0.081	0.044	0.890	0.022	0.064
	1-step	0.284	0.594	0.080	0.071	0.094	0.791	0.037	0.081	0.044	0.890	0.022	0.064	0.285	0.594	0.079	0.070	0.095	0.791	0.037	0.081	0.045	0.889	0.022	0.063
	1-step(sym)	0.285	0.594	0.079	0.070	0.095	0.791	0.037	0.081	0.045	0.889	0.022	0.063	0.295	0.592	0.076	0.067	0.100	0.787	0.036	0.071	0.048	0.888	0.021	0.054
$t_7$	ML	0.296	0.595	0.075	0.060	0.100	0.782	0.037	0.079	0.047	0.885	0.021	0.063	0.298	0.592	0.086	0.070	0.101	0.776	0.042	0.094	0.047	0.882	0.024	0.076
	QML	0.298	0.592	0.086	0.070	0.101	0.776	0.042	0.094	0.047	0.882	0.024	0.076	0.291	0.597	0.078	0.064	0.096	0.784	0.038	0.082	0.045	0.886	0.022	0.068
	1-step	0.291	0.597	0.078	0.064	0.096	0.784	0.038	0.082	0.045	0.886	0.022	0.068	0.292	0.597	0.077	0.063	0.097	0.783	0.038	0.082	0.045	0.886	0.022	0.068
	1-step(sym)	0.292	0.597	0.077	0.063	0.097	0.783	0.038	0.082	0.045	0.886	0.022	0.068	0.298	0.592	0.076	0.060	0.098	0.783	0.037	0.077	0.047	0.887	0.022	0.058
$t_9$	ML	0.298	0.592	0.076	0.060	0.098	0.783	0.037	0.077	0.047	0.887	0.022	0.058	0.300	0.591	0.083	0.066	0.099	0.781	0.040	0.085	0.048	0.886	0.024	0.064
	QML	0.300	0.591	0.083	0.066	0.099	0.781	0.040	0.085	0.048	0.886	0.024	0.064	0.295	0.593	0.079	0.062	0.096	0.785	0.038	0.080	0.046	0.889	0.022	0.057
	1-step	0.295	0.593	0.079	0.062	0.096	0.785	0.038	0.080	0.046	0.889	0.022	0.057	0.295	0.593	0.077	0.062	0.096	0.784	0.038	0.079	0.046	0.889	0.022	0.057
	1-step(sym)	0.295	0.593	0.077	0.062	0.096	0.784	0.038	0.079	0.046	0.889	0.022	0.057	0.297	0.596	0.042	0.034	0.099	0.796	0.020	0.036	0.050	0.899	0.012	0.022
$\chi^2_6$	ML	0.297	0.596	0.042	0.034	0.099	0.796	0.020	0.036	0.050	0.899	0.012	0.022	0.299	0.589	0.091	0.073	0.101	0.780	0.042	0.096	0.048	0.884	0.024	0.072
	QML	0.299	0.589	0.091	0.073	0.101	0.780	0.042	0.096	0.048	0.884	0.024	0.072	0.283	0.603	0.062	0.051	0.092	0.801	0.030	0.061	0.045	0.898	0.017	0.047
	1-step	0.283	0.603	0.062	0.051	0.092	0.801	0.030	0.061	0.045	0.898	0.017	0.047	0.298	0.596	0.057	0.045	0.099	0.794	0.029	0.048	0.048	0.893	0.016	0.036
$\chi^2_{12}$	ML	0.298	0.596	0.057	0.045	0.099	0.794	0.029	0.048	0.048	0.893	0.016	0.036	0.299	0.592	0.084	0.064	0.100	0.782	0.041	0.079	0.047	0.881	0.023	0.071
	QML	0.299	0.592	0.084	0.064	0.100	0.782	0.041	0.079	0.047	0.881	0.023	0.071	0.289	0.598	0.065	0.051	0.095	0.796	0.032	0.061	0.045	0.891	0.018	0.049
	1-step	0.289	0.598	0.065	0.051	0.095	0.796	0.032	0.061	0.045	0.891	0.018	0.049												

losses of QMLE with respect to MLE show up in Table 1 and we see that the semiparametric methods regain most of the loss caused by the inefficient QMLE method. For light-tailed alternatives, as in the mixture case, the situation is less clear cut. There the efficiency is approximately 0.94 and the performance of the estimators is not much different. For the normal distribution MLE and

QMLE are of course equivalent. The use of the additional symmetry information hardly improves the estimated standard deviation of the semiparametric estimator (maximal 0.002), just as expected from our general theory. In empirical data sets, one often observes outlier-type innovation distributions with high kurtoses. Therefore, it seems worthwhile to apply the semiparametric estimation programs in these situations.

We conclude this section with a simple empirical example based on daily data. We applied our estimation methods to 15 logarithmic differenced exchange rate series for the period 1 January, 1980 to 1 April, 1994 ( $n = 3719$ ): Austrian Schilling (AS), Australian Dollar (AD), Belgium Franc (BF), British Pound (BP), Canadian Dollar (CD), Dutch Guilder (DG), Danish Kroner (DK), French Franc (FF), German Mark (GM), Italian Lire (IL), Japanese Yen (JY), Norwegian Kroner (NK), Swiss Franc (SF), Swedish Kroner (SK), and Spanish Peseta (SP), all with respect to US Dollar. These data are taken from Datastream. To facilitate the interpretation of the autoregression parameters we have standardized the series such that the QMLE of  $\sigma$  equals 1. In all series both the QMLE method and the semiparametric procedure estimate the persistence  $\alpha\sigma^2 + \beta$  less than one (for the semiparametric estimates this cannot be inferred from Table 2 since the semiparametric estimate of  $\sigma$  is not constrained to equal 1). The estimates based on the original data sets are given in the first four columns of Table 2. Of course, we used the variance formula  $A^{-1}BA^{-1}$  for the direct estimate of the standard deviation of the QMLE. As described above, the parameter estimates produced by the semiparametric procedure are not very sensitive to the choice of the bandwidth. However, it turns out that the direct variance estimates change dramatically (even for small changes of the bandwidth). Therefore, these estimates are not reliable and they have been deleted from the table.

For the simulation study above the situation was quite different since we estimated the variance of the semiparametric one-step estimators from independent parameter estimates in the replications. Here we have only one data set. Independent replications are not available. This leads to the following paradox. On the one hand, one may have the imprecise QML estimate with quite large estimated standard deviations. So it may be possible that the hypotheses of integrated GARCH or no conditional heteroskedasticity cannot be rejected. On the other, one has the improved semiparametric estimate which allows for more powerful tests. But since the estimated standard deviations are unreliable one can get any answer one wants by changing the bandwidth. To avoid this paradox, we propose to use the bootstrap. I.e. simulate replications of the original data set with the estimated parameter and the estimated innovation distribution as inputs and proceed as in the case of simulations described above. Then we have several parameter estimates available from which we calculate the straightforward sample estimate of the variance. In this manner we only rely upon the parameter estimates and not on direct estimates of the variance. Hence, the variability of the variance due to different bandwidth choices is greatly

Table 2

Comparison of QMLE and a semiparametric one-step estimator for several logarithmic differenced daily exchange rate series. Observation period 1 January 1980 to 1 April 1994 ( $n = 3719$ ). The first part of the table gives the estimates based on the original data set. Estimated standard deviations are deleted for the semiparametric estimators. The sample means and sample standard deviations of 500 bootstrap replications are given in the second half of the table

		Estimates based on				Bootstrap samples			
		Original data							
		$\hat{\alpha}$	$\hat{\beta}$	$\hat{\sigma}_x$	$\hat{\sigma}_\beta$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\sigma}_x$	$\hat{\sigma}_\beta$
AD	QMLE	0.129	0.843	0.075	0.034	0.116	0.828	0.044	0.063
	1-step	0.253	0.867			0.112	0.844	0.027	0.024
AS	QMLE	0.075	0.891	0.013	0.016	0.073	0.888	0.018	0.018
	1-step	0.113	0.897			0.072	0.890	0.014	0.013
BF	QMLE	0.068	0.902	0.011	0.023	0.068	0.899	0.019	0.018
	1-step	0.093	0.906			0.065	0.903	0.013	0.014
BP	QMLE	0.052	0.932	0.008	0.013	0.051	0.931	0.015	0.012
	1-step	0.055	0.932			0.050	0.931	0.012	0.010
CD	QMLE	0.138	0.798	0.042	0.062	0.139	0.793	0.032	0.031
	1-step	0.169	0.809			0.133	0.797	0.021	0.021
DG	QMLE	0.078	0.888	0.013	0.016	0.077	0.886	0.018	0.018
	1-step	0.107	0.916			0.076	0.887	0.015	0.014
DK	QMLE	0.067	0.902	0.011	0.016	0.065	0.898	0.015	0.016
	1-step	0.095	0.920			0.064	0.901	0.012	0.014
FF	QMLE	0.088	0.873	0.016	0.017	0.088	0.869	0.023	0.022
	1-step	0.119	0.913			0.085	0.872	0.017	0.017
GM	QMLE	0.073	0.894	0.012	0.013	0.073	0.891	0.017	0.018
	1-step	0.095	0.925			0.072	0.893	0.014	0.015
IL	QMLE	0.093	0.869	0.016	0.031	0.092	0.864	0.022	0.020
	1-step	0.109	0.896			0.090	0.867	0.019	0.017
JY	QMLE	0.059	0.891	0.017	0.025	0.060	0.888	0.016	0.026
	1-step	0.078	0.912			0.057	0.891	0.012	0.021
NK	QMLE	0.080	0.907	0.009	0.014	0.078	0.906	0.023	0.015
	1-step	0.092	0.916			0.075	0.908	0.017	0.010
SF	QMLE	0.059	0.904	0.012	0.013	0.058	0.901	0.014	0.019
	1-step	0.064	0.922			0.057	0.903	0.012	0.015
SK	QMLE	0.221	0.754	0.035	0.119	0.210	0.751	0.049	0.038
	1-step	0.185	0.839			0.209	0.756	0.032	0.020
SP	QMLE	0.106	0.871	0.014	0.032	0.104	0.868	0.027	0.020
	1-step	0.171	0.912			0.100	0.870	0.021	0.014

reduced. Some simulation experiments show that this procedure works quite well. We apply the bootstrap procedure to our data sets and we report the sample means and sample standard deviations in the final four columns of Table 2. Observe that the estimated standard deviations of the semiparametric estimators of the heteroskedastic parameters are between four tenth (AD) and nine tenth (IL) of the estimated ones for the QMLE method. This implies the efficiency of the QMLE method lies approximately in the interval (0.15, 0.80) in

these special examples. The efficiency gain is also supported by the plots in Fig. 1 of the nonparametric density estimates and the corresponding score estimates which are far away from the normal density and score. Although these figures suggest some skewness of the exchange rate densities, they are close to the densities of Student's  $t_\nu$ -distributions with  $\nu$  between 4.1 and 5.4. If the true underlying density would be symmetric, we expect from the simulation study that the symmetric nonparametric procedure performs slightly better in moderate samples. However, in the exchange rate applications the latter procedure yields somewhat larger standard deviations (0.003 for AS and less than

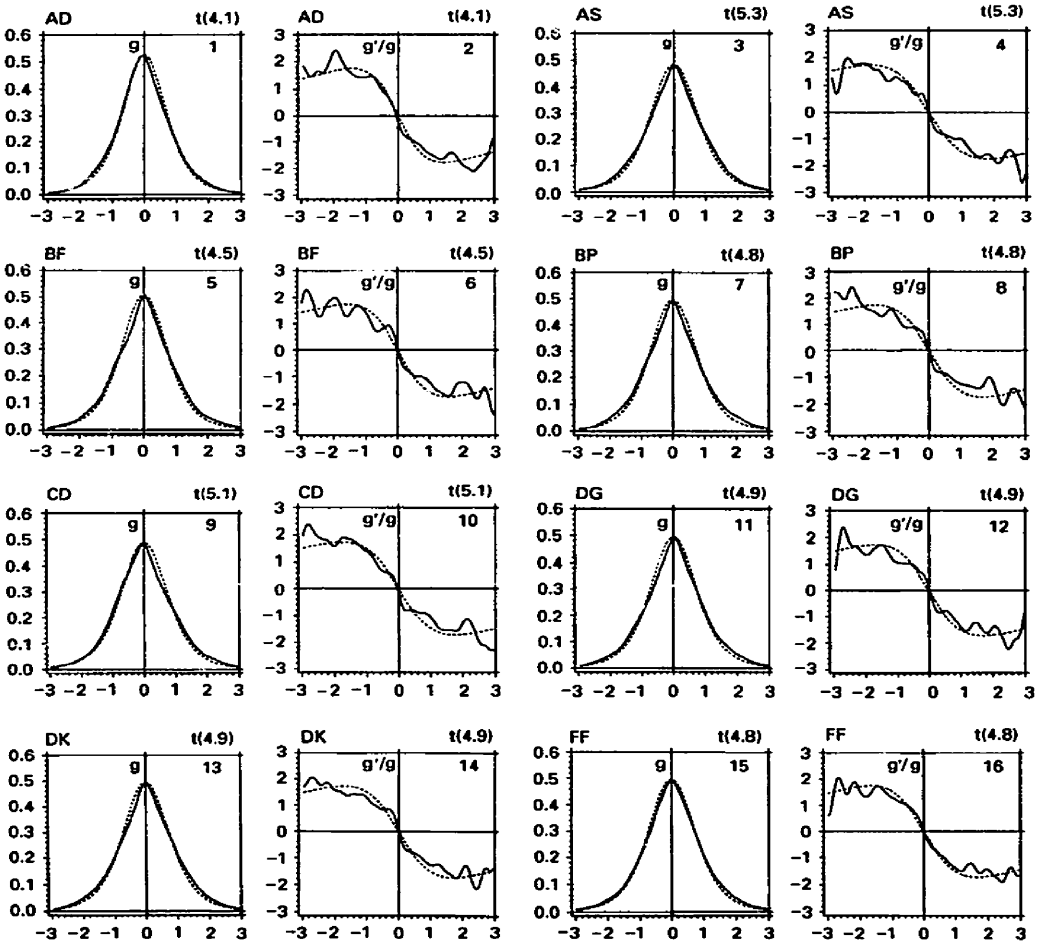


Fig. 1. Comparison of estimated densities and scores with  $t_\nu$ -densities and scores for several logarithmic differenced daily exchange rate series. Observation period 800101-940401 ( $n = 3719$ ).

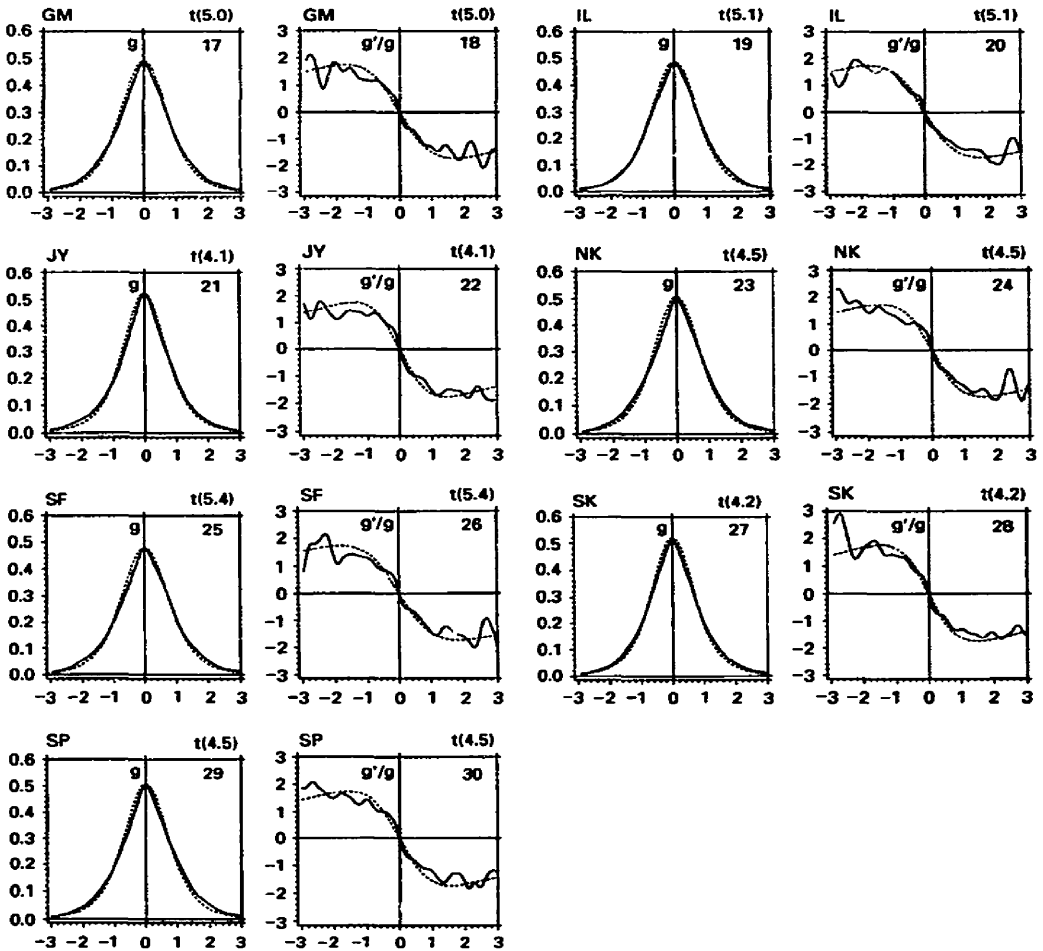


Fig. 1. Continued.

0.002 for the others, these values are not reported here). This indicates that the true densities are not fully symmetric and hence the symmetric semiparametric approach may lead to wrong conclusions. Since the possible moderate sample loss is very small it seems to be safer to use the ordinary non-symmetric improvement.

Finally, we note that the simulation results of Table 1 show that all estimators, even the unfeasible MLE, tend to underestimate the heteroskedasticity parameters. This negative bias explains why in Table 2, on average, the bootstrap estimates are less in value than the original estimates.

## 5. Conclusions

In this paper we studied the semiparametric properties of (integrated) GARCH-M-type models. In this model, adaptive estimation is not possible. This fact is completely caused by a location-scale parameter. After a suitable reparametrization of the model we showed that the estimation problem of the parameters characterizing the conditional heteroskedastic character of the process is equally difficult in cases where the innovation distribution is known or unknown, respectively. In that sense we may call these parameters still adaptively estimable. This property is derived in a general GARCH context avoiding moment conditions and including integrated GARCH models. The simulations showed that this property is not only interesting from a theoretical point of view. In moderate sample sizes with  $n = 1000$  observations, usually available in financial time-series, the semiparametric procedures work reasonably well. Most of the loss caused by the QMLE method (instead of the infeasible MLE method) is regained by the one-step estimator in case of the interesting group of heavy-tailed alternatives. Moreover, the empirical example showed that the efficiency loss caused by the QMLE method may be considerable.

It is clear from the exposition in this paper that the adaptivity results carry over to complicated models with time dependent mean and variance structures, e.g., ARMA with GARCH errors. The basic conditions given in DKW(1997) do not seem to put serious restrictions on the models. However, a complete verification of the technical details may be much more demanding.

## Appendix

*Proof of the LAN Theorem 2.1.* Since the general GARCH model (2.1), (2.2) is a location-scale model in which the location-scale parameter only depends on the past, our model fits into the general time-series framework of DKW(1997), especially Section 4. Therefore, it suffices to verify the conditions (2.3'), (A.1), and (2.4) of DKW(1997). In passing we also prove (3.3') of DKW(1997) which we will need in the proof of Theorem 3.1. I.e., with the notation introduced in (2.10), (2.14), and (2.15), and  $I_{ls}(g)$  the expectation under  $\theta$  of the product  $\psi(\theta)\psi(\theta)'$ , we have to show, under  $\theta_0$ ,

$$n^{-1} \sum_{i=1}^n W_i(\theta_0) I_{ls}(g) W_i(\theta_0)' \xrightarrow{P} I(\theta_0) > 0,$$

$$n^{-1} \sum_{i=1}^n |W_i(\theta_0)|^2 1_{\{n^{-1/2}|W_i(\theta_0)| > \delta\}} \xrightarrow{P} 0, \quad (\text{A.1})$$

$$n^{-1} \sum_{i=1}^n W_i(\theta_0) \xrightarrow{P} W(\theta_0), \quad (\text{A.2})$$

$$n^{-1} \sum_{t=1}^n |W_t(\theta_n) - W_t(\theta_0)|^2 \xrightarrow{P} 0, \tag{A.3}$$

and, under  $\theta_n$ ,

$$\sum_{t=1}^n |n^{-1/2}(M_{nt}, S_{nt})' - W_t(\theta_n)'(\tilde{\theta}_n - \theta_n)|^2 \xrightarrow{P} 0, \tag{A.4}$$

for some positive-definite matrix  $I(\theta_0)$  and some random matrix  $W(\theta_0)$ . Together with their Lemma A.1, these four relations yield the desired conclusions. We prepare the proof by deriving some helpful results.

Although  $W_t(\theta_0)$  is not stationary under  $\theta_0$ , the following proposition shows that these variables can be approximated by a stationary sequence.

**Proposition A.1.** Let  $h_t(\theta)$ ,  $H_t(\theta)$ , and  $W_t(\theta)$  be given by (2.9), (2.13), and (2.14), respectively, and let  $h_{st}(\theta)$ ,  $H_{st}(\theta)$ , and  $W_{st}(\theta)$  be their corresponding stationary solutions under  $\theta$ , i.e.

$$h_{st}(\theta) = \sum_{j=0}^{\infty} \prod_{k=1}^j \{\beta + \alpha \xi_{t-k}^2\}, \quad H_{st}(\theta) = \sum_{i=0}^{\infty} \beta^i h_{s,t-1-i}(\theta) \begin{pmatrix} \xi_{t-1-i}^2 \\ 1 \end{pmatrix},$$

$$W_{st}(\theta) = \sigma^{-1} \begin{pmatrix} \frac{1}{2} h_{st}^{-1}(\theta) H_{st}(\theta) (\mu, \sigma) \\ I_2 \end{pmatrix}.$$

Then, under  $\theta_0$ ,

$$n^{-1} \sum_{t=1}^n |W_t(\theta_0) - W_{st}(\theta_0)|^2 \rightarrow 0 \text{ (a.s.) as } n \rightarrow \infty. \tag{A.5}$$

*Proof.* We adopt the convention that empty sums are equal to zero while empty products are equal to one. Iterating  $h_t(\theta)$  yields

$$\begin{aligned} h_t(\theta) &= 1 + \beta h_{t-1}(\theta) + \alpha y_{t-1}^2 = 1 + h_{t-1}(\theta) \{\beta + \alpha \xi_{t-1}^2(\theta)\} \\ &= \sum_{j=0}^{t-1} \prod_{k=1}^j \{\beta + \alpha \xi_{t-k}^2(\theta)\} + h_{t-i}(\theta) \prod_{k=1}^i \{\beta + \alpha \xi_{t-k}^2(\theta)\}, \quad 0 \leq i \leq t-1, \end{aligned} \tag{A.6}$$

and hence

$$h_{t-i}(\theta)/h_t(\theta) \leq \prod_{k=1}^i \{\beta + \alpha \xi_{t-k}^2(\theta)\}^{-1}, \quad 0 \leq i \leq t-1. \tag{A.7}$$

Under  $\theta$ , the calculated variables  $\bar{\xi}_t(\theta)$  simply are the true innovations  $\xi_t$  in (A.6) and (A.7). For the stationary random variables  $h_{st}(\theta)$  we obtain similar relations,

$$h_{st}(\theta) = \sum_{j=0}^{i-1} \prod_{k=1}^j \{\beta + \alpha \xi_{t-k}^2\} + h_{s,t-i}(\theta) \prod_{k=1}^i \{\beta + \alpha \xi_{t-k}^2\}, \quad 0 \leq i,$$

$$h_{s,t-i}(\theta)/h_{st}(\theta) \leq \prod_{k=1}^i \{\beta + \alpha \xi_{t-k}^2\}^{-1}, \quad 0 \leq i,$$

and hence, under  $\theta$ , we obtain

$$\begin{aligned} |h_{st}(\theta)h_{t-i}(\theta) - h_t(\theta)h_{s,t-i}(\theta)| &= |h_{s,t-i}(\theta) - h_{t-i}(\theta)| \sum_{j=0}^{i-1} \prod_{k=1}^j \{\beta + \alpha \xi_{t-k}^2\} \\ &\leq h_t(\theta) |h_{s1}(\theta) - h_1(\theta)| \prod_{k=1}^{t-1-i} \{\beta + \alpha \xi_{t-i-k}^2\} \\ &= h_t(\theta) |h_{s1}(\theta) - h_1(\theta)| \prod_{k=1}^i \{\beta + \alpha \xi_{t-k}^2\}^{-1} \prod_{k=1}^{t-1} \{\beta + \alpha \xi_k^2\}, \quad 0 \leq i \leq t-1. \end{aligned}$$

With  $C$  some generic constant only depending on  $\theta$  we obtain, under  $\theta$ ,

$$\begin{aligned} |W_t(\theta) - W_{st}(\theta)| &\leq C |H_t(\theta)/h_t(\theta) - H_{st}(\theta)/h_{st}(\theta)| \\ &\leq C \sum_{i=0}^{t-2} \beta^i \left| \frac{h_{t-1-i}(\theta)}{h_t(\theta)} - \frac{h_{s,t-1-i}(\theta)}{h_{st}(\theta)} \right| \left| \begin{pmatrix} \xi_{t-1-i}^2 \\ 1 \end{pmatrix} \right| \\ &\quad + C \sum_{i=t-1}^{\infty} \beta^i \frac{h_{s,t-1-i}(\theta)}{h_{st}(\theta)} \left| \begin{pmatrix} \xi_{t-1-i}^2 \\ 1 \end{pmatrix} \right| \\ &\leq C |h_{s1}(\theta) - h_1(\theta)| \prod_{k=1}^{t-1} \{\beta + \alpha \xi_k^2\} \sum_{i=0}^{t-2} \prod_{k=1}^i \frac{\beta}{\beta + \alpha \xi_{t-k}^2} \\ &\quad + C \sum_{i=t-1}^{\infty} \prod_{k=1}^i \frac{\beta}{\beta + \alpha \xi_{t-k}^2} \\ &\leq C |h_{s1}(\theta) - h_1(\theta)| (t-1) \prod_{k=1}^{t-1} \{\beta + \alpha \xi_k^2\} \\ &\quad + C \prod_{k=1}^{t-1} \frac{\beta}{\beta + \alpha \xi_k^2} \sum_{i=-1}^{\infty} \prod_{k=0}^i \frac{\beta}{\beta + \alpha \xi_{t-k}^2}. \end{aligned}$$

By (2.3) the right-hand side tends to zero (a.s.), as  $t \rightarrow \infty$ . Cesàro’s theorem completes the proof of the proposition.  $\square$

Intuitively it is clear that slight perturbations of the parameters yield solutions of Eqs. (2.1) and (2.2) that are close. The following proposition makes this more precise. Just as expected from heuristic formal calculations, the leading term of

$\tilde{h}_{nt}/h_{nt} - 1$  is a linear combination of the components of  $H_{nt}/h_{nt}$  which appears in the score  $\dot{l}_{nt}$ .

**Proposition A.2.** Let  $h_t(\theta)$  and  $H_t(\theta)$  be given by (2.9) and (2.13), respectively, and define

$$\begin{aligned} Q_t(\theta) &= H_t(\theta)/h_t(\theta) = \sum_{i=0}^{t-2} \beta^i \left( \frac{y_{t-1-i}^2}{h_{t-1-i}(\theta)} \right) / h_t(\theta) \\ &= \sum_{i=0}^{t-2} \beta^i \frac{h_{t-1-i}(\theta)}{h_t(\theta)} \left( \frac{\xi_{t-1-i}^2}{1} \right), \end{aligned}$$

$$R_t(\theta, \tilde{\theta}) = h_t(\tilde{\theta})/h_t(\theta) - 1 - (\tilde{\alpha} - \alpha, \tilde{\beta} - \beta)Q_t(\theta).$$

Let  $\theta_n$  and  $\tilde{\theta}_n$  satisfy the conditions just above (2.6). Put  $Q_{nt} = Q_t(\theta_n)$  and  $R_{nt} = R_t(\theta_n, \tilde{\theta}_n)$ . Then, under  $\theta_n$ ,

$$\begin{aligned} n^{-1} \sum_{t=1}^n |Q_{nt}|^2 &= O_P(1), \quad n^{-1} \sum_{t=1}^n |Q_{nt}|^2 1_{\{n^{-1/2}|Q_{nt}| > \delta\}} \rightarrow 0, \\ &\text{(a.s.) as } n \rightarrow \infty, \end{aligned} \tag{A.8}$$

$$\sum_{t=1}^n R_{nt}^2 \rightarrow 0, \text{ (a.s.) as } n \rightarrow \infty. \tag{A.9}$$

*Proof.* By Eq. (A.7) we obtain

$$Q_t(\theta) \leq \beta^{-1} \sum_{i=0}^{t-2} \prod_{k=1}^{i+1} \frac{\beta}{\beta + \alpha \xi_{t-k}^2} \left( \frac{\xi_{t-1-i}^2}{1} \right).$$

For  $n$  sufficiently large, this latter relation shows that, under  $\theta_n$ ,  $|Q_{nt}|$  may be bounded by the product of a constant depending on  $\theta_0$  only and the stationary sequence

$$S_t = \sum_{i=0}^{\infty} \prod_{k=1}^i \frac{\beta_0}{\beta_0 + \frac{1}{2}\alpha_0 \xi_{t-k}^2}.$$

Note that all moments of  $S_t$  exist. The relations concerning  $Q_{nt}$  are easily obtained.

To prove the result concerning the remainder term  $R_{nt}$  note that an explicit relationship for the difference of  $h_t(\tilde{\theta})$  and  $h_t(\theta)$  is given by [compare (2.3) of Kreiss, 1987a]

$$h_t(\tilde{\theta}) - h_t(\theta) = \sum_{i=0}^{t-2} \tilde{\beta}^i h_{t-1-i}(\theta) (\tilde{\alpha} - \alpha, \tilde{\beta} - \beta) \left( \frac{\xi_{t-1-i}^2}{1} \right).$$

Hence, the remainder term  $R_t(\theta, \tilde{\theta})$  is given by

$$R_t(\theta, \tilde{\theta}) = \sum_{i=0}^{t-2} (\tilde{\beta}^i - \beta^i) \frac{h_{t-1-i}(\theta)}{h_t(\theta)} (\tilde{\alpha} - \alpha, \tilde{\beta} - \beta) \binom{\xi_{t-1-i}^2}{1}(\theta).$$

Choose  $c > 1$  such that  $Ec\beta_0/(\beta_0 + \frac{1}{2}\alpha_0\xi_1^2) < 1$ . By the mean value theorem, there exists a  $\tilde{\beta}_{ni}$  in between  $\tilde{\beta}_n$  and  $\beta_n$  such that, for  $n$  sufficiently large,

$$|\tilde{\beta}_n^i - \beta_n^i| = |\tilde{\beta}_n - \beta_n| i \tilde{\beta}_{ni}^{i-1} \leq |\tilde{\beta}_n - \beta_n| ic^{i-1} \beta_n^{i-1}, \quad i \geq 0.$$

Just as for  $Q_{nt}$ , we may bound  $R_{nt}$  by the product of a constant times  $n^{-1}$  and the stationary sequence

$$S_t^* = \sum_{i=0}^{\infty} i \prod_{k=1}^i \frac{c\beta_0}{\beta_0 + \frac{1}{2}\alpha_0\xi_{t-k}^2}.$$

The proof of the proposition can be easily completed.  $\square$

Now we are ready to prove (A.1)–(A.4). Define  $I(\theta_0) = E_{\theta_0} W_{s1}(\theta_0) I_{ts}(g) W_{s1}(\theta_0)'$  and  $W(\theta_0) = E_{\theta_0} W_{s1}(\theta_0)$  (the existence of these quantities can be obtained along the lines of the proof of Proposition A.2 since  $|W_{st}(\theta_0)|$  is bounded by the product of  $S_t$  and a constant depending on  $\theta_0$ , only). Obviously, the relations (A.1) and (A.2) hold true if  $W_t(\theta_0)$  is replaced by the stationary ergodic sequence  $W_{st}(\theta_0)$ . Consequently, Proposition A.1 implies the validity of these relations for  $W_t(\theta_0)$  itself.

To prove (A.4) we will use Proposition A.2. Writing  $\lambda_n = (\lambda'_{1n}, \lambda'_{2n})'$  with  $\lambda_{1n}$  ( $\lambda_{2n}$ ) the first (latter) two components of  $\lambda_n$ , and defining

$$\chi(x) = \{-1 + 2(\sqrt{1+x} - 1)/x\} 1_{\{x \geq -1\}},$$

we see

$$\begin{aligned} & \sum_{t=1}^n |n^{-1/2}(M_{nt}, S_{nt})' - W_t(\theta_n)'(\tilde{\theta}_n - \theta_n)|^2 \\ &= \sigma_n^{-2} n^{-1} \sum_{t=1}^n \left| (\tilde{\mu}_n, \tilde{\sigma}_n)' \right. \\ & \times \left\{ \frac{1}{2}(\lambda'_{1n} Q_{nt} + \sqrt{n} R_{nt}) \chi(n^{-1/2} \lambda'_{1n} Q_{nt} + R_{nt}) + \frac{1}{2} \sqrt{n} R_{nt} \right\} \\ & \left. + n^{-1/2} \lambda_{2n} \frac{1}{2} \lambda'_{1n} Q_{nt} \right|^2. \end{aligned}$$

Together with Proposition A.2, Lemma 2.1 of DKW(1997) (with  $Y_{nt} = \lambda'_{1n} Q_{nt}$ ,  $X_{nt} = \lambda'_{1n} Q_{nt} + \sqrt{n} R_{nt}$ , and the function  $\phi = \chi^2$  as above) yields (A.4).

Finally, we have to prove (A.3). Note that

$$|W_t(\theta_n) - W_t(\theta_0)|^2 \leq C|Q_t(\theta_n) - Q_t(\theta_0)|^2 + C|Q_t(\theta_0)|^2|\theta_n - \theta_0|^2$$

and obtain contiguity of  $P_{\theta_n}$  and  $P_{\theta_0}$  from (A.1) and (A.4), and Theorem 2.1 of DKW(1997). Then the required result is easily obtained from

$$\begin{aligned} Q_t(\tilde{\theta}) - Q_t(\theta) &= \sum_{i=0}^{t-2} (\tilde{\beta}^i - \beta^i) \frac{h_{t-1-i}(\tilde{\theta})}{h_t(\tilde{\theta})} \begin{pmatrix} \xi_{t-1-i}^2(\tilde{\theta}) \\ 1 \end{pmatrix} \\ &+ Q_t(\theta) \{(\theta_1 - \tilde{\theta}_1)' Q_t(\tilde{\theta}) + R_t(\tilde{\theta}, \theta)\} \\ &- \begin{pmatrix} 0 \\ \sum_{i=0}^{t-2} \beta^i \frac{h_{t-1-i}(\tilde{\theta})}{h_t(\tilde{\theta})} \{(\theta_1 - \tilde{\theta}_1)' Q_{t-1-i}(\tilde{\theta}) + R_{t-1-i}(\tilde{\theta}, \theta)\} \end{pmatrix} \end{aligned}$$

along the lines of the proofs of the propositions above. This completes the proofs of the theorems in Section 2.  $\square$

*Proof of Theorem 3.1.* It suffices to verify the conditions of DKW(1997). These reduce to (A.1)–(A.4) above, which are verified there, and the existence of an estimator  $\hat{\psi}_n(\cdot)$ , based on  $\varepsilon_1, \dots, \varepsilon_n$ , of  $\psi(\cdot) = -(l'(\cdot), 1 + \cdot l'(\cdot))'$ , from (2.15), satisfying the consistency condition

$$\int |\hat{\psi}_n(x) - \psi(x)|^2 g(x) dx \xrightarrow{P} 0, \text{ under } g.$$

Indeed, such an estimator exists in view of Proposition 7.8.1 of BKRW (1993) p. 400, with  $k = 0$  and  $k = 1$ ; see also Lemma 4.1 of Bickel (1982). The estimator  $\hat{\psi}_n(\cdot)$  in Section 3 is based on these constructions.  $\square$

### References

Baillie, R.T., Bollerslev, T., 1989. The message in daily exchange rates: a conditional-variance tale. *Journal of Business and Economic Statistics* 7, 297–305.

Bickel, P.J., 1982. On adaptive estimation. *Annals of Statistics* 10, 647–671.

Bickel, P.J., Klaassen, C.A.J., Ritov, Y., Wellner, J.A., 1993. *Efficient and Adaptive Estimation for Semiparametric Models*. Johns Hopkins University Press, Baltimore.

Bollerslev, T., 1986. Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics* 31, 307–327.

Bollerslev, T., Chou, R.Y., Kroner, K.F., 1992. ARCH modeling in finance: a review of the theory and empirical evidence. *Journal of Econometrics* 52, 5–59.

Diebold, F.X., 1988. *Empirical Modeling of Exchange Rate Dynamics*. Springer, New York.

- Drost, F.C., Klaassen, C.A.J., Werker, B.J.M., 1994. Adaptiveness in time-series models. In: Mandl, P., Hušková, M. (Eds.), *Asymptotic Statistics*. Physica Verlag, New York, pp. 203–212.
- Drost, F.C., Klaassen, C.A.J., Werker, B.J.M., 1997. Adaptive estimation in time-series models, *Annals of Statistics*, forthcoming.
- Drost, F.C., Nijman, T.E., 1993. Temporal aggregation of GARCH processes. *Econometrica* 61, 909–927.
- Drost, F.C., Werker, B.J.M., 1996. Closing the GARCH gap: continuous time GARCH modeling. *Journal of Econometrics* 74, 31–57.
- Engle, R.F., 1982. Autoregressive conditional heteroskedasticity with estimates of the variance of United Kingdom inflation. *Econometrica* 50, 987–1007.
- Engle, R.F., Bollerslev, T., 1986. Modelling the persistence of conditional variances. *Econometric Reviews* 5, 1–50, 81–87.
- Engle, R.F., González-Rivera, G., 1991. Semiparametric ARCH models. *Journal of Business and Economic Statistics* 9, 345–359.
- Engle, R.F., Lilien, D.M., Robins, R.P., 1987. Estimating time varying risk premia in the term structure: the ARCH-M model. *Econometrica* 55, 391–407.
- Gourieroux, C., Monfort, A., 1992. Qualitative threshold ARCH models. *Journal of Econometrics* 52, 159–199.
- Hájek, J., 1970. A characterization of limiting distributions of regular estimates. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete* 14, 323–330.
- Hájek, J., Šidák, Z., 1967. *Theory of Rank Tests*. Academia, Prague.
- Huber, P.J., 1981. *Robust Statistics*. Wiley, New York.
- Jeganathan, P., 1995. Some aspects of asymptotic theory with applications to time series models. *Econometric Theory* 11, 818–887.
- Klaassen, C.A.J., 1987. Consistent estimation of the influence function of locally asymptotically linear estimators. *Annals of Statistics* 15, 1548–1562.
- Koul, H.L., Schick, A., 1995. Efficient estimation in nonlinear time series models. Technical Report, Department of Statistics and Probability, Michigan State University.
- Kreiss, J.-P., 1987a. On adaptive estimation in stationary ARMA processes. *Annals of Statistics* 15, 112–133.
- Kreiss, J.-P., 1987b. On adaptive estimation in autoregressive models when there are nuisance functions. *Statistics and Decisions* 5, 59–76.
- Lee, S.-W., Hansen, B.E., 1994. Asymptotic theory for the GARCH(1,1) quasi-maximum likelihood estimator. *Econometric Theory* 10, 29–52.
- Linton, O., 1993. Adaptive estimation in ARCH models. *Econometric Theory* 9, 539–569.
- Lumsdaine, R.L., 1989. Asymptotic properties of the maximum likelihood estimator in GARCH(1,1) and IGARCH(1,1) models. Manuscript, Department of Economics, Harvard University.
- Nelson, D.B., 1990a. ARCH models as diffusion approximations. *Journal of Econometrics* 45, 7–38.
- Nelson, D.B., 1990b. Stationarity and persistence in the GARCH(1,1) model. *Econometric Theory* 6, 318–334.
- Nelson, D.B., 1991. Conditional heteroskedasticity in asset returns: a new approach. *Econometrica* 59, 347–370.
- Newey, W.K., 1990. Semiparametric efficiency bounds. *Journal of Applied Econometrics* 5, 99–135.
- Pötscher, B.M., 1995. Comment on ‘Adaptive estimation in time series regression models’ by D.G. Steigerwald. *Journal of Econometrics* 66, 123–129.
- Robinson, P.M., 1988. Semiparametric econometrics: a survey. *Journal of Applied Econometrics* 3, 35–51.
- Schick, A., 1986. On asymptotically efficient estimation in semiparametric models. *Annals of Statistics* 14, 1139–1151.
- Scholz, F.-W., 1971. Comparison of optimal location estimators. Ph.D. Thesis, University of California, Berkeley.

- Steigerwald, D.G., 1992. Adaptive estimation in time-series regression models. *Journal of Econometrics* 54, 251–275.
- Steigerwald, D.G., 1995. Reply to B.M. Pötscher's comment on 'Adaptive estimation in time series regression models'. *Journal of Econometrics* 66, 131–132.
- Tapia, R.A., Thompson, J.R., 1978. *Nonparametric probability density estimation*. Johns Hopkins University Press, Baltimore.
- Weiss, A.A., 1986. Asymptotic theory for ARCH models: estimation and testing. *Econometric Theory* 2, 107–131.