

## BOOTSTRAP MISSPECIFICATION TESTS FOR ARCH BASED ON THE EMPIRICAL PROCESS OF SQUARED RESIDUALS

LAJOS HORVÁTH<sup>a,\*</sup>, PIOTR KOKOSZKA<sup>b,†</sup> and GILLES TEYSSIERE<sup>c,‡</sup>

<sup>a</sup>Department of Mathematics, University of Utah, 155 South 1440 East, Salt Lake City, UT 84112-0090, USA; <sup>b</sup>Department of Mathematics and Statistics, Utah State University, 3900 Old Main Hill, Logan, UT 84322-3900, USA; <sup>c</sup>NBG Bank (France), GREQAM & CORE, 65, Avenue Roosevelt, F-75008 Paris, France

(Received 13 June 2002; In final form 20 February 2003)

We propose and study by means of simulations and graphical tools a class of goodness-of-fit tests for ARCH models. The tests are based on the empirical distribution function of squared residuals and smooth (parametric) bootstrap. We examine empirical size and power by means of a simulation study. While the tests have overall correct size, their power strongly depends on the type of alternative and is particularly high when the assumption of Gaussian innovations is violated. As an example, the tests are applied to returns on Foreign Exchange rates.

*Keywords:* ARCH model; Empirical process; Goodness-of-fit tests; Size–power curves; Smooth bootstrap; Squared residuals

### 1 INTRODUCTION

Suppose a sample  $y_1, \dots, y_n$  was observed and an ARCH type model has been postulated as an adequate description of the data. In this paper we are concerned with the problem of verifying if the postulated model fits the data. Tests of this kind are known as goodness-of-fit or misspecifications tests and play a central role in time series analysis.

In the classical “linear” time series analysis, such tests, known also as diagnostic checks, fall roughly into four categories: (1) examination of the residual plot, (2) portmanteau type tests based on weighted sums of covariances of residuals, (3) Lagrange-multiplier tests, (4) tests based on the empirical distribution function of the estimated residuals or on the spectral empirical distribution function (integrated periodogram).

Goodness-of-fit tests for non-linear time series models have only recently become an object of a more systematic research even though basic tools like residual plots have, of course, been used in a more or less intuitive way for a long time. Li and Mak (1994), Horváth and Kokoszka (2001) and Berkes *et al.* (2003, 2004) studied Portmanteau type

---

\* E-mail: horvath@math.utah.edu

† Corresponding author. E-mail: piotr@math.usu.edu

‡ E-mail: gilles@ehess.cnrs-mrs.fr/gteyssiere@nbg-france.com

statistics based on autocorrelations of residuals and squared residuals. Davis and Mikosch (1998) and Mikosch and Stáricá (2000) studied the autocorrelation function of the observations following, respectively, ARCH(1) and GARCH(1, 1) models. Lundbergh and Teräsvirta (2002) proposed a broad class of Lagrange-multiplier type specification tests whereas Mikosch and Stáricá (1999) proposed a goodness-of-fit test for GARCH based on the integrated periodogram of the observations.

Recently, Horváth *et al.* (2001) found the asymptotic distribution of the empirical process of ARCH( $p$ ) squared residuals and Berkes and Horváth (2003) extended these results to general GARCH( $p, q$ ) models. Relevant theoretical results are also presented in Chapter 8 of Koul (2002).

In this paper we follow up on this theoretical work by proposing and examining the finite sample performance of several goodness-of-fit tests based on the empirical distribution of squared residuals. These tests are particularly suited to detect departures from the postulated distribution of the unobservable innovations. To explain the idea, we focus on a simple ARCH( $p$ ) model. Suppose then that the postulated model is defined by

$$y_t = \sigma_t \varepsilon_t, \quad \sigma_t^2 = b_0 + \sum_{j=1}^p b_j y_{t-j}^2, \quad \varepsilon_t \sim N(0, 1) \quad (1.1)$$

and we want to check if this model is an acceptable approximation to the data. Having obtained estimates  $\hat{b}_0, \hat{b}_1, \dots, \hat{b}_p$ , we compute the squared residuals

$$\hat{\varepsilon}_k^2 = y_k^2 \left[ \hat{b}_0 + \sum_{j=1}^p \hat{b}_j y_{k-j}^2 \right]^{-1}, \quad k = p+1, \dots, n. \quad (1.2)$$

The empirical process of the squared residuals is then

$$\hat{\alpha}_n(x) = (n-p)^{-1/2} \sum_{k=p+1}^n (\mathbb{I}[\hat{\varepsilon}_k^2 \leq x] - (2\Phi(\sqrt{x}) - 1)), \quad (1.3)$$

where  $\Phi(\cdot)$  is the standard normal cumulative distribution function and  $\mathbb{I}[\cdot]$  is the indicator function. If a distribution different from normal is postulated for the  $\varepsilon_t$  in (1.1), then  $2\Phi(\sqrt{x}) - 1$  in (1.3) should be replaced by the cumulative distribution function of the  $\varepsilon_t^2$ .

It is intuitively clear that if (1.1) is an adequate description of the data, then the  $\hat{\varepsilon}_k^2$  should have distribution which is close to that of  $\varepsilon_k^2$ , so the empirical distribution function

$$\hat{F}(x) = \frac{1}{n-p} \sum_{k=p+1}^{\infty} \mathbb{I}[\hat{\varepsilon}_k^2 \leq x] \quad (1.4)$$

should be close to  $2\Phi(\sqrt{x}) - 1$  and consequently the empirical process (1.3) should be “small”. On the other hand, if (1.1) is a poor approximation to the data, then  $\hat{F}(x)$  will be very different from  $2\Phi(\sqrt{x}) - 1$  and so  $\hat{\alpha}_n$  will be “large”. This may be due not only to the fact that non-Gaussian innovations yield a better approximation, but also to the misspecified functional form of the model. For example, if ARCH( $p'$ ) with  $p' > p$  better fits to the data, then the estimates  $\hat{b}_j$  will have large bias and the  $\hat{\varepsilon}_k^2$  will have empirical distribution very different from  $2\Phi(\sqrt{x}) - 1$ , even when the  $\varepsilon_k$  can be assumed normal.

In Section 6 we present some mathematical and historical background on goodness-of-fit tests based on the empirical process of the residuals. Here we only mention that, unlike for many linear models, in the case of ARCH models the empirical process has an asymptotic distribution which depends in an intricate way on the distribution of the innovations and model parameters. For this reason asymptotically pivotal statistics based on the empirical process cannot be readily constructed. Nevertheless, by employing some form of bootstrap, standard goodness-of-fit tests might be appropriately modified. In this paper we propose a simple procedure of this kind and investigate its validity by means of a simulation study. We focus on three functionals of the empirical process  $\hat{\alpha}_n$  of squared residuals:

- The Cramér–von Mises (CVM) statistic

$$\hat{T}_n = \int |\hat{\alpha}_n(x)|^2 dx. \tag{1.5}$$

- The normalized Cramér–von Mises (NCVM) statistic

$$\widehat{NT}_n = \int |\hat{\alpha}_n(x)|^2 (2\pi)^{-1/2} \exp\left(\frac{-x^2}{2}\right) dx. \tag{1.6}$$

- The Kolmogorov–Smirnov (KS) statistic

$$\widehat{KS}_n = \max_{1 \leq i \leq n} |\hat{\alpha}_n(\hat{\varepsilon}_i^2)|. \tag{1.7}$$

In the case of model (1.1), our goodness-of-fit procedure works as follows. Denote by  $\hat{\tau}_n$  the statistic of interest, where  $\hat{\tau}_n = \hat{T}_n, \widehat{NT}_n, \widehat{KS}_n$ .

1. Having observed the sample  $y_1, \dots, y_n$ , obtain an estimate  $\hat{b}_n = (\hat{b}_0, \hat{b}_1, \dots, \hat{b}_p)$  using a pseudo maximum likelihood method (PML) described in Section 2.2.
2. Generate  $B$  independent sequences of i.i.d. standard normal random variables, each of length  $n + w$ . (The extra  $w$  values are used to avoid an initialization effect.) Denote each of the  $B$  sequences by  $\varepsilon_{-w+1}^*, \dots, \varepsilon_0^*, \varepsilon_1^*, \dots, \varepsilon_n^*$ .
3. Generate  $B$  bootstrap ARCH( $p$ ) realizations  $y_1^*, \dots, y_n^*$ , where the  $y_t^*$  satisfy

$$y_t^* = \sigma_t^* \varepsilon_t^*, \quad (\sigma_t^*)^2 = \hat{b}_0 + \sum_{j=1}^p \hat{b}_j (y_{t-j}^*)^2. \tag{1.8}$$

Thus  $\{y_t^*\}$  is a smooth bootstrap version of the sequence  $\{y_t\}$ . The choice of  $B$  is discussed in Section 2.1.

4. For each of the  $B$  samples  $y_1^*, \dots, y_n^*$  obtain estimates  $\hat{b}_0^*, \hat{b}_1^*, \dots, \hat{b}_p^*$  and construct bootstrap residuals

$$(\hat{\varepsilon}_k^2)^* = (y_k^*)^2 \left[ \hat{b}_0^* + \sum_{j=1}^p \hat{b}_j^* (y_{k-j}^*)^2 \right]^{-1}, \quad k = p + 1, \dots, n. \tag{1.9}$$

From the  $(\hat{\varepsilon}_k^2)^*$  construct the bootstrap empirical process  $\hat{\alpha}_n^*$  and compute the statistic  $\hat{\tau}_n^*$ . Then find, say, the 5% critical region which corresponds to the upper tail of the empirical distribution of the  $B$  numbers  $\hat{\tau}_n^*$ , or the  $P$ -value.

The paper is organized as follows. In Section 2 we describe our simulation study. Section 3 focuses on the presentation of the simulation results with graphical methods, whereas Section 4 contains an application of the tests proposed above to returns on Foreign exchange rates. We conclude in Section 5 and provide some mathematical and historical background in Section 6.

## 2 SIMULATION STUDY

In this section, we describe the details of simulation study.

### 2.1 Bootstrap Tests and the Choice of $B$

We consider  $N = 10,000$  replications of the DGP (data generating process), and for each replication  $j$ , we follow the test procedure described in Section 1. For a test of size  $\alpha$ , we calculate the bootstrap critical values  $\hat{C}_\alpha^*(j)$  as the  $(1 - \alpha)$ th percentiles of the  $\hat{\tau}_n^*$ . We consider  $\alpha = 1\%$ ,  $5\%$ , and  $10\%$ . For each replication  $j$ , we also calculate the bootstrap  $P$ -value  $\hat{P}_j^*$ ,  $j = 1, \dots, N$ , which is the empirical probability of observing a statistic  $\hat{\tau}^*$  greater than  $\hat{\tau}_n$ , *i.e.*

$$\hat{P}_j^* = \hat{P}_j^*(\hat{\tau}_n) = \frac{1}{B} \sum_{i=1}^B \mathbb{I}[\hat{\tau}_i^* > \hat{\tau}_n]. \quad (2.1)$$

The Gaussian errors have been generated using the Box-Muller method, which uses uniform deviates randomly drawn from two different random number generators, initialized with two different seeds. The resulting sequence of uniform deviates succeeds the whole set of Marsaglia's (1996) DIEHARD tests.

An important point is the choice of the number of draws  $B$  from the bootstrap DGP, as it involves a trade-off between power loss and computing time. Given that the computation of the CVM and NCVM bootstrap tests involve numerical integration, which is the most computing time consuming operation for these tests, we have to choose  $B$  in an optimal way to avoid excessive computing time. To explain how we intend to achieve this goal, we follow below the exposition of Davidson and MacKinnon (1999; 2000).

Let  $P^*$  be the ideal bootstrap  $P$ -value, *i.e.* the probability that  $\hat{\tau}_n^* > \hat{\tau}_n$  under the bootstrap DGP. The ideal bootstrap test rejects the null hypothesis at level  $\alpha$  if  $P^* < \alpha$ , while the feasible bootstrap test rejects the null hypothesis at level  $\alpha$  if  $\hat{P}^* < \alpha$ . As  $B \rightarrow \infty$ ,  $\hat{P}^* \rightarrow P^*$ , but since the number of draws  $B$  is finite, bootstrap tests involve some loss of power. The loss of power caused by a finite number of draws has been studied for Monte Carlo tests. These results are of interest in our case since bootstrap tests are Monte Carlo tests for pivotal statistics. For a Monte Carlo test to be exact,  $\alpha(B + 1)$  must be an integer, see Dufour and Kiviet (1998) and Dufour and Khalaf (2001). Even though the statistics considered in this paper are not pivotal, we follow the recommendation of Davidson and MacKinnon (2000) and always choose  $B$  so that  $\alpha(B + 1)$  is an integer.

In our simulations we used a pretesting procedure suggested by Davidson and MacKinnon (2000) which, they argue, is more efficient than the three-step method of Andrews and

Buchinsky (2000). This procedure consists of starting by choosing  $B$  relatively small, e.g.  $B = 99$ , and increasing it until we are confident at some level  $\beta$ , say  $\beta = 0.001$ , that the estimated bootstrap  $P$ -value  $\hat{P}^*(\hat{\tau}_n)$  does not differ from the ideal bootstrap  $P$ -value  $P^*(\hat{\tau}_n)$ . More precisely, the procedure works as follows:

1. Compute  $\hat{\tau}_n$ , set  $B = B' = 99$ , and the statistics  $\tau_j^*, j = 1, \dots, B$ .
2. Compute the estimated bootstrap  $P$ -value  $\hat{P}^*(\hat{\tau}_n)$  based on the  $B$  bootstrap samples. Depending whether  $\hat{P}^*(\hat{\tau}_n) < \alpha$  or  $\hat{P}^*(\hat{\tau}_n) > \alpha$ , test either the hypothesis that  $P^*(\hat{\tau}_n) \geq \alpha$  or  $P^*(\hat{\tau}_n) \leq \alpha$ .<sup>1</sup> If the hypothesis is rejected, then stop, else go to step 3.
3. Set  $B = 2B' + 1$ . If  $B$  is too large, e.g.  $B > 12,799$  then stop, else calculate  $\tau_j^*$  for a further  $B' + 1$  bootstrap samples, set  $B' = B$  and return to step 2.

**2.2 Data Generating Processes and Estimation Procedure**

In the simulations the Data Generating Process (DGP) follows

1. either an ARCH(2) model

$$y_t = \sigma_t \varepsilon_t, \quad \sigma_t^2 = b_0 + b_1 y_{t-1}^2 + b_2 y_{t-2}^2 \tag{2.2}$$

with

$$b_0 = 0.1, \quad b_1 = 0.2, \quad b_2 = 0.1 \tag{2.3}$$

or with (c.f. DGP 5)

$$b_0 = 0.1, \quad b_1 = 0.1, \quad b_2 = 0.4. \tag{2.4}$$

2. or a GARCH(1, 1) model

$$y_t = \sigma_t \varepsilon_t, \quad \sigma_t^2 = b_0 + b_1 y_{t-1}^2 + a_1 \sigma_{t-1}^2 \tag{2.5}$$

with

$$b_0 = 0.1, \quad b_1 = 0.3, \quad a_1 = 0.3. \tag{2.6}$$

*Remark 2.1* A sufficient condition for the stationarity of the process (2.2) and  $Ey_t^4 < \infty$  is  $(E\varepsilon_t^4)^{1/2}(b_1 + b_2) < 1$ , see e.g. Giraitis *et al.* (2000). For the parameters (2.3) we thus require that  $E\varepsilon_t^4 < (10/3)^2 = 11.11$  (for standard normal  $\varepsilon_t, E\varepsilon_t^4 = 3$ ). An analogous sufficient condition for the process (2.5) with standard normal  $\varepsilon_t$  is  $a_1^2 + 2a_1b_1 + 3b_1^2 < 1$ , see e.g. He and Teräsvirta (1999).

We consider the following Data Generating Processes. For each process we also report in Table V (see Appendix) the approximate skewness and kurtosis which has been computed by generating a realization of length  $10^6$  and finding the ratio of the fourth sample moment over the square of the second sample moment.

---

<sup>1</sup>This is done by using the Normal approximation of the Binomial distribution.

- DGP 0 The  $\varepsilon_t$  in (2.2) and (2.5) are i.i.d. standard normal. The models (2.2) and (2.5) are respectively denoted DGP 0A and DGP 0B.
- DGP 1 Consider i.i.d.  $\eta_t \sim \chi^2(k)$ . Thus  $E(\eta_t) = k$ , and  $\text{Var}(\eta_t) = 2k$ . We set  $k = 5$ , and define the  $\varepsilon_t$  in (2.2) by  $\varepsilon_t = (\eta_t - k)/\sqrt{2k}$ . The distribution of the  $\varepsilon_t$  has mean zero, but is asymmetric. Since the empirical distribution of asset returns is skewed, see *e.g.* Singleton and Wingender (1986), Badrinath and Chatterjee (1988), Engle and González-Rivera (1991), this case may be of practical relevance; it can be expected that the skewness of the  $\varepsilon_t$  will be reflected in the distribution of the  $y_t$ . Perhaps a more appropriate approach to modeling the skewness of the returns is to use a model which transforms symmetric innovations into asymmetric observations, see *e.g.* the asymmetric power ARCH of Ding *et al.* (1993) and the threshold GARCH of Zakoian (1994), but since in this study we use only standard GARCH models as archetypes, we experiment with skewed innovations.
- DGP 2 The  $\varepsilon_t$  in (2.2) are independent but  $\varepsilon_t \sim N(0, 1)$  for the first half of the sample, and  $\varepsilon_t \sim N(0, 2)$  for the second half. This alternative corresponds to a change-point (a break) in the distribution of the  $\varepsilon_t$ .
- DGP 3 The  $\varepsilon_t$  in (2.2) follow a linear ARCH(1), LARCH(1), process introduced by Giraitis *et al.* (2000), *i.e.*

$$\varepsilon_t = \zeta_t \xi_t, \quad \zeta_t = a_0 + a_1 \varepsilon_{t-1} \quad (2.7)$$

Thus  $y_t$  is the composition of an ARCH process and a LARCH process. For the LARCH(1) which forms the innovations sequence of the ARCH(2) defined by (2.2), we set  $a_0 = 1.0$ ,  $a_1 = 0.2$  and used standard normal innovations for the sequence  $\{\xi_t\}$ . This alternative corresponds to a remaining heteroskedastic structure in the innovations process. Under the null hypothesis  $a_1 = 0$ , we obtain the DGP 0A. We choose a LARCH(1) process for the nested process as taking a standard ARCH(1) process yields a process  $\{y_t\}$  with a very large kurtosis for  $a_1 \geq 0.2$ . Further details are given in the Appendix. For this DGP, we use as a benchmark the test for correct specification of the heteroskedastic functional form  $Q(M)$  by Li and Mak (1994) which is asymptotically equivalent to the test by Lundbergh and Teräsvirta (2002) defined for the surimposition of standard ARCH/GARCH processes.

- DGP 4 Consider i.i.d.  $u_t \sim t(k)$ , *i.e.*  $t$  distributed with  $k$  degrees of freedom. We set in (2.2)  $\varepsilon_t = \sqrt{k/(k-2)}u_t$  with  $k = 6$ . Then the  $\varepsilon_t$  have a symmetric unit variance distribution which is leptokurtic, *i.e.* its tails are thicker than the ones of a normal distribution.
- DGP 5 We also consider the case of misspecification of the functional form of the model: the DGP is still ARCH(2) (2.2) but with  $b_0 = 0.1$ ,  $b_1 = 0.1$ , and  $b_2 = 0.4$ . The estimates  $\hat{b}_j$  are however computed assuming that the observed series follows an ARCH(1) model. For this DGP, we also use as a benchmark the test for correct specification of the conditional variance by Li and Mak (1994).
- DGP 6 We consider a GARCH(1, 1) DGP, with the assumption that the error terms are  $t$  distributed with 7 degrees of freedom. The parameters  $b_1$  and  $a_1$  are the same as for DGP 0B. This choice has been motivated by the fact that Bollerslev (1987) and Teräsvirta (1996) observed that the observed kurtosis of financial data is better fitted with  $t(7)$  distributed error terms. These simulation results are of interest for our application on real data, as we will consider series of returns on Foreign Exchange (FX) rates at daily frequency.
- DGP 7 We consider that the error terms in (2.2) are i.i.d. and follow a Laplace distribution  $\text{Lap}(\sigma)$ , the density of which is  $f(t) = (2\sigma)^{-1/2} \exp\{-|t|/\sigma\}$ . For this DGP, we set

$\sigma = 2^{-1/2}$ , thus  $\varepsilon_t \sim \text{Lap}(2^{-1/2})$ , which is the two-sided exponential distribution with mean zero and variance equal to one.

DGP 8 We consider that the  $\varepsilon_t$  in (2.2) are independent but  $\varepsilon_t \sim N(0, 1)$  for the first half of the sample, and  $\varepsilon_t \sim \text{Lap}(2^{-1/2})$  for the second half.

DGP 9 We consider that the  $\varepsilon_t$  in (2.2) are independent but  $\varepsilon_t \sim N(0, 1)$  for the first half of the sample, and  $\varepsilon_t \sim t(7)$  for the second half.

DGP 10 We consider that the  $\varepsilon_t$  in (2.2) are independent but  $\varepsilon_t \sim N(0, 1)$  for the first half of the sample, and  $\varepsilon_t$  follow a centralized and standardized  $\chi^2(5)$  for the second half.

The parameters  $\hat{b}_j$  are estimated by the pseudo maximum likelihood (PML) method, *i.e.* the estimates are the parameter values which maximize the likelihood computed under the assumption that the innovations  $\varepsilon_t$  are independent and normally distributed. In that case, the log-likelihood function is:

$$\mathcal{L}(b; y) = -\frac{1}{2} \sum_{1 \leq t \leq n} \left( \ln(2\pi) + \ln(\sigma_t^2) + \frac{y_t^2}{\sigma_t^2} \right) \tag{2.8}$$

where  $n$  denotes the number of observations. In this paper we consider  $n = 200$  and  $n = 400$ .

As explained in Subsection 2.1, for a number of bootstraps  $B = 1499$ , an experiment for each DGP above involves 15 million PML estimation procedures. For each of them we check whether the Hessian matrix  $\partial^2 \mathcal{L}(b; y) / (\partial b \partial b^T)$  evaluated at the PML estimates  $\hat{b}_j$  is positive definite, and we reject the estimation results if this condition is not satisfied. The estimates were obtained using the sequential QP algorithm from the NAG library, which allows us to maximize a function subject to linear and nonlinear constraints, *e.g.* the constraint of fourth moment existence. Although the optimization procedure makes use of the analytic gradient of the log-likelihood function, see *e.g.* Bollerslev (1986), the whole set of simulations for the first version of this work required over 15,000 hours of computing time on several computers. A recent experiment for a sample of 200 observations, carried out on a Pentium IV with a 2.2 GHz clock, required around 55 seconds for  $B$  equal to 1499: the testing procedure is thus very affordable on modern computers.

### 3 PRESENTATION OF THE RESULTS WITH GRAPHICAL METHODS

The availability of the bootstrap  $P$ -values allows us to present the simulation results using the graphical methods advocated by Davidson and MacKinnon (1998). For a given test, define the empirical distribution function (EDF) of the  $P$ -values  $P_j, j = 1, \dots, N$  as

$$\hat{F}(x_i) \equiv \frac{1}{N} \sum_{j=1}^N \mathbb{I}[P_j \leq x_i], \tag{3.1}$$

where  $\{x_i\}$  is a fine grid of points in  $[0, 1]$ . To compare the size, we display the  $P$ -value discrepancy plots, *i.e.* the plots of  $\hat{F}(x_i) - x_i$  against  $x_i$ . Under the null hypothesis, the  $P$ -value is uniformly distributed on the interval  $[0, 1]$ , therefore the  $P$ -value discrepancy plots should be close to the  $0^\circ$  line. We estimate  $\hat{F}(\cdot)$  on a grid of 215 points,  $x_i = 0.001, 0.002, \dots, 0.010, 0.015, \dots, 0.990, 0.991, \dots, 0.999$ . Even with a large number of replications, *e.g.* 100,000,  $P$ -value discrepancy plots are jagged due to the experimental randomness. Therefore, when the number of replications is smaller because, as in our case, the statistics are computing time consuming, a smooth version of the  $P$ -value discrepancy plots is desirable. Denote by  $v_i$  the discrepancy, *i.e.*  $v_i = F(x_i) - x_i$ . A smoothed

curve is obtained by regressing  $v_i$  on smooth function of  $x_i$ . Davidson and MacKinnon (1998) suggested the following regression:

$$v_i = \sum_{l=1}^L \gamma_l \sin(l\pi x_i) + u_i \tag{3.2}$$

which can be written in matrix notation as

$$\mathbf{v} = \mathbf{Z}\boldsymbol{\gamma} + \mathbf{u} \text{ with } E(\mathbf{u}\mathbf{u}^T) =: \boldsymbol{\Omega},$$

where the vector  $\boldsymbol{\gamma}$  is obtained by using the feasible Generalized Least Squares estimator. The truncation order  $L$  is selected using the Akaike Information Criterion, the log-likelihood function associated with the GLS estimation being:

$$-\frac{m}{2} \log(2\pi) + \frac{1}{2} |\boldsymbol{\Omega}^{-1}| - \frac{1}{2} (\mathbf{v} - \mathbf{Z}\boldsymbol{\gamma})^T \boldsymbol{\Omega}^{-1} (\mathbf{v} - \mathbf{Z}\boldsymbol{\gamma})$$

Confidence bands are obtained by adding to the smoothed curve twice the square root of the diagonal elements of the covariance matrix  $\mathbf{Z}(\mathbf{Z}^T \boldsymbol{\Omega}^{-1} \mathbf{Z})^{-1} \mathbf{Z}^T$ . The exact form of  $\boldsymbol{\Omega}^{-1}$  and other details are presented in Davidson and MacKinnon (1998).

Figure 1 (see Appendix) shows the  $P$ -value discrepancy plots and the smoothed  $P$ -value discrepancy plots for the three tests for the DGP 0B (The graphs for the DGP 0A are very similar). Since these plots are close to the  $0^\circ$  line, which is inside the confidence bands, the bootstrap tests are seen to have correct size even for small sample sizes.

For representing the power of the test for each DGP, we use a graphical representation of the size–power function, denoted by  $\eta(\alpha)$ , which is the probability under that DGP that the test will reject the null hypothesis when the nominal rejection probability under the null hypothesis is  $\alpha$ . For a powerful test,  $\eta(\alpha)$ , is much larger than  $\alpha$  for small  $\alpha$ 's and so the curve  $\eta(\alpha)$  is  $\Gamma$  shaped. For a test with small power the curve  $\eta(\alpha)$  is close to the  $45^\circ$  line.

The size–power curve graphs the power of a test against its nominal level. The EDF of the  $P$ -values,  $\hat{F}(\alpha)$ , which is an estimator of the cumulative distribution function of the  $P$ -values, is an approximation to  $\eta(\alpha)$ . However, it is more informative to plot  $\hat{F}(\alpha)$  against the real size of the test. Thus the size–power trade-off curves in Figures 2–5 (see Appendix) are the graphs of the EDF  $\hat{F}(\alpha)$  of the  $P$ -values under the alternative hypothesis against the EDF  $\hat{F}(\alpha)$  of the  $P$ -values under the null hypothesis. The null hypothesis is always DGP 0A, except for DGP 6, where the null hypothesis is DGP 0B.

Figures 2–5 (see Appendix) plot the size–power functions for selected DGP's: DGP 4, 5, 6, 8 with the  $45^\circ$  line.<sup>2</sup> For all DGP's considered in our simulation study, the CVM test has more power than the NCVm test, which in turn is more powerful than the KS test. The only exception are DGP 7 and DGP 8, (Laplace innovations) where the NCVm test is slightly more powerful than the CVM test.

For the sample sizes considered, the tests have a high power for the DGP 1, 4, 6, 7, 8, *i.e.* for the asymmetric and leptokurtic alternatives. The curves  $\eta(\alpha)$  approach a  $\Gamma$  shape when the sample size increases from 200 to 400 observations.

We also calculated the empirical size and power of the test for three levels: 10%, 5%, and 1%. Table VI (see Appendix) shows that even for small samples of 200 observations, the empirical size is very close to the nominal size.

The comparison with the test by Li and Mak (1994) is of interest: while the latter test has no power at all for the sample sizes and the parameter chosen for DGP 3, *i.e.* the power of the test

<sup>2</sup>The full set of size–power curves is available upon request.



is equal to its size, even when considering bootstrap tests for this statistic, the portmanteau test is more powerful than our tests based on the EDF for DGP 5. For a sample of 200 observations, the empirical frequencies of rejection of the bootstrapped  $Q(M)$  statistic for a test of size 5% with  $M = 3, 4, 6$ , are respectively equal to 86.90%, 84.10%, and 78.65%. For sample sizes of 400 observations, the frequencies of rejection increase to 99.14%, 98.82%, and 97.84%.

A comparison with the tests for change-point developed in Kokoszka and Teysnière (2002), shows that the tests developed here have more power for detecting changes in the distribution of the innovations, *i.e.* DGP 8, 9, 10.

#### 4 APPLICATION TO FOREIGN EXCHANGE RATES RETURNS SERIES

We consider here an application of the three bootstrap tests to the series of squared residuals obtained from the estimation of a GARCH(1, 1) on series of Foreign Exchange (FX) rates returns  $R_t = 100 \times \log(P_t/P_{t-1})$ , where  $P_t$  denotes the spot FX rate at time  $t$ . We consider the series of daily US Dollar/Deutschmark (USD–DEM) and US Dollar/British Pound (USD–GBP) FX rates returns from April 1979, *i.e.* after the inception of the European Monetary System. At first sight, this application is not relevant since empirical evidence has shown that the sums of the coefficients of a GARCH(1, 1) are close to one, and thus violate the condition for stationarity of a GARCH(1, 1) given above and in He and Teräsvirta (1999). However, Mikosch and Stărică (1999; 2002a,b) have argued that this property can be due to the concatenation of several GARCH(1, 1) processes, which can also be viewed as a non-homogeneous GARCH(1, 1) process, *i.e.* with changing coefficients. Indeed, estimating a GARCH(1, 1) model on shorter samples gives results which satisfy the stationarity condition.

We consider 2 samples of 1000 observations. The model specification is:

$$R_t = \mu + \sigma_t \varepsilon_t, \quad \sigma_t^2 = b_0 + a_1 \varepsilon_{t-1}^2 + b_1 \sigma_{t-1}^2, \quad \varepsilon_t \sim \text{i.i.d. } N(0, 1). \tag{4.1}$$

Tables I and III below display the estimated parameters of the GARCH(1, 1) while Tables II and IV display the bootstrap CVM, NCVM and KS statistics  $\hat{\tau}$ , and the associated bootstrap

TABLE I Estimated Parameters USD–DEM. Heteroskedastic Robust  $t$ -Statistics are in Parentheses.

Series	$b_0$	$b_1$	$a_1$	$\mu$
1	0.0148 (0.9596)	0.7599 (7.2343)	0.1844 (2.3757)	0.0204 (0.9842)
2	0.0139 (1.2151)	0.8537 (15.2439)	0.1231 (2.0036)	-0.0128 (-0.3785)

TABLE II Bootstrap Test Statistics and their Critical Values. USD–DEM.

	Statistics			Statistics		
	CVM	NCVM	KS	CVM	NCVM	KS
		Series 1			Series 2	
$\hat{\tau}$	17.2911	4.4413	3.7807	7.66543	2.1729	2.88938
1%	0.9835	0.2392	1.3809	1.0497	0.2363	1.38856
5%	0.6999	0.1561	1.1436	0.6720	0.1545	1.14692
10%	0.5810	0.1219	1.0255	0.5666	0.1200	1.02789

TABLE III Estimated Parameters USD–GBP. Heteroskedastic Robust *t*-Statistics are in Parentheses.

<i>Series</i>	$b_0$	$b_1$	$a_1$	$\mu$
1	0.0172 (1.4735)	0.8725 (15.9346)	0.0965 (1.9942)	−0.0391 (−1.0204)
2	0.0253 (1.4374)	0.8942 (17.8774)	0.04715 (1.7436)	0.0430 (1.0703)

TABLE IV Bootstrap Test Statistics and their Critical Values. USD–GBP.

	<i>Statistics</i>			<i>Statistics</i>		
	<i>CVM</i>	<i>NCVM</i>	<i>KS</i>	<i>CVM</i>	<i>NCVM</i>	<i>KS</i>
		Series 1			Series 2	
$\hat{\tau}$	4.1374	1.1634	2.2537	8.5529	2.1097	2.8452
1%	0.9961	0.2449	1.3503	1.0470	0.2480	1.4074
5%	0.6918	0.1563	1.1468	0.6888	0.1559	1.1505
10%	0.5837	0.1249	1.0524	0.5496	0.1147	1.0213

critical values for a test of size 1%, 5%, and 10%. For all the considered series, the hypothesis of normality of the innovations is rejected, even for the size 1%.

## 5 CONCLUSIONS

We have examined several goodness-of-fit tests based on the residual empirical process of an ARCH sequence. Unlike for linear time series models, asymptotic tests cannot be readily constructed, but bootstrap tests are relatively easy to perform. The bootstrap tests we studied have correct size and can detect a misspecified probability distribution of the innovations (errors) in an ARCH model with high probability. They can detect an incorrect postulated order of the model with moderate success, and are practically unable to detect a structural break in the model parameters. For the detection of a misspecified distribution of the innovations, the tests have good power even for samples of size 200 observations. In situations where the tests are applicable, the CVM test performs better than the NCVM and KS tests, except when the error terms or part of the error terms follow a Laplace distribution. In that case, the NCVM test is slightly more powerful than the CVM test. An application to returns on FX rates shows that if a GARCH(1, 1) model is postulated, normal errors should not be assumed.

## 6 HISTORICAL AND MATHEMATICAL APPENDIX

In this section we provide some mathematical and historical background for the goodness-of-fit procedures considered in this paper.

Suppose we observe a sample  $\varepsilon_1, \dots, \varepsilon_n$  of i.i.d. observations with cumulative distribution function  $F$  which we assume for simplicity to be strictly increasing and absolutely continuous. Then the empirical process of the  $\varepsilon_i$  is

$$\alpha_n(x) = n^{-1/2} \sum_{i=1}^n \{\mathbb{I}[\varepsilon_i \leq x] - F(x)\}. \quad (6.1)$$

It is well known, see *e.g.* Shorack and Wellner (1986), that  $\alpha_n(x)$  converges in  $D(-\infty, \infty)$  to  $B^0(F(x))$ , where  $B^0$  is a Brownian bridge. The latter convergence is to be interpreted as the convergence of  $\alpha_n(F^{-1}(\cdot))$  to  $B^0(\cdot)$  in  $D[0, 1]$ . Relation (6.1) is used to construct classical goodness-of-fit tests. For example, a variant of the Cramér–von Mises test (NCVM in our notation) is based on the convergence

$$\int_{-\infty}^{\infty} |\alpha_n(x)|^2 F(dx) \xrightarrow{d} \int_{-\infty}^{\infty} |B^0(F(x))|^2 F(dx) \stackrel{d}{=} \int_0^1 |B^0(u)|^2 du, \tag{6.2}$$

the distribution of the last integral being tabulated.

In time series models, the innovations (noise sequence) are not observable and we have to work with residuals. To illustrate the idea, suppose we observe a sample  $x_1, \dots, x_n$  from an autoregressive model of order  $p$  which, in a sense, is a linear analogue of the ARCH( $p$ ) model. Thus, we assume that the observations satisfy

$$x_t = \sum_{j=1}^p \phi_j x_{t-j} + \varepsilon_t, \tag{6.3}$$

where the  $\varepsilon_t$  are as above and, in addition, have mean zero and unit variance. Suppose we estimate the coefficients  $\phi_j$  using estimators such that for each  $j$ ,  $n^{1/2}(\hat{\phi}_j - \phi_j)$  is bounded in probability (all reasonable estimators satisfy this condition) and construct the residuals

$$\hat{\varepsilon}_t = x_t - \sum_{j=1}^p \hat{\phi}_j x_{t-j}, \quad t = p + 1, \dots, n. \tag{6.4}$$

Then the empirical process of these residuals, defined analogously to (6.1) with the  $\varepsilon_i$  replaced by the  $\hat{\varepsilon}_i$ , still converges to  $B^0(F(\cdot))$ . This property holds for more general linear models and under weaker assumptions on the  $\varepsilon_i$ , see Chapter 7 of Koul (1992) and references therein.

The situation is different for ARCH models. It follows from Boldin (1998) that the difference between the empirical process of the residuals from an ARCH(1) and a transformed Brownian bridge tends to a normal random variable. Horváth *et al.* (2001) considered a general ARCH( $p$ ) models and squared residuals and showed the process (1.3) converges in  $D[0, \infty)$  to  $B^0(G(t)) + tg(t)\xi$ , where  $G$  and  $g$  are, respectively, the cumulative distribution function and the density of  $\varepsilon_0^2$  and  $\xi$  is a normal random variable which is correlated in an intricate way with the process  $B^0(G(\cdot))$ . Even though the joint covariance structure is known, it does not lead to a simple convergence like (6.2) and so, unlike for linear models, asymptotic tests cannot be readily constructed, but the existence of a limiting distribution suggests the applicability of bootstrap tests. For extensions to GARCH( $p, q$ ) processes and more general ARCH processes we refer, respectively, to Berkes and Horváth (2003) and Section 8.3 of Koul (2002).

**Acknowledgement**

We thank the referee for constructive advice which led to a significant improvement of the paper, and Jeroen Rombouts for a careful reading of the paper. This text presents research results of the Belgian Program on Interuniversity Poles of Attraction initiated by the

Belgian State, Prime Minister's Office, Science Policy Programming. The scientific responsibility is assumed by the authors.

## References

- Andrews, D. W. K. and Buchinsky, M. (2000). A three-step method for choosing the number of bootstrap repetitions. *Econometrica*, **68**, 23–51.
- Badrinath, S. G. and Chatterjee, S. (1988). On measuring skewness and elongation in common stock returns distributions: The case of market index. *Journal of Business*, **61**, 451–472.
- Berkes, I. and Horváth, L. (2003). Limit results for the empirical process of squared residuals in GARCH models. *Stochastic Processes and their Applications*, **105**, 279–298.
- Berkes, I., Horváth, L. and Kokoszka, P. S. (2003). Asymptotics for GARCH squared residual correlations. *Econometric Theory*, **19**, 515–540.
- Berkes, I., Horváth, L. and Kokoszka, P. S. (2004). A weighted goodness-of-fit test for GARCH (1,1) specification. *Lithuanian Mathematics Journal*, **44**, 1–17.
- Boldin, M. V. (1998). On residual empirical distribution functions in ARCH models with applications to testing and estimation. *Mitt. Math. Sem. Giessen*, **235**, 49–66.
- Bollerslev, T. (1986). Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, **31**, 307–327.
- Bollerslev, T. (1987). A conditional heteroskedastic time series model for speculative prices and rates of return. *Review of Economics and Statistics*, **69**, 542–547.
- Davidson, R. and MacKinnon, J. G. (1998). Graphical methods for investigating the size and power of hypothesis tests. *The Manchester School*, **66**, 1–26.
- Davidson, R. and MacKinnon, J. G. (1999). The size distortion of bootstrap tests. *Econometric Theory*, **15**, 361–376.
- Davidson, R. and MacKinnon, J. G. (2000). Bootstrap tests: How many bootstraps? *Econometric Reviews*, **19**, 55–68.
- Davis, R. A. and Mikosch, T. (1998). Limit theory for the sample act of stationary process with heavy tails with applications to ARCH. *The Annals of Statistics*, **26**, 2049–2080.
- Ding, Z., Granger, C. W. J. and Engle, R. (1993). A long memory property of stock marked returns and a new model. *Journal of Empirical Finance*, **1**, 83–106.
- Dufour, J. M. and Khalaf, L. (2001). Monte Carlo tests in econometrics. In: Baltagi, B. (Ed.), *Companion to Theoretical Econometrics*. Blackwell, Oxford, pp. 494–519.
- Dufour, J. M. and Kiviet, J. F. (1998). Exact inference methods for first-order autoregressive distributed lag models. *Econometrics*, **66**, 79–104.
- Engle, R. F. and González-Rivera, G. (1991). Semiparametric ARCH models. *Journal of Business and Economic Statistics*, **9**, 345–359.
- Giraitis, L., Kokoszka, P. S. and Leipus, R. (2000). Stationary ARCH models: Dependence structure and Central Limit Theorem. *Econometric Theory*, **16**, 3–22.
- Giraitis, L., Robinson, P. and Surgailis, D. (2000). A model for long memory conditional heteroskedasticity. *The Annals of Applied Probability*, **10**, 1002–1024.
- He, C. and Teräsvirta, T. (1999). Fourth moment structure of the GARCH( $p, q$ ) process. *Econometric Theory*, **15**, 824–846.
- Horváth, L. and Kokoszka, P. S. (2001). Large sample distribution of ARCH( $p$ ) squared residual correlations. *Econometric Theory*, **17**, 283–295.
- Horváth, L., Kokoszka, P. S. and Teyssiére, G. (2001). Empirical process of the squared residuals of an ARCH sequence. *The Annals of Statistics*, **29**, 445–469.
- Kokoszka, P. S. and Teyssiére, G. (2002). Change point detection in GARCH models: Asymptotic and bootstrap tests, *Technical Report*. Utah State University, Logan. Working paper available at <http://stat.usu.edu/~piotr/research.html>.
- Koul, H. L. (1992). *Weighted Empirical and Linear Models*. IMS, Hayward, California.
- Koul, H. L. (2002). *Weighted Empirical Processes in Dynamic Nonlinear Models*, Lecture Notes in Statistics, Springer Verlag.
- Li, W. K. and Mak, T. K. (1994). On the squared residual autocorrelations in non-linear time series with conditional heteroskedasticity. *Journal of Time Series Analysis*, **15**, 627–636.
- Lundbergh, S. and Teräsvirta, T. (2002). Evaluating GARCH models. *Journal of Econometrics*, **110**, 417–435.
- Marsaglia, G. (1996). DIEHARD: A battery of tests of randomness. Available at <http://stat.fsu.edu/pub/diehard>.
- Mikosch, T. and Stărică, C. (1999). Change of structure in financial time series, long range dependence and the GARCH model, *Technical Report*. University of Groningen, preprint available at <http://www.math.ku.dk/~mikosch/preprint.html>.
- Mikosch, T. and Stărică, C. (2000). Limit theory for the sample autocorrelation and extremes of a GARCH(1, 1) process. *The Annals of Statistics*, **28**, 1427–1451.
- Mikosch, T. and Stărică, C. (2002a). Long-range dependence effects and ARCH modeling. In: Doukhan, P., Oppenheim, G. and Taqqu, M. S. (Eds.), *Theory and Applications of Long-range Dependence*. Birkhäuser, Boston, pp. 439–459.

- Mikosch, T. and Stărică, C. (2002b). Non-stationarities in financial time series: The long range dependence and the IGARCH effects, *Technical Report*, University of Copenhagen.
- Shorack, G. R. and Wellner, J. A. (1986). *Empirical Processes with Applications to Statistics*. Wiley, New York.
- Singleton, J. C. and Wingender, J. (1986). Skewness persistence in common stock returns. *Journal of Financial and Quantitative Analysis*, **21**, 335–341.
- Teräsvirta, T. (1996). Two stylized facts and the GARCH(1, 1) model. Stockholm School of Economics, Working Paper Series in Economics and Finance 96.
- Zakoian, J. M. (1994). Threshold heteroskedastic models. *Journal of Economic Dynamics and Control*, **18**, 931–955.

## APPENDIX

### A Tables

- DGP 0A: ARCH(2) with  $N(0, 1)$  errors, *i.e.* the null hypothesis
- DGP 0B: GARCH(1, 1) with  $N(0, 1)$  errors, *i.e.* the null hypothesis
- DGP 1: ARCH(2) with centralized and standardized  $\chi^2(5)$  errors
- DGP 2: ARCH(2) with  $N(0, 1)$  errors for the first half of observations, and  $N(0, 2)$  errors for the second half
- DGP 3: surimposition of two ARCH processes, *i.e.* a LARCH(1) nested in an ARCH(2)
- DGP 4: ARCH(2) with standardized  $t(6)$  errors
- DGP 5: ARCH(1) model estimated on an ARCH(2) DGP
- DGP 6: GARCH(1, 1) with standardized  $t(7)$  errors
- DGP 7: ARCH(2) with Lap( $2^{-1/2}$ ) errors
- DGP 8: ARCH(2) with  $N(0, 1)$  errors for the first half of observations, and Lap( $2^{-1/2}$ ) errors for the second half
- DGP 9: ARCH(2) with  $N(0, 1)$  errors for the first half of observations, and standardized  $t(7)$  errors for the second half
- DGP 10: ARCH(2) with  $N(0, 1)$  errors for the first half of observations, and centralized and standardized  $\chi^2(5)$  errors for the second half

We have also considered for DGP 3 the case of an ARCH(1) nested in an ARCH(2), *i.e.*

$$\varepsilon_t = \zeta_t \check{\zeta}_t, \quad \check{\zeta}_t^2 = a_0 + a_1 \varepsilon_{t-1}^2 \tag{A.1}$$

When  $a_1 = 0.1$ , the kurtosis of the process computed on a sample of  $10^6$  observation is equal to 6.6745, but increases to 4012.12 for  $a_1 = 0.2$ . For  $a_1 \geq 0.3$ , the kurtosis is huge. Thus, we stick to the case of a LARCH(1) nested in an ARCH(2) as for  $a_1 = 0.2$  the kurtosis is equal to 7.20, and becomes huge for  $a_1 \geq 0.4$ .

TABLE V Empirical Skewness and Kurtosis for the Different DGP.

<i>DGP</i>	<i>Skewness</i>	<i>Kurtosis</i>
DGP 0A	0.0004	3.4111
DGP 0B	0.0005	4.0630
DGP 1	1.3694	6.8832
DGP 2	0.0015	8.1750
DGP 3	0.0108	7.2059
DGP 4	-0.0227	7.2830
DGP 5	-0.0027	5.0889
DGP 6	0.0369	6.9694
DGP 7	-0.0073	8.8038
DGP 8	0.0002	6.0825
DGP 9	-0.0185	5.9228
DGP 10	0.6812	5.2007

TABLE VI Empirical Frequencies of Rejection (200 Observations).

<i>DGP</i>	<i>CVM</i>			<i>NCVM</i>			<i>KS</i>		
	10%	5%	1%	10%	5%	1%	10%	5%	1%
DGP 0A	0.1020	0.0508	0.0108	0.0995	0.0505	0.0110	0.0984	0.0501	0.0101
DGP 0B	0.1032	0.0518	0.0113	0.1021	0.0518	0.0107	0.1015	0.0505	0.0105
DGP 1	0.8590	0.7698	0.5417	0.5577	0.4216	0.1959	0.4530	0.3119	0.1277
DGP 2	0.1320	0.0764	0.0180	0.1265	0.0647	0.0133	0.1028	0.0502	0.0123
DGP 3	0.1868	0.1109	0.0383	0.1790	0.1086	0.0349	0.1431	0.0788	0.0221
DGP 4	0.7227	0.6323	0.4417	0.6907	0.5889	0.3602	0.5590	0.4330	0.2224
DGP 5	0.3897	0.3139	0.1909	0.3442	0.2587	0.1459	0.2626	0.1824	0.0849
DGP 6	0.5826	0.4762	0.2893	0.5422	0.4416	0.2548	0.4173	0.3108	0.1503
DGP 7	0.9897	0.9814	0.9453	0.9918	0.9841	0.9526	0.9840	0.9682	0.9063
DGP 8	0.6580	0.5616	0.3655	0.6821	0.5847	0.3849	0.6096	0.4916	0.2886
DGP 9	0.2564	0.1765	0.0812	0.2383	0.1578	0.0669	0.1731	0.1030	0.0376
DGP 10	0.3987	0.2850	0.1365	0.2403	0.1496	0.0530	0.1733	0.1000	0.0311

TABLE VII Empirical Frequencies of Rejection (400 Observations).

<i>DGP</i>	<i>CVM</i>			<i>NCVM</i>			<i>KS</i>		
	10%	5%	1%	10%	5%	1%	10%	5%	1%
DGP 0A	0.1019	0.0526	0.0104	0.1017	0.0540	0.0104	0.1015	0.0519	0.0133
DGP 0B	0.1015	0.0487	0.0103	0.0984	0.0494	0.0105	0.0987	0.0484	0.0114
DGP 1	0.9915	0.9801	0.9280	0.8285	0.7157	0.4641	0.8104	0.6820	0.4110
DGP 2	0.2006	0.1175	0.0374	0.1540	0.0870	0.0258	0.1262	0.0658	0.0191
DGP 3	0.3171	0.2241	0.0979	0.3183	0.2153	0.0948	0.2421	0.1572	0.0591
DGP 4	0.9323	0.8970	0.7979	0.9172	0.8720	0.7521	0.8371	0.7598	0.5665
DGP 5	0.6009	0.5181	0.3637	0.5414	0.4535	0.3095	0.4303	0.3348	0.1942
DGP 6	0.8497	0.7869	0.6287	0.8215	0.7460	0.5695	0.7030	0.5992	0.3819
DGP 7	0.9999	0.9998	0.9996	1.0000	1.0000	0.9998	0.9999	0.9999	0.9993
DGP 8	0.9255	0.8757	0.7398	0.9381	0.8962	0.7661	0.8976	0.8371	0.6592
DGP 9	0.5363	0.4366	0.2700	0.5031	0.4002	0.2355	0.3932	0.2834	0.1348
DGP 10	0.5604	0.4379	0.2364	0.3018	0.1916	0.0781	0.2212	0.1257	0.0382

**B** *P*-value discrepancy plots

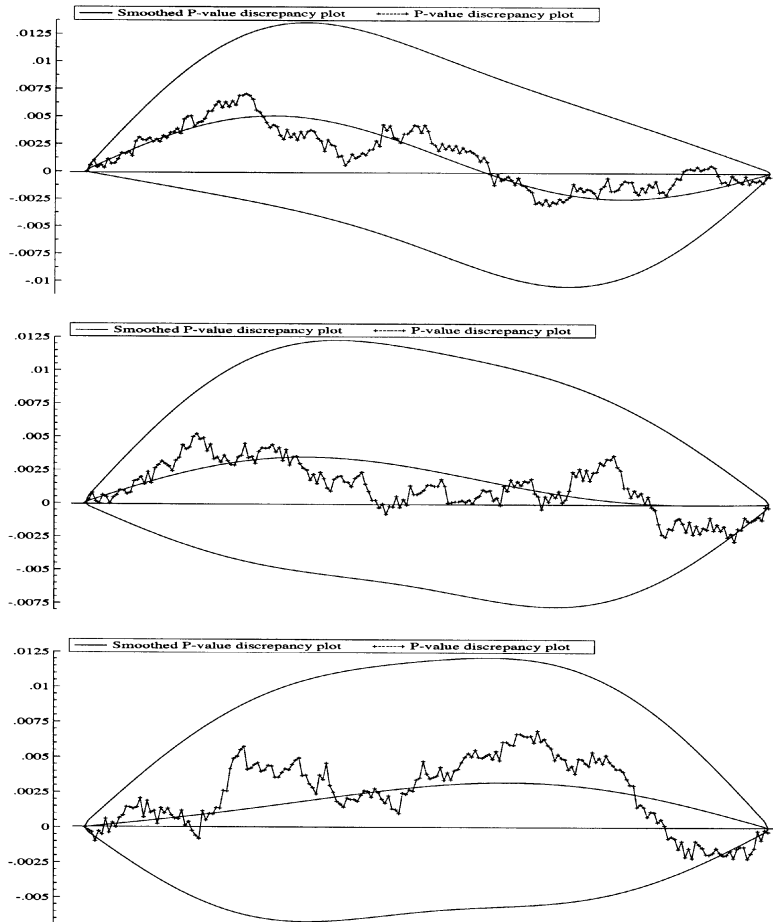


FIGURE 1 *P*-value discrepancy plots and their smoothed versions with confidence bands. Sample size is 200. Top – CVM, middle – NCVM, bottom – KS test.

### C Size–power curves

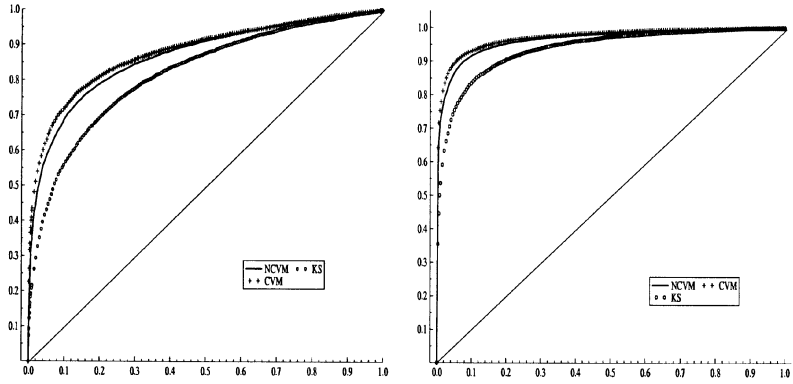


FIGURE 2 DGP 4 (leptokurtic errors). Left: 200 observations; right: 400 observations.

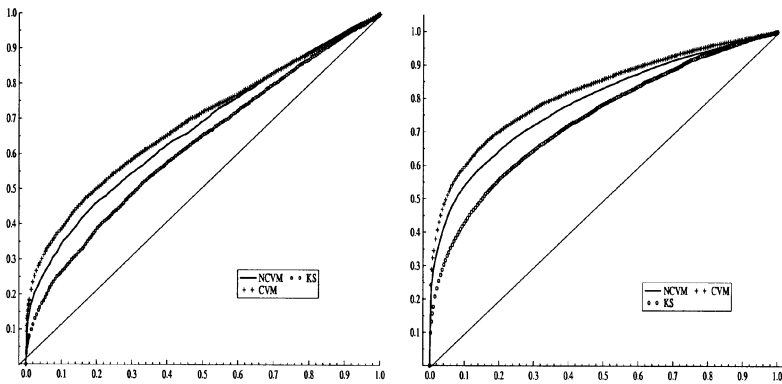


FIGURE 3 DGP 5 (incorrect order). Left: 200 observations; right: 400 observations.

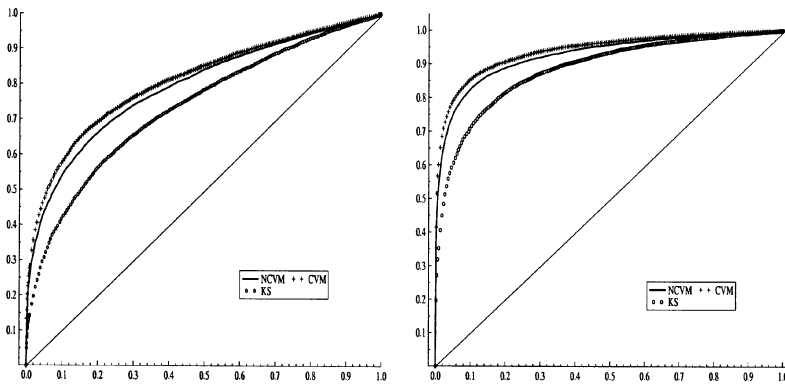


FIGURE 4 DGP 6 (GARCH(1, 1) with  $t(7)$  innovations). Left: 200 observations; right: 400 observations.



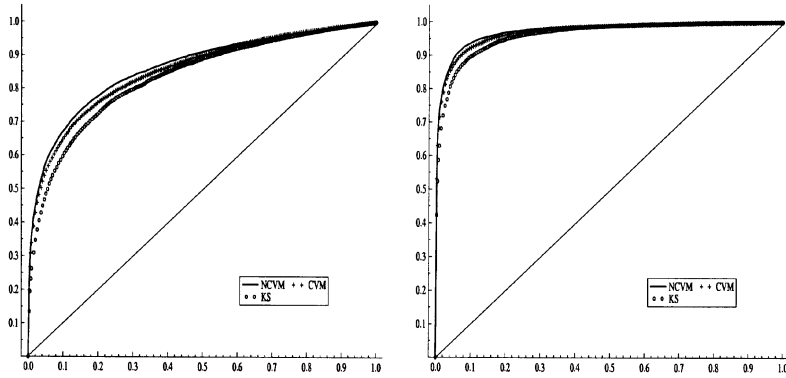


FIGURE 5 DGP 8 (change point in the innovations). Left: 200 observations; right: 400 observations.

