

# Efficient Estimation of Conditional Variance Functions in Stochastic Regression

By JIANQING FAN

Department of Statistics, University of California, Los Angeles, CA 90095, USA

e-mail: jfan@stat.unc.edu

AND QIWEI YAO

Institute of Mathematics and Statistics, University of Kent

Canterbury, Kent CT2 7NF, UK

e-mail: Q.Yao@ukc.ac.uk

## SUMMARY

Conditional heteroscedasticity has been often used in modelling and understanding the variability of statistical data. Under a general setup which includes the nonlinear time series model as a special case, we propose an efficient and adaptive method for estimating the conditional variance. The basic idea is to apply a local linear regression to the squared residuals. We demonstrate that without knowing the regression function, we can estimate the conditional variance asymptotically as well as if the regression were given. This asymptotic result, established under the assumption that the observations are made from a strictly stationary and absolutely regular process, is also verified via simulation. Further, the asymptotic result paves the way for adapting an automatic bandwidth selection scheme. An application with financial data illustrates the usefulness of the proposed techniques.

*Some key words:* Absolutely regular; ARCH; Conditional variance; Efficient estimator; Heteroscedasticity; Local linear regression; Nonlinear time series; Volatility.

# 1 INTRODUCTION

Many scientific studies depend on understanding the local variability of the data, which is often featured as the conditional variance or the volatility function in a statistical model. It is of common interest to estimate conditional variance functions in a variety of statistical applications such as measuring the volatility or risk in finance (Andersen and Lund, 1997; Gallant and Tauchen, 1997), monitoring the reliability in nonlinear prediction (Yao and Tong, 1994), identifying homoscedastic transforms in regression (Carroll and Ruppert, 1988), choosing optimal design and understanding residual pattern (Müller and Stadtmüller, 1987; Gasser et al., 1986), monitoring the signal-to-noise ratios in quality control of experimental design (Box, 1988) and so on. The problem can be mathematically formulated as follows.

Let  $\{(Y_i, X_i)\}$  be a two-dimensional strictly stationary process having the same marginal distribution as  $(Y, X)$ . Let  $m(x) = E(Y|X = x)$  and  $\sigma^2(x) = \text{var}(Y|X = x) \neq 0$ . We write a regression model of  $Y_i$  on  $X_i$  as

$$Y_i = m(X_i) + \sigma(X_i)\epsilon_i. \quad (1.1)$$

Then  $E(\epsilon_i|X_i) = 0$  and  $\text{var}(\epsilon_i|X_i) = 1$ , although the conditional distribution of  $\epsilon_i$  given  $X_i = x$  may still depend on  $x$ . For  $X_i = Y_{i-1}$ , (1.1) is an autoregressive conditional heteroscedastic (ARCH) time series model, and  $\sigma(\cdot)$  is called the volatility function (Engle, 1982). Connections of model (1.1) with one-factor diffusion model in finance will be discussed in Section 4. The aim of this paper is to derive an efficient fully-adaptive procedure for estimating  $\sigma^2(\cdot)$ .

Due to the simple decomposition  $\sigma^2(x) = E(Y^2|X = x) - m^2(x)$ , the following obvious and direct estimator is used:

$$\hat{\sigma}_d^2(x) = \hat{\nu}(x) - \{\hat{m}(x)\}^2, \quad (1.2)$$

where  $\hat{m}(\cdot)$  and  $\hat{\nu}(x)$  are respectively a regression estimator for  $m(\cdot)$  and  $\nu(x) \equiv E\{Y^2|X = x\}$  (Yao and Tong, 1994; Härdle and Tsybakov, 1997). However,  $\hat{\sigma}_d^2(\cdot)$  is not always non-negative, especially if different smoothing parameters are used in estimating  $m(\cdot)$  and  $\nu(\cdot)$ . Furthermore, such a direct method can create a very large bias (§3.1 below). Härdle and Tsybakov (1997) recognized these problems and used a common bandwidth and a common kernel to reduce the bias. While their idea is useful, the approach is still not fully adaptive

to the unknown regression function  $m(\cdot)$ . An alternative regression-adaptive approach is to apply the difference-based estimator (Rice, 1984; Gasser et al., 1986; Müller and Stadtmüller, 1987; also §3.2 below), which uses a high-pass filter to remove the regression function from the data sequence  $\{Y_i\}$ . Hall et al. (1990) demonstrated that the resulting estimator was inefficient even in homoscedastic models with optimal filters.

In this paper, we consider a residual-based estimator of the conditional variance. While the idea is not new (Hall and Carroll, 1989; Neumann, 1994), its implications and implementations are novel. In particular, we show that our estimator is fully regression-adaptive in the sense that without knowing  $m(\cdot)$ , we can estimate the conditional variance function  $\sigma^2(\cdot)$  asymptotically as well as if  $m(\cdot)$  were known. After we have completed this paper, we find that this phenomenon is observed independently by Ruppert et al. (1997) in regression models with independent and identically distributed observations.

One interesting feature of our approach is that we do not need to undersmooth the regression function  $m(\cdot)$  in order to obtain a regression-adaptive estimator for the conditional variance  $\sigma^2(\cdot)$ . In practice, this implies that we can use a data-driven bandwidth selector in estimating  $m(\cdot)$ , then apply the same bandwidth selector with the the squared residuals to estimate  $\sigma^2(\cdot)$ . This is in marked contrast with the previous methods, where new bandwidth (or filter length) selection problems are encountered.

The paper is organized as follows. In §2, we propose and study the residual-based estimator of the conditional variance based on local linear regression. In §3, we compare the performance of our estimator with various procedures in the literature and discuss their mutual relationship. In §4, we present numerical applications with a financial data set and two simulated models. All the technical proofs are given in the Appendix.

## 2 MAIN RESULTS

### 2.1 Estimator

If the regression function  $m(\cdot)$  is given, we can regard the problem of estimating  $\sigma^2(\cdot)$  as a nonparametric regression problem due to the relation

$$E(r|X = x) = \sigma^2(x), \quad \text{where } r = \{Y - m(X)\}^2.$$

Given the observations  $\{(Y_i, X_i), 1 \leq i \leq n\}$  from model (1.1), we write  $r_i = \{Y_i - m(X_i)\}^2$ . Then the local linear estimator of  $\sigma^2(\cdot)$  is  $\hat{\sigma}_b^2(x) = \hat{\alpha}$  (the subscript  $b$  stands for

‘benchmark’), where

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n \{r_i - \alpha - \beta(X_i - x)\}^2 W\left(\frac{X_i - x}{h_1}\right), \quad (2.1)$$

and  $W(\cdot)$  is a density function on  $R$  and  $h_1 > 0$  is a bandwidth (Fan and Gijbels, 1996, p.58). The local linear estimators have several nice properties. They possess high statistical efficiency in an asymptotic minimax sense and are design-adaptive (Fan, 1993). Further, they automatically correct edge effects (Fan and Gijbels, 1992; Ruppert and Wand, 1994; Hastie and Loader, 1995). Therefore,  $\hat{\sigma}_b^2(\cdot)$  provides a benchmark to our problem.

In practice,  $m(\cdot)$  is typically unknown. A natural approach is to substitute  $m(\cdot)$  by a nonparametric regression estimator. We choose the local linear estimator because of its aforementioned optimal properties. Let  $\hat{m}(x) = \hat{a}$  be the local linear estimator that solves the following weighted least-squares problem:

$$(\hat{a}, \hat{b}) = \arg \min_{a, b} \sum_{i=1}^n \{Y_i - a - b(X_i - x)\}^2 K\left(\frac{X_i - x}{h_2}\right), \quad (2.2)$$

where  $K(\cdot)$  is a density function on  $R$  and  $h_2 > 0$  is a bandwidth. Denote the squared residuals by  $\hat{r}_i = \{Y_i - \hat{m}(X_i)\}^2$ . This leads to the residual-based estimator  $\hat{\sigma}^2(x) = \hat{a}$  with kernel  $W$  and bandwidth  $h_1$ , where

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n \{\hat{r}_i - \alpha - \beta(X_i - x)\}^2 W\left(\frac{X_i - x}{h_1}\right). \quad (2.3)$$

Although the above idea appears somewhat ad hoc, it has interesting implications. Specifically, while the bias for  $\hat{m}$  itself is of order  $O(h_2^2)$ , its contribution to  $\hat{\sigma}^2(\cdot)$  is only of  $o(h_2^2)$ . This can be intuitively explained as follows: Observe that

$$\hat{r}_i - r_i = 2\{m(X_i) - \hat{m}(X_i)\}\sigma(X_i)\epsilon_i + \{m(X_i) - \hat{m}(X_i)\}^2,$$

It is intuitively clear that the biases of the residuals are of order  $O\{h_2^4 + (nh_2)^{-1}\}$  and this is the effect of the estimated regression function on the estimated variance. See Theorem 1 and Remark 1 below. This result also paves the way for adapting a fully data-driven bandwidth procedure in our estimation.

## 2.2 ASYMPTOTIC NORMALITY

**THEOREM 1.** Suppose that conditions (C1) — (C5) in the Appendix hold. Then,  $\sqrt{nh_1}\{\hat{\sigma}^2(x) - \sigma^2(x) - \theta_n\}$  is asymptotically normal with mean 0 and variance

$$p^{-1}(x)\sigma^4(x)\lambda^2(x) \int W^2(t)dt,$$

where  $p(\cdot)$  denotes the marginal density function of  $X$ ,  $\lambda^2(x) = E\{(\epsilon^2 - 1)^2 | X = x\}$ ,  $\epsilon = \{Y - m(X)\}/\sigma(X)$ ,  $\sigma_W^2 = \int t^2 W(t) dt$ , and

$$\theta_n = \frac{h_1^2}{2} \sigma_W^2 \ddot{\sigma}^2(x) + o(h_1^2 + h_2^2). \quad (2.4)$$

The above adaptive result is obtained under the assumption that  $m$  is twice differentiable. This is not a minimal condition. The function  $\sigma(\cdot)$  can be estimated with optimal rates under weaker smoothness conditions on  $m(\cdot)$ . See Hall and Carroll (1989) and Müller and Stadtmüller (1993). Note that the above asymptotic normality result complements the asymptotic approximations for conditional mean square error obtained by Ruppert et al. (1997) for regression models with independent observations.

*Remark 1.* The bias and variance expressions given in Theorem 1 are exactly those which arise in the usual nonparametric regression analysis, considering the regression function to be  $\sigma^2(x)$ . In the bias of  $\hat{\sigma}^2(x)$ , the contribution from the error caused by estimating  $m(x)$  is of smaller order than  $h_2^2$ , namely the order of the bias of  $\hat{m}(x)$  itself. This permits us to use the optimal bandwidth to smooth  $\hat{m}$  — no undersmooth of  $\hat{m}$  is needed. Our proof shows further that the second term on the RHS of (2.4) is  $o(n^{-0.6})$  if the bandwidths with optimal rates (i.e.  $h_1 = O(n^{-1/5})$  and  $h_2 = O(n^{-1/5})$ ) are used.

### 2.3 EFFICIENCY

It follows from the local linear regression theory (§6.2.2 of Fan and Gijbels 1996), the benchmark estimator  $\hat{\sigma}_b^2(\cdot)$  derived from (2.1) is asymptotically normal. The leading terms in asymptotic bias and variance are exactly the same as those given in Theorem 1, provided  $h_2$  used in estimator  $\hat{\sigma}^2(\cdot)$  converges to 0 not slower than  $h_1$ . This is a very minor requirement. It is well known that the best  $h_2$  for estimating  $m(\cdot)$  should be of the order  $n^{-1/5}$ . Substituting such an  $h_2$  in (2.4), the optimal  $h_1$  which minimizes the asymptotic mean squared error is also of the order  $n^{-1/5}$ . Therefore, the estimator  $\hat{\sigma}^2(\cdot)$  behaves asymptotically as well as  $\hat{\sigma}_b^2(\cdot)$  and hence is adaptive to the unknown regression function  $m(\cdot)$ . Since the local linear estimator  $\hat{\sigma}_b^2(\cdot)$  is efficient in the sense of Fan (1993), so is  $\hat{\sigma}^2(\cdot)$ .

### 2.4 BANDWIDTH SELECTION

Bandwidth parameter is important to virtually any kernel estimators. The results given in §2.2 permit us to take advantage of existing bandwidth selection methods for

the local linear fit. Let  $\hat{h}(X_1, \dots, X_n; Y_1, \dots, Y_n)$  be a data-driven bandwidth selection rule for the local linear regression based on the data  $(X_1, Y_1), \dots, (X_n, Y_n)$ . This can be in one case the cross-validation bandwidth rule, and in another case the pre-asymptotic substitution method of Fan and Gijbels (1995) or the plug-in approach of Ruppert et al. (1995). The latter two methods have been demonstrated to be less variable and more effective. In all cases,  $\hat{h}(X_1, \dots, X_n; Y_1, \dots, Y_n)$  is a consistent estimate of the asymptotic optimal bandwidth, which is of order  $O(n^{-1/5})$ . Our bandwidth selection rule reads as follows:

1. Use bandwidth  $h_2 = \hat{h}(X_1, \dots, X_n; Y_1, \dots, Y_n)$  in local linear regression (2.2) to obtain the estimate  $\hat{m}(X_i)$  for  $i = 1, \dots, n$ .
2. Compute squared residuals  $\hat{r}_i = \{Y_i - \hat{m}(X_i)\}^2$ ,  $i = 1, \dots, n$ .
3. Apply bandwidth  $h_1 = \hat{h}(X_1, \dots, X_n; \hat{r}_1, \dots, \hat{r}_n)$  in local linear regression (2.3) to obtain  $\hat{\sigma}^2(\cdot)$ .

In the above algorithm, we keep the bandwidth selection method flexible. In our implementation, we use the pre-asymptotic substitution method by Fan and Gijbels (1995), since it has been demonstrated that the resulting estimator possesses fast relative rate of convergence (Huang, 1995, a PhD dissertation).

### 3 OTHER ESTIMATORS

#### 3.1 DIRECT ESTIMATORS

Härdle and Tsybakov (1997) proposed an improved version of the direct estimator  $\hat{\sigma}_d^2(\cdot)$ , as given in (1.2), with local polynomial regression estimators  $\hat{m}(\cdot)$  and  $\hat{\nu}(\cdot)$  using the same kernel function and the same bandwidth, where  $\hat{\nu}(x)$  is an estimate for  $E(Y^2|X = x)$ . They also established the asymptotic normality of the estimator. If the local linear estimators are used with kernel  $W(\cdot)$  and bandwidth  $h_1$ , the leading terms in the asymptotic bias and the asymptotic variance of  $\hat{\sigma}_d^2(x)$  are

$$\text{bias}\{\hat{\sigma}_d^2(x)\} : \frac{h_1^2}{2}\sigma_W^2[\hat{\sigma}^2(x) + 2\{\hat{m}(x)\}^2], \quad \text{var}\{\hat{\sigma}_d^2(x)\} : \frac{1}{nh_1}\sigma^4(x)\lambda^2(x)p^{-1}(x)\int W^2(t)dt.$$

On comparing this with Theorem 1, the direct estimator has the same asymptotic variance as the benchmark  $\hat{\sigma}_b^2(\cdot)$  and the residual-based estimator  $\hat{\sigma}^2(\cdot)$ , but admits one more term

in the bias. This extra term  $h_1^2 \sigma_W^2 \{\dot{m}(x)\}^2$  could lead to an adverse effect on the quality of estimation. For example, even when  $m(\cdot)$  is a linear function with a large slope, this direct method would have a large bias. Thus, the residual-based estimator  $\hat{\sigma}^2(\cdot)$  appears more appealing.

The existence of one more term in the bias of the direct estimator can be understood through the following heuristic arguments. Note that

$$\hat{\sigma}_d^2(x) - \sigma^2(x) = \{\hat{\nu}(x) - \nu(x)\} - 2m(x)\{\hat{m}(x) - m(x)\} - \{\hat{m}(x) - m(x)\}^2. \quad (3.1)$$

The first term on the RHS has the bias

$$\frac{h_1^2}{2} \sigma_W^2 \ddot{\nu}(x) = \frac{h_1^2}{2} \sigma_W^2 \ddot{\sigma}^2(x) + h_1^2 \sigma_W^2 \{\dot{m}(x)\}^2 + h_1^2 \sigma_W^2 m(x) \ddot{m}(x), \quad (3.2)$$

in which the last term on the RHS will cancel the bias of the second term on the RHS of (3.1). Note that the bias from the third term on the RHS of (3.1) is of the order  $h_1^4$ . Therefore, the term involving  $\{\dot{m}(x)\}^2$  stays. This argument also shows that using different kernels or bandwidths in the estimators  $\hat{m}(\cdot)$  and  $\hat{\nu}(\cdot)$  could further increase the bias of  $\hat{\sigma}_d^2(\cdot)$ .

Why can the residual-based estimator  $\hat{\sigma}^2(\cdot)$  give smaller bias? To gain some insight, let us consider the local constant smoother, namely setting  $\beta$  equal 0 in (2.3). Then the resulting estimator is

$$\sum_{i=1}^n \{Y_i - \hat{m}(X_i)\}^2 W\left(\frac{X_i - x}{h_1}\right) \bigg/ \sum_{i=1}^n W\left(\frac{X_i - x}{h_1}\right).$$

This estimator will reduce to the direct estimator  $\hat{\sigma}_d^2(x)$  if all the  $\hat{m}(X_i)$ 's in the above expression are replaced by  $\hat{m}(x)$ . Clearly,  $\{Y_i - \hat{m}(x)\}^2$  is more biased for  $E\{Y - m(X)\}^2$  than  $\{Y_i - \hat{m}(X_i)\}^2$ . This explains why the residual-based estimator inherits less bias from  $\hat{m}(\cdot)$  than the direct estimator.

### 3.2 DIFFERENCE-BASED ESTIMATORS

For a fixed design model

$$Y_i = m(x_i) + \sigma(x_i)\epsilon_i,$$

in which  $x_1 \leq \dots \leq x_n$  are fixed,  $E(\epsilon_i) = 0$  and  $E(\epsilon_i^2) = 1$ , Müller and Stadtmüller (1987) proposed to estimate  $\sigma^2(\cdot)$  through a difference sequence. Their approach can be *briefly*

described as follows. Form an initial local variance estimate

$$\tilde{\sigma}^2(x_i) = \left( \sum_{j=-m}^m w_j Y_{i+j} \right)^2, \quad (3.3)$$

where  $m > 0$  is a prescribed integer, and the difference sequence  $\{w_j\}$  satisfies the conditions

$$\sum_{j=-m}^m w_j = 0, \quad \sum_{j=-m}^m w_j^2 = 1. \quad (3.4)$$

By writing  $\sigma^2(x_i) = \tilde{\sigma}^2(x_i) + \tilde{\epsilon}_i$ , a kernel smoother is applied to obtain the final estimator for  $\sigma^2(\cdot)$  based on the above regression relationship.

Estimators of this type have a long history in the time series context; see, for example, Anderson (1971, p.66). The application in nonparametric homoscedastic regression includes Rice (1984), Gasser et al. (1986), and Hall et al. (1990). It is shown by Hall et al. (1990) that if the optimal difference sequence of  $\{w_i\}$  is employed for a Gaussian model, the efficiency of the estimator is  $4m/(4m+1)$ .

In fact, the residual-based estimator  $\hat{\sigma}^2(\cdot)$  can be regarded as a generalized difference-based estimator, and  $\hat{r}_i$  serves as a crude estimate of  $\sigma^2(X_i)$ . To make such a connection, we express the local linear estimator of  $m(\cdot)$  as

$$\hat{m}(x) = \sum_{i=1}^n w_i(x) Y_i.$$

Then, it can be shown that  $w_1(x) + \dots + w_n(x) = 1$ . Write

$$\hat{r}_i = \{Y_i - \hat{m}(X_i)\}^2 = \left( \sum_{j=1}^n w_{i,j} Y_j \right)^2,$$

where  $w_{i,i} = 1 - w_i(X_i)$  and  $w_{i,j} = -w_j(X_i)$  for  $i \neq j$ . Obviously,  $\{w_{i,j}\}$  is a difference sequence satisfying  $w_{i,1} + \dots + w_{i,n} = 0$ . However, such a sequence of  $\{w_{i,j}\}$  does not exactly fulfill the second condition in (3.4), but

$$\sum_j w_{i,j}^2 = 1 + O_p\{(nh_2)^{-1}\}.$$

The effective length of the sequence is  $2m = 2nh_2$ , which tends to infinity. This also explains why the estimator  $\hat{\sigma}^2(\cdot)$  is efficient in contrast to the aforementioned results of Hall et al. (1990).

Estimation of variance functions with more general weights was discussed by Müller and Stadtmüller (1993). The rates of convergence for this class of estimators were thoroughly investigated. In particular, Müller and Stadtmüller (1993) find that it requires



only very mild smoothness condition on the regression function in order to obtain the optimal rates for the variance estimation.

### 3.3 MAXIMUM LOCALLY LIKELIHOOD ESTIMATORS

If the distribution of  $\epsilon$  is known, the locally maximum likelihood approach could be more efficient. See §4.9 of Fan and Gijbels (1996) and the references therein. For example, if  $\{\epsilon_i\}$  are independent and normal, the log likelihood function can be expressed as

$$-\frac{1}{2} \sum_{i=1}^n L(\sigma^2(X_i), Y_i - m(X_i)),$$

where  $L(\alpha, y) = \alpha^{-1}y^2 + \log \alpha$ . The local maximum likelihood approach with the local constant smoother leads to the direct estimator  $\hat{\sigma}_d^2(\cdot)$  with both  $\hat{m}(\cdot)$  and  $\hat{\nu}(\cdot)$  being the local constant estimators. The approach with the local linear smoother needs to estimate four functions. To make it more tractable, we substitute  $m(\cdot)$  directly by its local linear estimator  $\hat{m}(\cdot)$ , derived from (2.2). Let  $\hat{\alpha}$  and  $\hat{\beta}$  be the minimizer of the residual-based likelihood function:

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n L\{\alpha + \beta(X_i - x), Y_i - \hat{m}(X_i)\} W\left(\frac{X_i - x}{h_1}\right).$$

Then, the maximum local likelihood estimator is defined by  $\hat{\sigma}_{ml}^2(x) = \hat{\alpha}$ . The estimator is also residual-based, and is adaptive to all unknown regression functions in a similar way to what  $\hat{\sigma}^2(\cdot)$  does. The local maximum likelihood estimator share the same leading terms as those for  $\hat{\sigma}^2(x)$  given in §2.1, but is more computationally involved.

## 4 APPLICATIONS AND SIMULATIONS

In this section, we first apply the adaptive estimator  $\hat{\sigma}^2(\cdot)$  derived from (2.3) to an interest-rate data set. The finding from the application includes the validation of an existing structural model. Then, extensive simulations are carried out to confirm the theoretical claim that the adaptive estimator works almost as well as the ideal estimator  $\hat{\sigma}_b^2(\cdot)$  defined in (2.1). We use two simulated models, one with independent observations and one with nonlinear time series.

Throughout this section, the two dashed curves around a solid curve always indicate the two standard deviations above and below the estimated curve. The conditional variance functions are always estimated by the adaptive estimator  $\hat{\sigma}^2(\cdot)$  derived from (2.3)

unless specified otherwise. We always use the Epanechnikov kernel in our calculation. All bandwidths are selected using the pre-asymptotic substitution method by Fan and Gijbels (1995).

*Example 1.* This example concerns the yields of the three month Treasury Bill from the secondary market rates (on Fridays). The secondary market rates are annualized using a 360-day year of bank interest and quoted on a discount basis. The data consist of 1,735 weekly observations, from January 5, 1962 to March 31, 1995, and are presented in Figure 1(a). The data were previously analyzed by various authors, including Andersen and Lund (1997) and Gallant and Tauchen (1997).

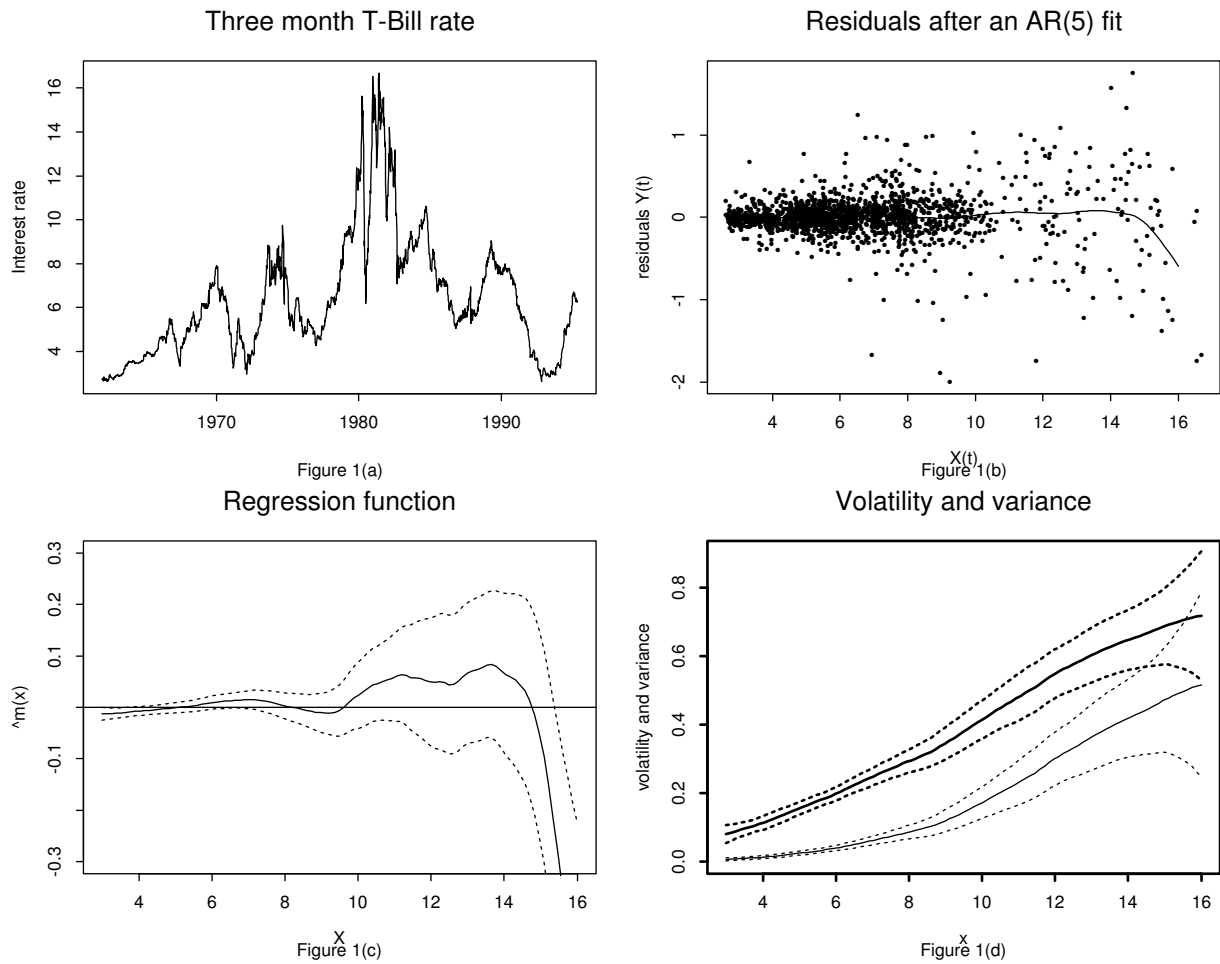


Figure 1: *Three-month Treasury Bill data.* (a) Raw data  $z_t$ . (b) Residuals after an AR(5) fit is plotted against  $X_t \equiv z_{t-1}$ ; solid curve is the regression curve. (c) The regression curve for the data in (b). (d) The estimated volatility curve (thick curve) and the conditional variance function (thin curve).

Let  $z_t$  denote the time series presented in Figure 1(a). We first fitted an AR model with order selected by the Akaike information criterion. This yields the following AR(5)

model:

$$z_t = 1.0733z_{t-1} - 0.0423z_{t-2} + 0.0165z_{t-3} + 0.0228z_{t-4} - 0.0773z_{t-5} + Y_t.$$

The selection of an AR(5) model coincides with that used by Andersen and Lund in a technical report. The ‘residuals’  $Y_t$  are plotted against  $X_t \equiv z_{t-1}$  in Figure 1(b). Figure 1(c) depicts the estimator of the mean regression function  $m(x) \equiv E\{Y_t | z_{t-1} = x\}$ . The nonlinearity with a slightly increasing trend (up to  $z_{t-1} = 14$ ) can be noted. The bandwidth selected by our software is 1.9535. The residual-based estimator for the conditional variance of  $Y_t$  given  $z_{t-1} = x$  is denoted as  $\hat{\sigma}^2(x)$  with the automatically selected bandwidth 3.1461. The estimated volatility function  $\hat{\sigma}(x)$  is presented in Figure 1(d). The overall fitted model is

$$z_t = \hat{m}(z_{t-1}) + 1.0733z_{t-1} - 0.0423z_{t-2} + 0.0165z_{t-3} + 0.0228z_{t-4} - 0.0773z_{t-5} + \hat{\sigma}(z_{t-1})\epsilon_t,$$

in which  $E(\epsilon_t | z_{t-1}) = 0$ , and  $\text{var}(\epsilon_t | z_{t-1}) = 1$ . Note that the correlation coefficient between the logarithm of  $z_{t-1}$  and logarithm of  $\hat{\sigma}(z_{t-1})$  is 0.999. This lends a strong support to the structural volatility model

$$\sigma(z_{t-1}) = \alpha z_{t-1}^\beta,$$

which was considered by Andersen and Lund in a technical report. Applying the least-square fit to the log-transformed data, we found that  $\alpha = 0.0169$  and  $\beta = 1.380$ .

A commonly-used model for asset pricing admits the following form: The value  $S_t$  of an underlying asset at time  $t$  satisfies

$$dS_t = \mu(S_t)dt + \sigma(S_t)dW_t, \quad (4.1)$$

where  $\mu$  is the (instantaneous) expected rate of return,  $\sigma$  is the price volatility, and  $W_t$  is the standard Wiener process. This nonparametric model was recently used to model term structure dynamics by for example Aït-Sahalia (1996) and Stanton (1998). It includes the famous interest rate models of Cox, Ingersoll and Ross (1985), Chan, Karolyi, Longstaff and Sanders (1992), among others. We now briefly connect this continuous model with our nonparametric regression model. Suppose that the data are sampled at time  $i\Delta$  for  $i = 1, \dots, T-1$ . Set  $Y_i = (S_{(i+1)\Delta} - S_{i\Delta})/\Delta$  and  $X_i = S_{i\Delta}$ . Model (4.1) can be understood as

$$Y_i \approx \mu(X_i) + \sigma(X_i)\varepsilon_i/\sqrt{\Delta}, \quad (4.2)$$

where  $\{\varepsilon_i\}$  are Gaussian white noise. Therefore, our method can be directly used to estimate functions  $\mu(\cdot)$  and  $\sigma(\cdot)$ . For the short interest rate data set, our  $Y_t$  is basically the same as the difference  $z_t - z_{t-1}$ . Therefore, functions in Figures 1(c) and 1(d) are respectively a scaled version of the estimated expected rate of return and price volatility in the continuous model (4.1). In fact, similar estimates to Figures 1 (c) and (d) were independently obtained by Stanton (1998). Our method differs from that of Stanton (1988) in the following three important aspects: Squared residuals instead of the squared responses  $Y_i^2$  are used to estimate the volatility; local linear approach instead of kernel method is used for nonparametric regression; more sophisticated bandwidth selection techniques are implemented. The first two aspects can reduce considerably biases in the estimate and the last aspect enables one to conduct correct amount of smoothing.

*Example 2.* We simulated 400 random samples of size  $n = 200$  from the model

$$Y_i = a\{X_i + 2 \exp(-16X_i^2)\} + \sigma(X_i)\epsilon_i, \quad \text{with } \sigma(x) = 0.4 \exp(-2x^2) + 0.2,$$

where  $\{X_i\}$  and  $\{\epsilon_i\}$  are two independent sequences of independent random variables, and  $X_i \sim U[-2, 2]$  and  $\epsilon_i \sim N(0, 1)$ . Four different values of  $a$ , namely  $a = 0.5, 1, 2, 4$ , are used in the simulation. For each simulated sample, the performance of the estimator is evaluated by the Mean Absolute Deviation Error (MADE):

$$\text{MADE} = n_{\text{grid}}^{-1} \sum_{i=1}^{n_{\text{grid}}} |\hat{\sigma}(x_j) - \sigma(x_j)|,$$

where  $\{x_j, j = 1, \dots, n_{\text{grid}}\}$  are the grid points on  $[-1.8, 1.8]$  with  $n_{\text{grid}} = 101$ . The results are summarized in Figure 2. Figure 2(a) compares the adaptive variance estimator with the ideal variance estimator  $\hat{\sigma}_b^2(\cdot)$  which does not vary with different values of  $a$ . Presented there are the boxplots of MADEs based on 400 simulations. The first four boxplots are the MADEs of the adaptive estimator  $\hat{\sigma}^2(\cdot)$  for  $a = 0.5, 1, 2, 4$  in order, and the last one is that of the ideal estimator  $\hat{\sigma}_b^2(\cdot)$ . As anticipated, the adaptive estimator performs almost as well as the ideal one.

To get further insights, we consider the specific case  $a = 1$ . The scenario is similar for other cases. Figure 2(b) plots the MADE based on the adaptive estimator versus the MADE based on the ideal estimator, using the same sample data. Clearly, there is about equal chance that one estimator beats the other. The marginal densities of MADE of the adaptive estimator (thick curve) and of the ideal estimator (thin curve) are also depicted in Figure 2(b). This shows again that the performance of the two estimators

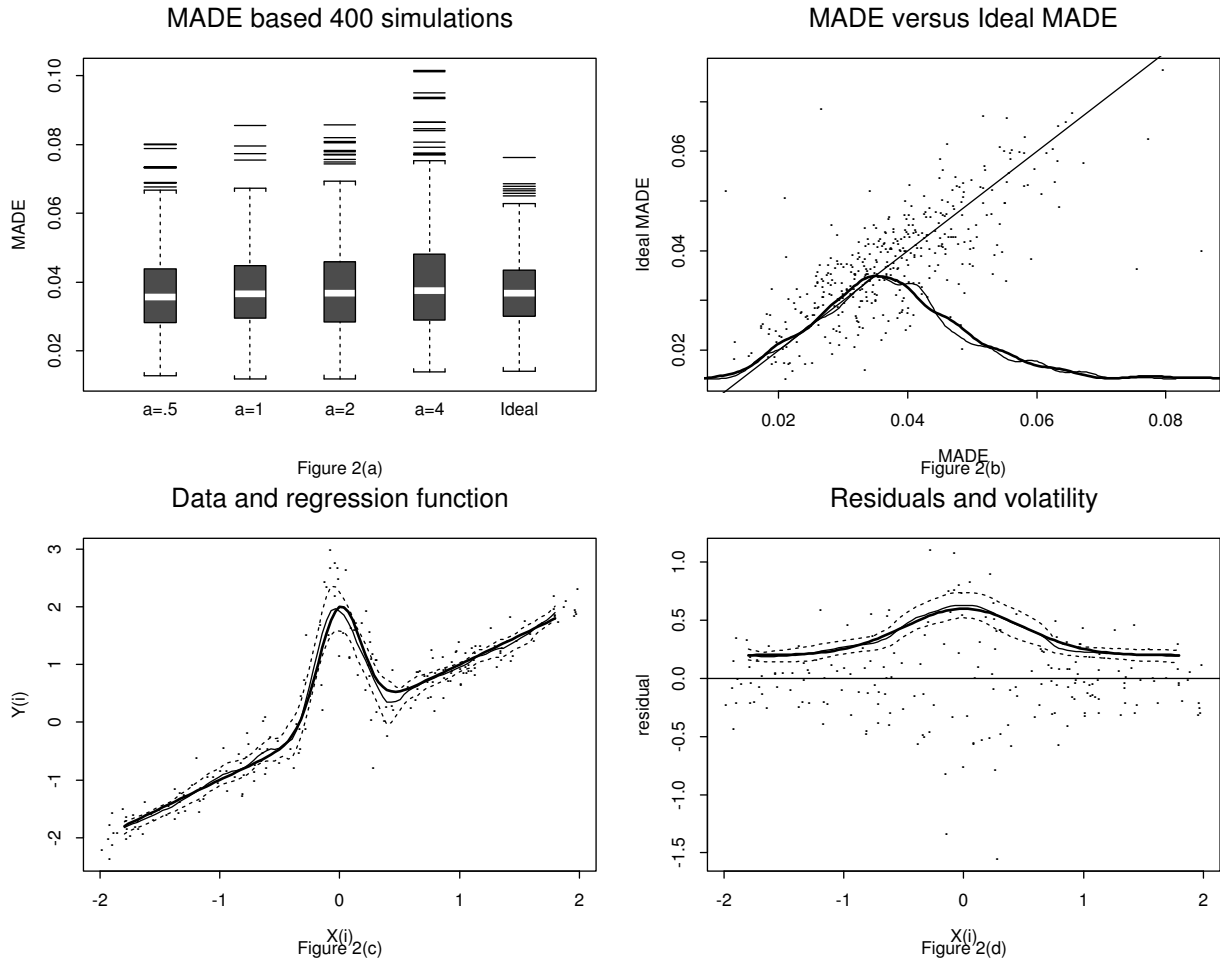


Figure 2: *Simulation results for Example 2. (a) Boxplots of the MADEs for the adaptive estimator with  $a = 0.5, 1, 2, 4$ , and for the ideal estimator (from left to right). (b) The scatter plot of the MADE of  $\hat{\sigma}^2(\cdot)$  versus the MADE of  $\hat{\sigma}_b^2(\cdot)$ ; the straight line marks the position where the two MADEs are equal. Thick curve — the estimated density function of the MADE of  $\hat{\sigma}^2(\cdot)$ ; thin curve — the estimated density function of the MADE of  $\hat{\sigma}_b^2(\cdot)$ . (c) A representative sample, the corresponding estimated regression curve (thin curve), and the true regression curve (thick curve). (d) The sample residuals from (c), the estimated volatility (thin curve), and the true volatility (thick curve).*

is comparable. Figure 2(c) presents a typical simulated sample with its corresponding estimated regression function. The typical sample was selected in such a way that the corresponding MADE is equal to its median among the 400 simulations. The sample residuals and the estimated conditional standard deviations are plotted in Figure 2(d). The bandwidths are automatically selected by the procedure outlined in §2.4 and are 0.1867 for the mean regression and 0.4841 for the conditional variance function respectively.

*Example 3.* Consider the following nonlinear time series model

$$X_{t+1} = 0.235X_t(16 - X_t) + e_t,$$

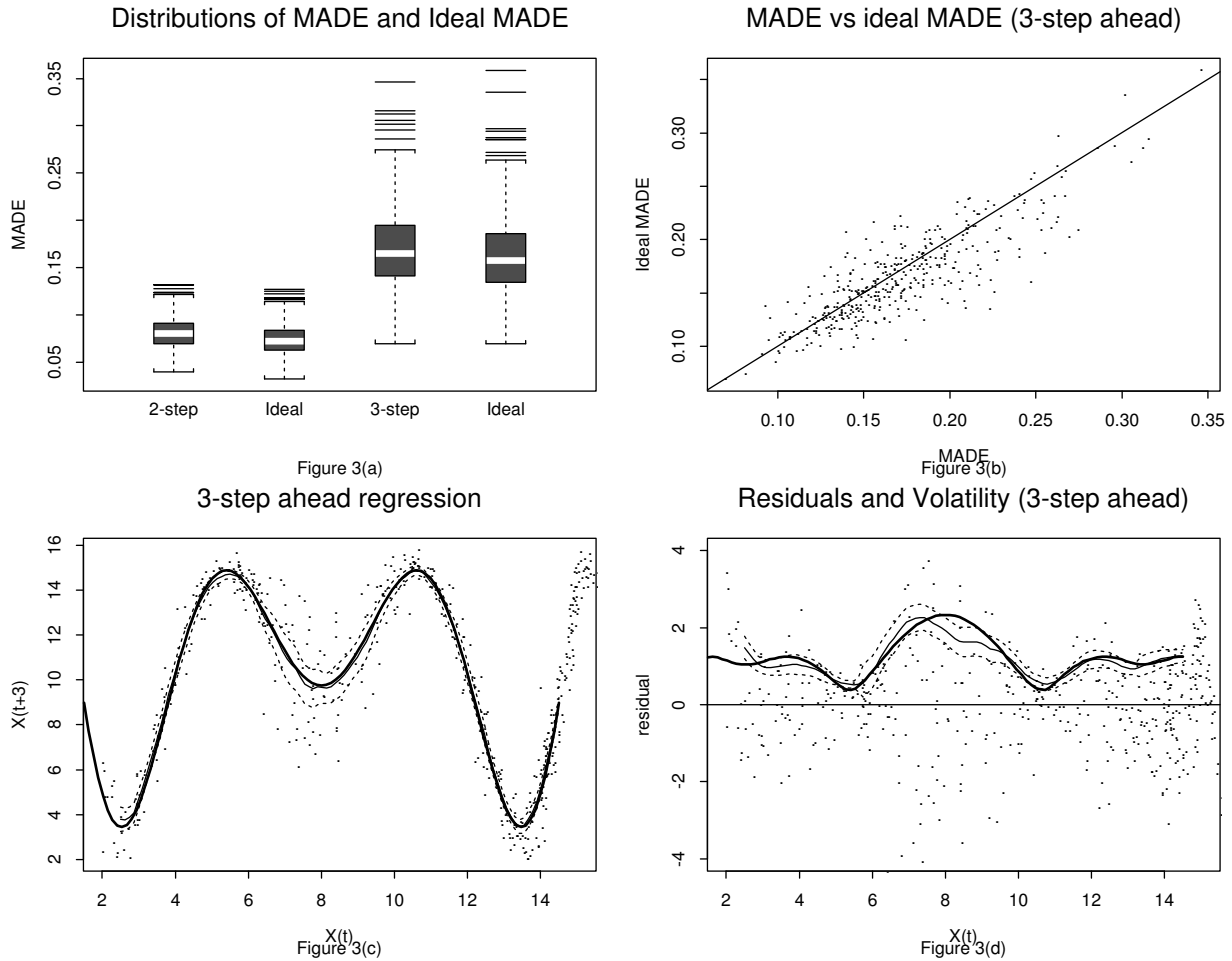


Figure 3: *Simulation results for Example 3. (a) Boxplots of the MADEs for the adaptive estimator and for the ideal estimator. (b) The scatter plot of the MADE versus the ideal MADE for 3-step prediction. The straight line marks the position where the two MADEs are equal. (c) A representative sample and its estimated 3-step ahead regression curve. (d) The sample residuals of (c), the true volatility function (thick curve), and estimated volatility function.*

where  $e_1, e_2, \dots$ , are independent with the common distribution  $N(0, 0.3^2)$ . The skeleton of this model exhibits chaos and has been used by Yao and Tong (1994) to illustrate the influence of the initial values on nonlinear prediction.

For this nonlinear time series, we consider the two-step and three-step ahead prediction by taking respectively  $Y_t = X_{t+2}$  and  $Y_t = X_{t+3}$ . Note that the conditional variance functions concerned are not constant. On the other hand, the conditional variance of the one-step prediction is a constant, and is therefore not presented here.

Figure 3(a) compares the ideal estimator with the adaptive estimator based on 400 simulations with  $n = 500$ . As we can see, the adaptive estimator works almost as well

as the ideal estimator. Figure 3(b) gives the scatter plot of MADE for the adaptive estimator and the ideal estimator for three-step ahead prediction. A typical simulated data set and the corresponding estimated curves are presented in Figures 3 (b) and (c). The criterion used to choose a typical sample is again the one for which the MADE is equal to its median among the 400 simulations. Figure 3(c) presents the estimated regression function for 3-step ahead prediction, where bandwidth 0.5577 was selected by our procedure. The estimated volatility function is presented in Figure 3(d) with data-driven bandwidth 0.8165. Similar results to Figures 3(b)–(d) were obtained for two-step prediction and are omitted for brevity.

## ACKNOWLEDGMENT

We would like to thank Professor Ron Gallant for kindly providing us the data analyzed in Example 1. We are also grateful to the reviewers for their insightful comments. JF's research was partially supported by an NSF grant and an NSA grant. QY's research was partially supported by an EPSRC grant.

## REFERENCES

- Aït-Sahalia, Y. (1996). Nonparametric pricing of interest rate derivative securities. *Econometrica*, **64**, 527-560.
- Andersen, T.G. and Lund, J. (1997). Estimating continuous time stochastic volatility models of the short term interest rate. *Journal of Econometrics*, **77**, 343-77.
- Anderson, T.W. (1971). *The Statistical Analysis of Time Series*. Wiley, New York.
- Box, G. (1988). Signal-to-noise ratios, performance criteria, and transformations. *Technometrics*, **30**, 1-17.
- Carroll, R. and Ruppert, D. (1988). *Transformations and Weighting in Regression*. Chapman & Hall, London.
- Chan, K.C., Karolyi, A.G., Longstaff, F.A. and Sanders, A.B. (1992). An empirical comparison of alternative models of the short-term interest rate. *J. Finance*, **47**, 1209-1227.
- Cox, J.C. Ingersoll, J.E. and Ross, S. A. (1985). A theory of the term structure of interest rates. *Econometrica*, **53**, 385-467.
- Engle, R.F. (1982). Autoregressive conditional heteroscedasticity with estimates of variance of U.K. inflation. *Econometrica*, **50**, 987-1008.

- Fan, J. (1993). Local linear regression smoothers and their minimax efficiency. *Ann. Statist.*, **21**, 196-216.
- Fan, J. and Gijbels, I. (1992). Variable bandwidth and local linear regression smoothers. *Ann. Statist.*, **20**, 2008-36.
- Fan, J. and Gijbels, I. (1995). Data-driven bandwidth selection in local polynomial fitting: variable bandwidth and spatial adaptation. *J. Royal Statist. Soc. B*, **57**, 371-94.
- Fan, J. and Gijbels, I. (1996). *Local Polynomial Modelling and Its Applications*. Chapman & Hall, London.
- Gasser, T., Sroka, L. and Jennen-Steinmetz, C. (1986). Residual variance and residual pattern in nonlinear regression. *Biometrika*, **73**, 625-33.
- Gallant, A. R. and Tauchen, G. (1997). Estimation of continuous time models for stock returns and interests rates. *Macroeconomic Dynamics*, **1**, 343-378.
- Hall, P. and Carroll, R.J. (1989). Variance function estimation in regression: the effect of estimation of the mean. *J. Roy. Statist. Soc. B*, **51**, 3-14.
- Hall, P., Kay, J. and Titterton, D.M. (1990). Asymptotically optimal difference-based estimation of variance in nonparametric regression. *Biometrika*, **77**, 521-28.
- Härdle, W. and Tsybakov, A. (1997). Local polynomial estimators of the volatility function in nonparametric autoregression. *Journal of Econometrics*, **81**, 223-242.
- Hastie, T.J. and Loader, C. (1993). Local regression: automatic kernel carpentry (with discussion). *Statist. Sci.*, **8**, 120-43.
- Müller, H.G. and Stadtmüller U. (1987). Estimation of heteroscedasticity in regression analysis. *Ann. Statist.*, **15**, 610-25.
- Müller, H.G. and Stadtmüller U. (1993). On variance function estimation with quadratic forms. *J. Statist. Plann. Inf.* **35**, 213-31.
- Neumann, M.H. (1994). Fully data-driven nonparametric variance estimators. *Statistics*, **25**, 189-212.
- Peligrad, M. (1986). Recent advances in the central limit theorem and its weak invariance principle for mixing sequences of random variables. *Dependence in Probability and Statistics*, Ed. E. Eberlein and M.S. Taquq. Birkhäuser, Boston, 193-223.
- Rice, J. (1984). Bandwidth choice for nonparametric kernel regression. *Ann. Statist.*, **12**, 1215-30.
- Ruppert, D., Sheather, S.J. and Wand, M.P. (1995). An effective bandwidth selector for local least squares regression. *J. Amer. Statist. Assoc.*, **90**, 1257-70.
- Ruppert, D. and Wand, M.P. (1994). Multivariate weighted least squares regression. *Ann. Statist.*, **22**, 1346-70.
- Ruppert, D., Wand, M.P., Holst, U. and Hössjer, O. (1997). Local Polynomial Variance Function Estimation. *Technometrics*, **39**, 262-73.



- Stanton, R. (1998). A nonparametric model of term structure dynamics and the market price of interest rate risk. *J. Finance*, to appear.
- Yao, Q. and Tong, H. (1994). Quantifying the influence of initial values on nonlinear prediction. *J. Roy. Statist. Soc. B*, **56**, 701-25.
- Yoshihara, K. (1976). Limiting behaviour of U-statistics for stationary absolutely regular processes. *Z. Wahr. v. Gebiete*, **35**, 237-52.

## APPENDIX

### *Regularity conditions*

We use the same notation as in §2. We always use  $c$  to denote a generic constant which may be different at different places. We introduce the following regularity conditions.

- (C1) For a given point  $x$ ,  $p(x) > 0$ ,  $\sigma^2(x) > 0$  and the functions  $E\{Y^k|X = z\}$  is continuous at  $x$  for  $k = 3, 4$ . Further,  $\ddot{m}(z) \equiv d^2 m(z)/dz^2$  and  $\ddot{\sigma}^2(z) \equiv d^2\{\sigma^2(z)\}/dz^2$  are uniformly continuous on an open set containing the point  $x$ .
- (C2)  $E\{Y^{4(1+\delta)}\} < \infty$ , where  $\delta \in [0, 1)$  is a constant.
- (C3) The kernel functions  $W$  and  $K$  are symmetric density functions each with a bounded support in  $R$ . Further,  $|W(x_1) - W(x_2)| \leq c|x_1 - x_2|$ ,  $|K(x_1) - K(x_2)| \leq c|x_1 - x_2|$  and also  $|p(x_1) - p(x_2)| \leq c|x_1 - x_2|$  for  $x_1, x_2 \in R$ .
- (C4) The strictly stationary process  $\{(X_i, Y_i)\}$  is absolutely regular, i.e.

$$\beta(j) \equiv \sup_{i \geq 1} E \left\{ \sup_{A \in \mathcal{F}_{i+1}^\infty} |pr(A|\mathcal{F}_1^i) - pr(A)| \right\} \rightarrow 0, \quad \text{as } j \rightarrow \infty,$$

where  $\mathcal{F}_i^j$  is the  $\sigma$ -field generated by  $\{(X_k, Y_k) : k = i, \dots, j\}$ , ( $j \geq i$ ). Further for the same  $\delta$  as in (C2),

$$\sum_{j=1}^{\infty} j^2 \beta^{\frac{\delta}{1+\delta}}(j) < \infty.$$

(We use the convention  $0^0 = 0$ .)

- (C5) As  $n \rightarrow \infty$ ,  $h_i \rightarrow 0$ , and  $\liminf_{n \rightarrow \infty} nh_i^4 > 0$ , for  $i = 1, 2$ .

We impose the boundedness on the supports of  $K(\cdot)$  and  $W(\cdot)$  for brevity of proofs; it may be removed at the cost of lengthier proofs. In particular, the Gaussian kernel is allowed. The assumption of the convergence rate of  $\beta(j)$  is also for technical convenience. The assumption on the convergence rates of  $h_1$  and  $h_2$  is not the weakest possible.

*Remark 2.* When  $\{(X_t, Y_t)\}$  are independent, (C4) holds with  $\delta = 0$  and condition (C2) reduces to  $E(Y^4) < \infty$ . On the other hand, if (C4) holds with  $\delta = 0$ , there are at most finitely many non-zero  $\beta(j)$ 's. This means that there exists an integer  $0 < j_0 < \infty$  for which  $(X_i, Y_i)$  is independent of  $\{(X_j, Y_j), j \geq i + j_0\}$ , for all  $i \geq 1$ .

### Proofs

In the sequel,  $\hat{m}(\cdot)$  denotes the local linear estimator derived from (2.2). We always assume that conditions (C1) — (C5) hold. We call that  $B_n(x) = B(x) + o_p(b_n)$  (or  $O_p(b_n)$ ) uniformly for  $x \in G$  if

$$\sup_{x \in G} |B_n(x) - B(x)| = o_p(b_n) \text{ (or } O_p(b_n)\text{)}.$$

We only present the proof for the cases with  $\delta > 0$ . The case with  $\delta = 0$  can be dealt in a more direct and simpler way. (See Remark 2.)

The proof is based on the following lemma which follows from Lemma 2 of Yao and Tong (1996, a technical report) directly.

LEMMA 1. Let  $G \subset \{p(x) > 0\}$  be a compact subset. As  $n \rightarrow \infty$ , uniformly for  $x \in G$ ,

$$\hat{\sigma}^2(x) - \sigma^2(x) = \frac{1}{nh_1 p(x)} \sum_{i=1}^n W\left(\frac{X_i - x}{h_1}\right) \{\hat{r}_i - \sigma^2(x) - \dot{\sigma}^2(x)(X_i - x)\} + O_p\{R_{n,1}(x)\}, \quad (\text{A.3})$$

$$\hat{m}(x) - m(x) = \frac{1}{nh_2 p(x)} \sum_{i=1}^n \sigma(X_i) \epsilon_i K\left(\frac{X_i - x}{h_2}\right) + \frac{h_2^2 \sigma_K^2}{2} \ddot{m}(x) + O_p\{R_{n,2}(x)\} \quad (\text{A.4})$$

where

$$\begin{aligned} R_{n,1}(x) &= \frac{1}{np(x)} \left[ \left| \sum_{i=1}^n W\left(\frac{X_i - x}{h_1}\right) \{\hat{r}_i - \sigma^2(x) - \dot{\sigma}^2(x)(X_i - x)\} \right| \right. \\ &\quad \left. + \left| \sum_{i=1}^n \frac{X_i - x}{h_1} W\left(\frac{X_i - x}{h_1}\right) \{\hat{r}_i - \sigma^2(x) - \dot{\sigma}^2(x)(X_i - x)\} \right| \right], \\ R_{n,2}(x) &= \frac{1}{np(x)} \left\{ \left| \sum_{i=1}^n K\left(\frac{X_i - x}{h_2}\right) \sigma(X_i) \epsilon_i \right| + \left| \sum_{i=1}^n \frac{X_i - x}{h_2} K\left(\frac{X_i - x}{h_2}\right) \sigma(X_i) \epsilon_i \right| \right\} + O(h_2^3). \end{aligned}$$

**Proof of Theorem 1.** Note that

$$\begin{aligned}\hat{r}_i &= \{Y_i - \hat{m}(X_i)\}^2 = \{\sigma(X_i)\epsilon_i + m(X_i) - \hat{m}(X_i)\}^2 \\ &= \sigma^2(X_i)\epsilon_i^2 + 2\sigma(X_i)\epsilon_i\{m(X_i) - \hat{m}(X_i)\} + \{m(X_i) - \hat{m}(X_i)\}^2.\end{aligned}$$

It follows from (A.3) that

$$\hat{\sigma}^2(x) - \sigma^2(x) = I_1 + I_2 - I_3 + I_4 + O_p(h_1)(|I_1 + I_2 - I_3 + I_4| + |I'_1 + I'_2 - I'_3 + I'_4|),$$

where

$$\begin{aligned}I_1 &= \frac{1}{nh_1p(x)} \sum_{i=1}^n W\left(\frac{X_i - x}{h_1}\right) \{\sigma^2(X_i) - \sigma^2(x) - \hat{\sigma}^2(x)(X_i - x)\}, \\ I_2 &= \frac{1}{nh_1p(x)} \sum_{i=1}^n W\left(\frac{X_i - x}{h_1}\right) \sigma^2(X_i)(\epsilon_i^2 - 1), \\ I_3 &= \frac{2}{nh_1p(x)} \sum_{i=1}^n W\left(\frac{X_i - x}{h_1}\right) \sigma(X_i)\epsilon_i\{\hat{m}(X_i) - m(X_i)\}, \\ I_4 &= \frac{1}{nh_1p(x)} \sum_{i=1}^n W\left(\frac{X_i - x}{h_1}\right) \{\hat{m}(X_i) - m(X_i)\}^2,\end{aligned}\tag{A.5}$$

and  $I'_j$  is defined in the same way as  $I_j$  with one more factor  $h_1^{-1}(X_i - x)$  in the  $i$ -th summand ( $1 \leq j \leq 4$ ). It is easy to see that the theorem follows from statements (a) — (d) below directly.

- (a)  $I_1 = \frac{1}{2}h_1^2\ddot{\sigma}^2(x)\sigma_W^2 + o_p(h_1^2)$ , and  $I'_1 = o_p(h_1^2)$ .
- (b)  $\sqrt{nh_1}I_2 \xrightarrow{d} N(0, \sigma^4(x)\lambda^2(x) \int W^2(t)dt/p(x))$ , and  $\sqrt{nh_1}I'_2 \xrightarrow{d} N(0, \sigma^4(x)\lambda^2(x) \int t^2W^2(t)dt/p(x))$ .
- (c)  $I_3 = o_p(h_1^2 + h_2^2)$ , and  $I'_3 = o_p(h_1^2 + h_2^2)$ .
- (d)  $I_4 = o_p(h_1^2 + h_2^2)$ , and  $I'_4 = o_p(h_1^2 + h_2^2)$ .

In the sequel, we establish the statements on  $I_j$  in (a) — (d) only. The cases with  $I'_j$  can be proved in the same manner.

It is easy to see that (a) follows from a Taylor's expansion, and a direct application of the ergodic theorem. Conditions (C2) and (C3) imply that  $E\{W\left(\frac{X_i - x}{h_1}\right) \sigma^2(X_i)(\epsilon_i^2 - 1)\}^{2+\delta/2} < \infty$ . Note that the condition of absolutely regular implies  $\alpha$ -mixing with  $\alpha(j) \leq \beta(j)$ . By (C4) and Theorem 1.7 of Peligrad (1986),  $I_2$  is asymptotically normal with mean 0 and variance  $\sigma_*^2/nh_1$ , where

$$\sigma_*^2 = \frac{1}{h_1} E \left\{ W\left(\frac{X - x}{h_1}\right) \frac{\sigma^2(X)}{p(X)} (\epsilon^2 - 1) \right\}^2$$

$$+ \frac{1}{h_1} \sum_{i=2}^n E \left\{ W \left( \frac{X_1 - x}{h_1} \right) \frac{\sigma^2(X_1)}{p(X_1)} (\epsilon_1^2 - 1) W \left( \frac{X_i - x}{h_1} \right) \frac{\sigma^2(X_i)}{p(X_i)} (\epsilon_i^2 - 1) \right\}. \quad (\text{A.6})$$

It is easy to see that the first term in the above expression converges to  $\sigma^4(x)\lambda^2(x) \int W^2(t)dt/p(x)$ .

Note that for any  $i \geq 2$ ,

$$E \left\{ W \left( \frac{X_1 - x}{h_1} \right) \frac{\sigma^2(X_1)}{p(X_1)} (\epsilon_1^2 - 1) W \left( \frac{X_i - x}{h_1} \right) \frac{\sigma^2(X_i)}{p(X_i)} (\epsilon_i^2 - 1) \right\}^{1+\delta} = O(h_1^2),$$

$$E \left\{ W \left( \frac{X - x}{h_1} \right) \frac{\sigma^2(X)}{p(X)} (\epsilon^2 - 1) \right\} = 0, \quad E \left| W \left( \frac{X - x}{h_1} \right) \frac{\sigma^2(X)}{p(X)} (\epsilon^2 - 1) \right|^{1+\delta} = O(h_1).$$

It follows from (C4) and Lemma 1 Yoshihara (1976) that the absolute value of the second term on the RHS of (A.6) is bounded above by  $ch_1^{(1-\delta)/(1+\delta)} \{\beta^{\frac{\delta}{1+\delta}}(1) + \dots + \beta^{\frac{\delta}{1+\delta}}(n-1)\} = o(1)$ . Hence (b) holds.

Note that  $W(\cdot)$  has a bounded support contained in the interval  $[-s_w, s_w]$ , say. Therefore in the summation on the RHS of (A.5), only those terms with  $X_i \in [x - h_2s_w, x + h_2s_w]$  might not be 0. It follows from (A.4) that we may write  $I_3 = I_{31} + I_{32} + I_{33}$ , where

$$\begin{aligned} I_{31} &= \frac{1}{n^2 h_1 h_2 p(x)} \sum_{i,j=1}^n K \left( \frac{X_i - X_j}{h_2} \right) \sigma(X_i) \sigma(X_j) \epsilon_i \epsilon_j \left\{ p^{-1}(X_i) W \left( \frac{X_i - x}{h_1} \right) \right. \\ &\quad \left. + p^{-1}(X_j) W \left( \frac{X_j - x}{h_1} \right) \right\} \\ &\equiv \frac{1}{n^2 h_1 h_2 p(x)} \sum_{i,j=1}^n \varphi_{ij} = \frac{2}{n^2 h_1 h_2 p(x)} \sum_{1 \leq i < j \leq n} \varphi_{ij} + O_p \left( \frac{1}{n h_2} \right), \end{aligned} \quad (\text{A.7})$$

$$I_{32} = \frac{h_2^2 \sigma_K^2}{n h_1 p(x)} \sum_{i=1}^n W \left( \frac{X_i - x}{h_1} \right) \sigma(X_i) \epsilon_i \ddot{m}(X_i) = o_p(h_2^2), \quad (\text{A.8})$$

$$\begin{aligned} |I_{33}| &\leq \frac{O_p(1)}{n^2 h_1} \left| \sum_{i,j=1}^n W \left( \frac{X_i - x}{h_1} \right) K \left( \frac{X_i - X_j}{h_2} \right) \sigma(X_i) \sigma(X_j) |\epsilon_i| |\epsilon_j| / p(X_i) \right| \\ &\quad + \frac{O_p(1)}{n^2 h_1} \left| \sum_{i,j=1}^n \frac{X_j - X_i}{h_2} W \left( \frac{X_i - x}{h_1} \right) K \left( \frac{X_i - X_j}{h_2} \right) \sigma(X_i) \sigma(X_j) |\epsilon_i| |\epsilon_j| / p(X_i) \right| + o_p(h_2^2). \end{aligned}$$

It follows from Lemma A(ii) of Hjellvik, Yao and Tjøstheim (1996, a technical report) that for any  $\varepsilon_0 > 0$  and  $\varepsilon > 0$ ,

$$pr\{n^{-1}(h_1 h_2)^{-(\frac{1}{1+\delta}-\varepsilon_0)/2} \left| \sum_{i < j} \varphi_{ij} \right| > \varepsilon\} \leq \frac{c(h_1 h_2)^{\varepsilon_0}}{n^2} E\{(h_1 h_2)^{-\frac{1}{2(1+\delta)}} \sum_{i < j} \varphi_{ij}\}^2 = o((h_1 h_2)^{\varepsilon_0}).$$

Therefore, the first term on the RHS of (A.7) is  $o_p\{n^{-1}(h_1 h_2)^{-(\frac{1+2\delta}{1+\delta}+\varepsilon_0)/2}\}$ . Thus

$$I_{31} = o_p(n^{-1}(h_1 h_2)^{-(\frac{1+2\delta}{1+\delta}+\varepsilon_0)/2}) + O_p(n^{-1}h_2^{-1}).$$

Condition (C5) implies that both terms on the RHS of the above expression is of the order  $o_p(h_1^2 + h_2^2)$  if we choose  $\varepsilon_0 < (1 + \delta)^{-1}$ . Performing Hoeffding's projection decomposition of  $U$ -statistics, we can prove  $I_{33} = o_p(h_1^2 + h_2^2)$  in the same way.

The proof of (d) is similar to that of (c), therefore is omitted here.