

NONPARAMETRIC ESTIMATION AND SYMMETRY TESTS FOR CONDITIONAL DENSITY FUNCTIONS

ROB J. HYNDMAN^a and QIWEI YAO^b

^a*Department of Econometrics and Business Statistics, Monash University, Clayton VIC 3800, Australia;* ^b*Department of Statistics, London School of Economics, Houghton Street, London WC2A 2AE, UK*

(Received 26 October 1999; In final form 6 June 2001)

We suggest two improved methods for conditional density estimation. The first is based on locally fitting a log-linear model, and is in the spirit of recent work on locally parametric techniques in density estimation. The second method is a constrained local polynomial estimator. Both methods always produce non-negative estimators. We propose an algorithm suitable for selecting the two bandwidths for either estimator. We also develop a new bootstrap test for the symmetry of conditional density functions. The proposed methods are illustrated by both simulation and application to a real data set.

Keywords: Bandwidth selection; Bootstrap; Conditioning; Density estimation; Kernel smoothing; Symmetry tests

1 INTRODUCTION

We have two goals in this paper. First, we propose two new methods for estimating the conditional density function of Y_t given X_t based on observations from a strictly stationary process $\{(X_t, Y_t)\}$. Second, we propose a new bootstrap method for testing the symmetry of conditional density functions.

Our new conditional density estimation methods improve on the local polynomial estimators proposed by Fan, Yao and Tong (1996) by restricting the estimator to be non-negative. The “double kernel” smoothing approach is similar to that adapted by Yu and Jones (1998) to estimate conditional quantiles.

Our first estimation method is locally parametric; it produces estimators of arbitrarily high order and is always non-negative. In spirit, this approach is related to recently-introduced local parametric methods for density estimation; see, for example, Copas (1995), Simonoff (1996, Section 3.4), Hjort and Jones (1997), Loader (1996) and Hall, Wolff and Yao (1999). Our second method is a constrained version of the estimator studied by Fan, Yao and Tong (1996). The simple constraint makes the estimator always non-negative while retaining the nice asymptotic properties of the local polynomial estimators.

We consider the mean square error properties of our estimators and show that the asymptotic optimal bandwidth in the x -direction is greater than that in ordinary kernel regression

estimation in order to compensate for the data sparseness due to the smoothing in y -direction. Similarly, the optimal bandwidth in the y -direction is greater than that for unconditional density estimation to compensate for the smoothing in the x -direction. Based on the mean-square error properties, we propose a practical bandwidth selection algorithm for the new estimators.

The symmetry of conditional density functions is of interest in modelling time series data in business and finance (Brännäs and De Gooijer, 1992) and in constructing predictive regions for nonlinear time series (Hyndman, 1995; Polonik and Yao, 2000; De Gooijer and Gannoun, 2000). As far as we know, the symmetry of *conditional* density functions has never been addressed in the literature before. However, various statistical methods have been proposed for testing the symmetry of *unconditional* density functions, which include, among others, Butler (1969), Hollander (1971), Rothman and Woodroffe (1972), Srinivasan and Godio (1974), Doksum *et al.* (1977), Hill and Rao (1977), Lockhart and McLaren (1985), Csörgö and Heathcote (1987), Zhu (1998) and Diks and Tong (1999).

The paper is organized as follows: we propose the two new estimators for conditional densities in Section 2. The asymptotic normality of the estimators is presented under some mixing conditions. Section 3 addresses the issue of bandwidth selection. The bootstrap tests for the symmetry are discussed in Section 4. Numerical illustration through two simulated examples and a real data set is reported in Section 5. In particular, we demonstrate via a repeated simulation that the bootstrap provides an adequate approximation for the null-distribution of the test statistic.

2 ESTIMATION OF CONDITIONAL DENSITIES

We assume that data are available in the form of a strictly stationary stochastic process $\{(X_i, Y_i)\}$, where Y_i and X_i are scalars. Naturally, this includes the case where the pairs (X_i, Y_i) are independent and identically distributed. In the time series context, X_i typically denotes a lagged value of Y_i . Let $g(y|x)$ be the conditional density of Y_i given $X_i = x$, which we assume to be smooth in both x and y . We are interested in estimating $g(y|x)$ and its derivatives from the data $\{(X_i, Y_i), 1 \leq i \leq n\}$.

Let $K(\cdot)$ be a symmetric density function on \mathbb{R} and $K_b(u) = b^{-1}K(u/b)$. Note that as $b \rightarrow 0$,

$$E\{K_b(Y_i - y)|X_i = x\} = g(y|x) + O(b^2).$$

This suggests that $g(y|x)$ can be regarded as a regression of $K_b(Y_i - y)$ on X_i . For example, Nadaraya–Watson kernel regression yields the kernel estimator

$$\tilde{g}(y|x) = \sum_{i=1}^n w_i(x) K_b(Y_i - y) \quad (2.1)$$

where

$$w_i(x) = \frac{W_h(X_i - x)}{\sum_{j=1}^n W_h(X_j - x)},$$

$W_h(u) = h^{-1}W(u/h)$, $W(\cdot)$ is a kernel function and $h > 0$ is a bandwidth. This estimator was proposed by Hyndman, Bashtannyk and Grunwald (1996) and is a modification of the esti-

mator proposed by Rosenblatt (1969). Hyndman, Bashtannyk and Grunwald (1996) derive some of its properties and Bashtannyk and Hyndman (2001) explore bandwidth selection rules. Note that there are two smoothing parameters: h controls the smoothness between conditional densities in the x direction (the smoothing parameter for the regression) and b controls the smoothness of each conditional density in the y direction.

The estimator has two desirable properties which match those of the density being estimated: (1) it is always non-negative; and (2) integrals of the estimator with respect to y equal 1. However, it does suffer from the bias problems often associated with kernel smoothers (see Hyndman, Bashtannyk and Grunwald, 1996).

If local polynomial regression is used we obtain the local polynomial estimator proposed by Fan, Yao and Tong (1996). Let

$$R(\theta; x, y) = \sum_{i=1}^n \left\{ K_b(Y_i - y) - \sum_{j=0}^r \theta_j (X_i - x)^j \right\}^2 W_h(X_i - x). \quad (2.2)$$

Then $\hat{g}(y|x) = \hat{\theta}_0$ is a local r th order polynomial estimator where $\hat{\theta}_{xy} = (\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_r)'$ is that value of θ which minimizes $R(\theta; x, y)$. For $r = 0$, this estimator is identical to (2.1). While this estimator has some nice properties such as smaller bias than (2.1) when $r > 0$, it is not restricted to be non-negative and it does not integrate to 1 except in the special case $r = 0$. In this paper, we propose two new estimators which are always non-negative.

2.1 Two New Non-negative Estimators

We replace $R(\theta; x, y)$ by

$$R_1(\theta; x, y) = \sum_{i=1}^n \{ K_b(Y_i - y) - A(X_i - x, \theta) \}^2 W_h(X_i - x). \quad (2.3)$$

where

$$A(x, \theta) = \ell \left(\sum_{j=0}^r \theta_j x^j \right)$$

and $\ell(\cdot)$ is a monotonic function mapping $\mathbb{R} \rightarrow \mathbb{R}^+$. Using $\ell(u) = \exp(u)$ seems a reasonable choice. Then $\hat{g}(y|x) \equiv A(0, \hat{\theta}_{xy}) = \ell(\hat{\theta}_0)$ where $\hat{\theta}_{xy}$ minimizes $R_1(\theta; x, y)$.

We call this the **local parametric estimator**. It is in the same spirit as the local logistic estimator for a conditional distribution function proposed by Hall, Wolff and Yao (1999), and is a conditional version of the density estimator proposed by Loader (1996). Further, it is equivalent to using local likelihood estimation (Tibshirani and Hastie, 1987) for the regression of $K_b(Y_i - y)$ against X_i with the Gaussian likelihood and link function ℓ^{-1} . Consequently, $\hat{\theta}_{xy}$ may be easily computed using local likelihood estimation software such as `locfit` (Loader, 1997). (Note that the `gam` function in S-Plus will not allow a non-identity link function with the Gaussian likelihood.) If an identity link is used ($\ell(u) = u$), we obtain the local polynomial estimator as a special case.

An alternative estimator is obtained by modifying the local linear estimator for $g(y|x)$ directly to force it to be positive. We constrain the minimization of (2.2) so that the coefficient θ_0 is positive. This is achieved by setting $\theta_0 = \ell(\alpha)$ where $\ell(u) = \exp(u)$. We shall denote

this estimator by $\hat{g}_2(y|x)$ and refer to it as the **constrained local polynomial estimator**. Obviously, this idea can also be applied to the problem of estimation of a conditional distribution function, addressed by Hall, Wolff and Yao (1999).

Depending on bandwidth choice, both of these estimators also furnish consistent estimators of the derivatives of the conditional density. Let

$$\begin{aligned} g^{(i)}(y|x) &\equiv \left(\frac{\partial}{\partial y}\right)^i g(y|x), & g^{(lj)}(y|x) &\equiv \left(\frac{\partial}{\partial x}\right)^j g(y|x), \\ g^{(li)}(y|x) &\equiv \left(\frac{\partial}{\partial y}\right)^i \left(\frac{\partial}{\partial x}\right)^j g(y|x), \\ \ell^{(j)}(u) &\equiv \left(\frac{\partial}{\partial u}\right)^j \ell(u) \quad \text{and} \quad A^{(j)}(x, \theta) \equiv \left(\frac{\partial}{\partial x}\right)^j A(x, \theta). \end{aligned}$$

For $j = 1, 2, \dots, r$ we can estimate the density derivatives:

$$\hat{g}_1^{(lj)}(y|x) = A^{(j)}(0, \hat{\theta}_{xy}) = \sum_{k=1}^j \hat{\theta}_k \binom{j-1}{k-1} \ell^{(k)}(\hat{\theta}_0)$$

and

$$\hat{g}_2^{(lj)}(y|x) = j! \hat{\theta}_j.$$

If $K(u)$ is at least q -times differentiable, then for $i = 1, 2, \dots, q$ we can also estimate the density derivatives $\hat{g}_1^{(i)}(y|x)$ and $\hat{g}_2^{(i)}(y|x)$. These are unavailable in closed form but they are easily obtained using numerical differentiation.

In practice, we rescale $\hat{g}(y|x)$, $\hat{g}_1(y|x)$ and $\hat{g}_2(y|x)$ to ensure they integrate to 1. Note that there is no need to rescale the kernel estimator $\hat{g}(y|x)$.

Based on an intentionally biased bootstrap argument of Hall and Presnell (1999), Hall, Wolff and Yao (1999) proposed a modified Nadaraya–Watson estimator for the conditional distribution function which is always non-negative and shares the same first order asymptotic properties as the local linear regression estimator. The same idea can be adapted to the estimation of conditional density functions although we have not pursued this idea here.

2.2 Asymptotic Properties

For the local parametric estimator $\hat{g}_1(y|x)$ we only consider functions A of type $A(x, \theta) = \exp(\theta_0 + \theta_1 x + \dots + \theta_r x^r)$, with $r \geq 1$. Let f denote the marginal density of X_i . We impose the following regularity conditions:

- (C1) For fixed y and x , $f(x) > 0$, $g(y|x) > 0$, f is continuous at x , and $g(y|\cdot)$ has $2[r/2] + 2$ continuous derivatives in a neighbourhood of x , where $[t]$ denotes the integer part of t .
- (C2) The kernels K and W are symmetric, compactly supported probability density functions. Further, $|W(x_1) - W(x_2)| \leq C|x_1 - x_2|$ for any x_1, x_2 .
- (C3) The process $\{(X_i, Y_i)\}$ is absolutely regular, that is

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

where \mathcal{F}_i^j denotes the σ -field generated by $\{(X_k, Y_k) : i \leq k \leq j\}$. Furthermore, $\sum_{j \geq 1} j^2 \beta(j)^{\delta/(1+\delta)} < \infty$ for some $\delta \in [0, 1)$. (We define $a^b = 0$ when $a = b = 0$.)
 (C4) As $n \rightarrow \infty$, $h \rightarrow 0$, $b \rightarrow 0$, $n b h \rightarrow \infty$ and $\liminf_{n \rightarrow \infty} n h^{2(r+1)} > 0$.

Condition (C3) holds with $\delta = 0$ if and only if the process $\{(X_i, Y_i)\}$ is m -dependent for some $m \geq 1$. The requirement of the kernels being compactly supported is imposed for the sake of brevity of proofs. In particular, the Gaussian kernel is allowed. The assumption on the mixing conditions is also not the weakest possible.

Theorem 1 below presents the asymptotic normality of the estimators. The asymptotic expressions for biases and variances are useful in development of the bandwidth selection procedures described in Section 3.

We introduce some notation first. Define

$$\kappa_j = \int u^j W(u) \, du, \quad v_j = \int u^j W^2(u) \, du, \quad \mu_j = \int u^j K(u) \, du,$$

and

$$\lambda_j = \int u^j K^2(u) \, du.$$

Let S denote the $(r+1) \times (r+1)$ matrix with (i, j) th element κ_{i+j-2} , and $\kappa^{(i,j)}$ be the (i, j) th element of S^{-1} . Let $r_1 = 2[r/2] + 2$,

$$\tau_r^2 = \lambda_0 \int \left(\sum_{i=1}^{r+1} \kappa^{(1,i)} v^{i-1} \right)^2 W^2(v) \, dv,$$

$$\eta_r = \frac{1}{(r+1)!} \sum_{i=1}^{r+1} \kappa^{(1,i)} \kappa_{r_1+i-1},$$

and let θ_{xy} be uniquely defined by

$$g(y|x) = A(0, \theta_{xy}), \quad \text{and} \quad g^{(j)}(y|x) = A^{(j)}(0, \theta_{xy}) \quad j = 1, \dots, r. \quad (2.4)$$

Let N_{n1} and N_{n2} denote random variables with the standard Normal distribution.

THEOREM 1 (i) Suppose $r \geq 1$ and conditions (C1)–(C4) hold. Then as $n \rightarrow \infty$,

$$\begin{aligned} \hat{g}_1(y|x) - g(y|x) &= (n b h)^{-1/2} \left\{ \frac{g(y|x)}{f(x)} \right\}^{1/2} \tau_r N_{n1} \\ &\quad + h^{r_1} \eta_r \{ g^{(r_1)}(y|x) - A^{(r_1)}(0, \theta_{xy}) \} \\ &\quad + b^2 \frac{\mu_2}{2} g^{(2)}(y|x) \\ &\quad + o\{(n b h)^{-1/2} + h^{r_1} + b^2\}. \end{aligned} \quad (2.5)$$

(ii) Assume conditions (C1)–(C4) with $r = 1$. Then as $n \rightarrow \infty$,

$$\begin{aligned}\hat{g}_2(y|x) - g(y|x) &= (nhb)^{-1/2} \left\{ \frac{\lambda_0 g(y|x)}{f(x)} \right\}^{1/2} N_{n2} \\ &\quad + h^2 \frac{\kappa_2}{2} g^{(2)}(y|x) + b^2 \frac{\mu_2}{2} g^{(2)}(y|x) \\ &\quad + o\{(nhb)^{-1/2} + h^2 + b^2\}.\end{aligned}\quad (2.6)$$

Remark 1 To the first order, the asymptotic variance of $\hat{g}_1(y|x)$ is exactly the same as in the case of local polynomial estimator $\hat{g}(y|x)$ of order r . This similarity extends also to the bias term, to the extent that for both \hat{g}_1 and local polynomial estimators the bias is of order $O(h^{r+1} + b^2)$ for odd r and $O(h^{r+2} + b^2)$ for even r . However, the form of bias as functionals of the ‘regression mean’ g are quite different. This is a consequence of the fact that $\hat{g}_1(y|x)$ is constrained to be non-negative. In fact, (2.5) would also hold for the local polynomial estimator with order r if we replace the term $A^{(r_1)}(0, \theta_{xy})$ by 0. See Fan and Gijbels (1996) §6.2 or Fan, Yao and Tong (1996). Note, however, that neither reference gives explicitly the bias term in the order h^{r_1} and that the expression they give for r_2^2 contains some typographical errors.

Remark 2 For the linear case ($r = 1$) we have $r_1^2 = \lambda_0 v_0$ and $\eta_1 = \kappa_2/2$. Because of the above remark, (2.6) also holds for the standard local linear estimator. On the other hand, when $\ell(u) = \exp(u)$ and $r = 1$, $A^{(r_1)}(0, \theta_{xy}) = [g^{(1)}(y|x)]^2/g(y|x)$.

Remark 3 For the quadratic case ($r = 2$), we have

$$\tau_2^2 = \frac{\lambda_0(\kappa_4^2 v_0 - 2\kappa_2 \kappa_4 v_2 + \kappa_2^2 v_4)}{(\kappa_4 - \kappa_2^2)^2} \quad \text{and} \quad \eta_2 = \frac{\kappa_4^2 - \kappa_6 \kappa_2}{6(\kappa_4 - \kappa_2^2)}.$$

Remark 4 It may be proved that, under conditions (C1)–(C4) and $r \geq 1$, $\hat{\theta}_{xy} \rightarrow \theta_{xy}$ (see Lemma 1 in the Appendix). Consequently, we may prove that $\hat{g}_1(y|x)$ is a consistent estimator. Similarly, $\hat{g}_2(y|x)$ is also consistent.

3 BANDWIDTH SELECTION

Using (2.5), we find the asymptotic mean square error of $\hat{g}_1(y|x)$ is

$$\begin{aligned}E\{\hat{g}_1(y|x) - g(y|x)\}^2 &\approx \frac{\tau_r^2 g(y|x)}{nhbf(x)} \\ &\quad + \left\{ h^{r_1} \eta_r [g^{(r_1)}(y|x) - A^{(r_1)}(0, \theta_{xy})] + b^2 \frac{\mu_2}{2} g^{(2)}(y|x) \right\}^2,\end{aligned}$$

and so the weighted integrated MSE is

$$\begin{aligned} \text{IMSE} &= \iint \mathbb{E}\{\hat{g}(y|x) - g(y|x)\}^2 f^2(x) \, dx \, dy \\ &= \left\{ \frac{\tau_r^2}{nhb} + \alpha_r h^{2r_1} + \beta_r h^{r_1} b^2 + \gamma b^4 \right\} \{1 + o(1)\} \end{aligned} \quad (3.1)$$

where

$$\alpha_r = \eta_r^2 \iint [g^{(r_1)}(y|x) - A^{(r_1)}(0, \theta_{xy})]^2 f^2(x) \, dx \, dy \quad (3.2)$$

$$\beta_r = \mu_2 \eta_r \iint g^{(2)}(y|x) [g^{(r_1)}(y|x) - A^{(r_1)}(0, \theta_{xy})] f^2(x) \, dx \, dy \quad (3.3)$$

and

$$\gamma = \frac{\mu_2^2}{4} \iint (g^{(2)}(y|x))^2 f^2(x) \, dx \, dy. \quad (3.4)$$

Bashtannyk and Hyndman (2001) used a similar weighted IMSE to derive bandwidths for the estimator (2.1).

Optimal bandwidths for $\hat{g}_1(y|x)$ can be derived by differentiating (3.1) with respect to h and b and setting the derivatives to zero. Solving the resulting equations gives

$$\hat{h} = \left(\frac{\tau_r^2}{nc_r r_1 (2\alpha_r + \beta_r c_r^2)} \right)^{2/(5r_1+2)} \quad \text{and} \quad \hat{b} = c_r (\hat{h})^{r_1/2} \quad (3.5)$$

where

$$c_r = \sqrt{\frac{(r_1 - 2)\beta_r + \sqrt{(r_1 - 2)^2 \beta_r^2 + 32r_1 \alpha_r \gamma}}{8\gamma}}.$$

When $r = 1$, this simplifies to $c_1 = (\alpha_1/\gamma)^{1/4}$. (Because of Remark 2, \hat{h} and \hat{b} in (3.5) are also optimal for $\hat{g}_2(y|x)$ with $r = 1$.) Substituting these optimal bandwidths into (3.1) shows that the IMSE is of order $n^{-4r_1/(5r_1+2)}$. Note that the optimal bandwidth \hat{h} is different from that in standard kernel regression estimation. For example, $\hat{h} = O(n^{-1/6})$ when $r = 1$ while the optimal bandwidth for local linear regression estimation is of order $n^{-1/5}$. Intuitively we need a larger bandwidth (in the order $n^{-1/6}$) to compensate the sparseness of data points due to the smoothing in the y -direction. Similarly, the optimal bandwidth \hat{b} is of order $O(n^{-1/6})$ when for unconditional density estimation, the optimal order is $O(n^{-1/5})$. The larger bandwidth for the conditional estimator is because of the local estimation due to smoothing in the x -direction.

We use these results in the following sections to develop a bandwidth selection strategy. Here we follow the approach of Bashtannyk and Hyndman (2001) in using a mixture of normal reference rules and a regression method. It may be preferable to derive a plug-in rule as Sheather and Jones (1991) have done for univariate density estimation, but this is more difficult to develop.

3.1 Normal Reference Rules

In the kernel estimation of marginal densities, a useful bandwidth selection procedure is to find the optimal bandwidth assuming the normal density (see Silverman, 1986). This has also been used successfully by Bashtannyk and Hyndman (2001) in conditional density estimation with the kernel estimator $\hat{g}(y|x)$ defined by (2.1). Even with non-normal densities, the bandwidths arising from these calculations are usually reasonable.

We shall follow a similar approach for the estimator $\hat{g}_1(y|x)$ and derive optimal bandwidths assuming the conditional distribution and the marginal distribution are both normal. We further assume the conditional distribution has quadratic conditional mean and constant variance σ^2 , and that the marginal distribution of X has mean μ and variance v^2 .

Then we can write

$$g(y|x) = \frac{1}{\sigma} \phi\left(\frac{y - d_0 - d_1(x - \mu) - d_2(x - \mu)^2}{\sigma}\right)$$

and

$$f(x) = \frac{1}{v} \phi\left(\frac{x - \mu}{v}\right).$$

Substituting these into (3.2)–(3.4), we obtain

$$\gamma = \frac{3\mu_2^2}{64\pi\sigma^5v}, \quad \alpha_1 = \frac{\kappa_2^2(2d_2^2\sigma^2 + d_1^4 + 12d_2^2v^2(d_1^2 + d_2^2v^2))}{16\pi\sigma^5v},$$

$$\beta_1 = \frac{\mu_2\kappa_2(d_1^2 + 2d_2^2v^2)}{16\pi\sigma^5v}$$

and $c_1 = (\alpha_1/\gamma)^{1/4}$ when the log link ($\ell(u) = \exp(u)$) is used. For the local linear estimator ($\ell(u) = u$), we obtain the same γ and c_1 values, with

$$\alpha_1 = \frac{\kappa_2^2(8d_2^2\sigma^2 + 3d_1^4 + 36d_2^2v^2(d_1^2 + d_2^2v^2))}{64\pi\sigma^5v}$$

and

$$\beta_1 = \frac{3\mu_2\kappa_2(d_1^2 + 2d_2^2v^2)}{32\pi\sigma^5v}.$$

The local quadratic estimator is more difficult and we only give the bandwidths for the identity link ($\ell(u) = u$) assuming the conditional mean is linear (*i.e.*, $d_2 = 0$). Then we obtain the same γ with

$$\alpha_2 = \frac{105\eta_2^2d_1^8}{64\pi\sigma^9v}, \quad \beta_2 = \frac{-15\eta_2\mu_2d_1^4}{32\pi\sigma^7v} \quad \text{and}$$

$$c_2^2 = \frac{|\eta_2|d_1^4(\sqrt{305} - 5 \operatorname{sign}(\eta_2))}{2\mu_2\sigma^2}$$

where $\operatorname{sign}(u) = u/|u|$.

In the special case where both $W(u)$ and $K(u)$ denote a standard normal kernel, and the conditional mean is linear ($d_2 = 0$), we substitute the above values into (3.5) to obtain the following simple rules:

- When $r = 1$ and $\ell(u) = \exp(u)$, $\hat{h} \approx 0.916(v\sigma^5/n|d_1|^5)^{1/6}$ and $\hat{b} = 1.05|d_1|\hat{h}$.
- When $r = 1$ and $\ell(u) = u$, $\hat{h} \approx 0.935(v\sigma^5/n|d_1|^5)^{1/6}$ and $\hat{b} = |d_1|\hat{h}$.

- When $r = 2$ and $\ell(u) = u$, $\hat{h} \approx 0.703(v\sigma^{10}/nd_1^{10})^{1/11}$ and $(\hat{b} \approx 2.37d_1^2/\sigma)(\hat{h})^2$.

3.2 A Bandwidth Selection Algorithm

For a given bandwidth b and a given value y , finding $\hat{g}(y|x)$ is a standard nonparametric problem of regressing $K_b(Y_i - y)$ on X_i . Therefore, we can adapt bandwidth selection methods used in regression for use in this problem. Let $M_b(h; y)$ denote a goodness-of-fit statistic for the regression of $K_b(Y_i - y)$ on X_i with bandwidth h . For example, $M_b(h; y)$ may denote the generalized cross-validation statistic (Fan and Gijbels, 1996, p. 45). We then define

$$M_b(h) = \sum_{j=1}^N M_b(h; y'_j)$$

where $\mathbf{y} = \{y'_j, \dots, y'_N\}$ are equally spaced in the sample space of Y . For a given value of b , $M_b(h)$ may be minimized to select a value of h . This approach was suggested by Bashtannyk and Hyndman (2001) for the kernel estimator with $M_b(h; y)$ denoting the penalized average square prediction error (see, for example, Härdle, 1991). Fan, Yao and Tong (1996) suggested a similar approach for the local polynomial estimator with $M_b(h)$ denoting the Residual Squares Criterion proposed by Fan and Gijbels (1995).

When this approach is combined with the normal reference rules, we have a useful algorithm for selecting the bandwidth parameters.

1. Select the smoothing parameter b using the normal reference rule.
2. Given this value of b , minimize $M_b(h)$ to find a value for h .

4 BOOTSTRAP TESTS FOR SYMMETRY

We are interested in testing for the symmetry of a conditional density function $g(y|x)$ at a particular value of x . If the conditional density is shown to be symmetric at x , then a more efficient estimator of $g(y|x)$ can be constructed (see Remark 6). Note that in interval forecasting of time series, the conditional (rather than unconditional) distributions are relevant; see Polonik and Yao (2000). Conditional symmetry is helpful in constructing predictive intervals as both tails of the density can be used to estimate the boundaries of the intervals.

For fixed x with $f(x) > 0$, we are interested in testing the hypothesis that the conditional distribution $g(\cdot|x)$ is symmetric, that is

$$H_0 : g(y|x) = g(2u(x) - y|x) \quad \text{for any } y,$$

where $u(x)$ is the centre of the conditional distribution of $g(\cdot|x)$. Under hypothesis H_0 , we would expect that the above equality also holds approximately for a good estimator of g , say \hat{g} . Therefore, we define the test statistic

$$T(x) = \min_u \int \{\hat{g}(y|x) - \hat{g}(2u - y|x)\}^2 dy$$

and reject H_0 for large values of T .

To derive the asymptotic distribution of T (under H_0) is a tedious matter. Typically the sample size n must be very large to ensure asymptotic results are adequately accurate in non-parametric tests (see, for example, Hjellvik, Yao and Tjøstheim, 1998). Therefore we adopt a bootstrap approach in this paper.

Note all the estimators described in Section 2 can be written as linear forms of $\{K_b(Y_i - y)\}$ as follows

$$\widehat{g}(y|x) = \sum_{i=1}^n m_i(x) K_b(Y_i - y),$$

where the weight $m_j(x)$ depends on $\{X_i\}$ and x only. Note the kernel function $K(\cdot)$ is symmetric. It is easy to see that

$$\widehat{g}(2u(x) - y|x) = \sum_{i=1}^n m_i(x) K_b(2u(x) - Y_i - y).$$

This means that the mirror reflection of the estimator $\widehat{g}(\cdot|x)$ with respect to $u(x)$ is \widehat{g} itself obtained with the sample $\{(Y_i, X_i)\}$ replaced by $\{(2u(x) - Y_i, X_i)\}$. This motivates the following resampling scheme.

1. We calculate

$$u(x) = \arg \min_u \int \{\widehat{g}(y|x) - \widehat{g}(2u - y|x)\}^2 dy. \quad (4.1)$$

2. We sample n independent observations $\{X_i^*, 1 \leq i \leq n\}$ from $\{X_i, 1 \leq i \leq n\}$ with replacement.
3. Suppose $X_i^* = X_{i_j}$. For each $1 \leq i \leq n$, sample Y_i^* from the uniform distribution on the two symmetric points Y_{i_j} and $2u(x) - Y_{i_j}$.
4. Form the statistic T^* in the same way as T with $\{X_i, Y_i\}$ replaced by $\{X_i^*, Y_i^*\}$.

We reject H_0 if T is greater than the upper α -point of the conditional distribution of T^* given $\{X_i, Y_i\}$. In fact, the p -value is the relative frequency of the event $\{T^* \geq T\}$ in the bootstrap replications.

We may let \widehat{g} be the local parametric estimator \widehat{g}_1 with $r = 1$ or the constrained local linear estimator \widehat{g}_2 . We use the same method to choose the bandwidth for the original data and bootstrap data.

Since we only test the symmetry of $g(\cdot|x)$ at fixed x , one would expect that we only sample Y_i^* from a symmetric distribution when X_i^* is close to x . This is effectively achieved in the nonparametric estimation of $g(\cdot|x)$, since the estimation is localized by the kernel function.

When we generate the bootstrap samples, we largely ignore the possible dependence in the data. Note that under the mixing condition (C3), the dependence does not enter the major terms (*i.e.*, first order terms) in the asymptotic expansions in Theorem 1. This is due to the fact that in nonparametric regression (with random design), we only use effectively the nh nearest neighbours in the *state space*, which are unlikely to be the neighbours in the *time space* under the mixing condition (C3). Those points could be regarded as asymptotically independent when $n \rightarrow \infty$. In fact we may prove that it holds almost surely that the conditional distribution of T^* given $\{X_i, Y_i\}$ is asymptotically equal to the null-hypothesis distribution of T (cf. Kreiss, Neumann and Yao, 1998).

Remark 5 Note that since $f(x) > 0$, the null hypothesis can be expressed equivalently as $H_0 : g(y|x)f(x) = g(2u(x) - y|x)f(x)$ for any y . Furthermore, the joint density function $p(x, y) \equiv g(y|x)f(x)$ can be easily estimated. For example, the simple product kernel estimator is $\hat{p}(x, y) = 1/n \sum_{i=1}^n W_h(X_i - x)K_b(Y_i - y)$. Therefore, an alternative test statistic can be defined as $T_1(x) = \min_u \int \{\hat{p}(x, y) - \hat{p}(x, 2u - y)\}^2 dy$. The bootstrap procedure described above can be applied to facilitate this alternative test.

Remark 6 When the density is symmetric, a symmetric estimator may be obtained as

$$\hat{\hat{g}}(y|x) = \frac{1}{2}(\hat{g}(y|x) + \hat{g}(2u(x) - y|x)). \quad (4.2)$$

See Kraft, Lepage and van Eeden (1985) and Meloche (1991) for further discussion on estimation of symmetric densities. Note that for most values of x and y , $\hat{\hat{g}}(y|x)$ will have smaller variance than $\hat{g}(y|x)$. In the numerical examples, we estimate the density by (4.2) if $\hat{g}(y|x)$ passes the symmetry test.

5 NUMERICAL EXAMPLES

We illustrate the symmetry tests through simulations and by application to some real data. In all cases, we have used a truncated Gaussian kernel,

$$K(u) = W(u) = \begin{cases} \exp(-u^2/2)/\sqrt{2\pi} & |u| < 10; \\ 0 & \text{otherwise.} \end{cases}$$

(The truncation is used to satisfy the finite domain requirement of C2, although in practice it has negligible effect.)

Example 1 Consider the model $Y_i = 5 + (1 + W_i)X_i + \varepsilon_i$ where $\{X_i\}$, $\{W_i\}$ and $\{\varepsilon_i\}$ are all independent with X_i uniformly distributed on $[0, 12]$, ε_i normally distributed with zero mean and variance 9, and W_i is a binary variable with $\Pr(W_i = 1) = 1 - \Pr(W_i = 0) = 0.3$. Figure 1 shows a scatterplot of 500 observations from this model. The line through the points is $u(x)$ calculated from (4.1). When $x = 0$, the density is symmetric, and it increases in skewness as x increases. For $x \leq 6$ the skewness is hardly visible from Figure 1 due to the masking effect from the large variance of ε_i .

We computed the p -value of the bootstrap test for symmetry for $0 \leq x \leq 12$ at steps of 0.5. For these tests, we used the local parametric estimator of $g(y|x)$ with $r = 1$ and bandwidths chosen using the algorithm of Section 3.2 to be $h = 1.35$ and $b = 1.59$. (For this example, the true optimal bandwidths calculated using (3.5) are $\hat{h} = 0.87$ and $\hat{b} = 1.25$.) Figure 2 shows the p -values. Each test involved 100 replications. The skewness is clearly detected by the tests for $x > 6$.

To demonstrate that the bootstrap method does provide an accurate approximation for the distribution of the test statistic under H_0 , we modify the above model in order (i) to make $x = 0$ an inner point in the sample space, and (ii) to reduce the masking effect for the asymmetry due to large errors. The modified model is

$$Y_i = 2.5 + (1 + W_i)X_i + \varepsilon_i,$$

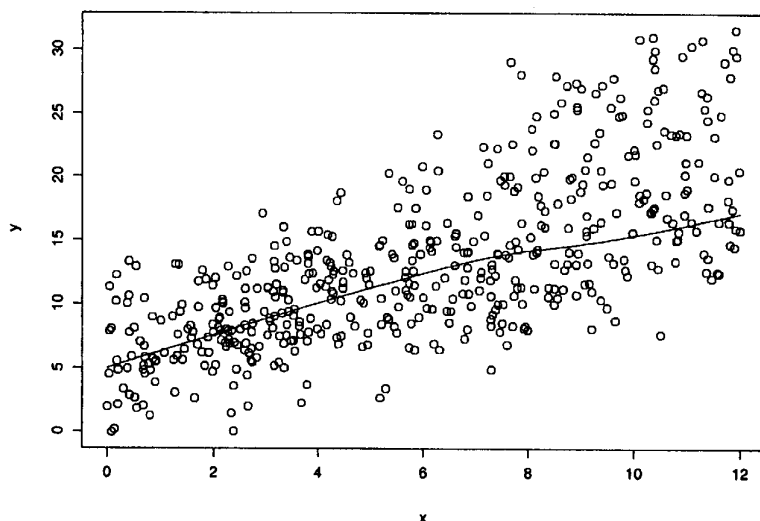


FIGURE 1 Scatterplot of 500 observations from Example 1. The line through the points is $u(x)$, the estimated centre of symmetry, calculated from (4.1).

where $X_i \stackrel{d}{=} U(-6, 6)$, $\varepsilon_i \stackrel{d}{=} N(0, 0.5^2)$, and W_i , unchanged. Note that the conditional distribution of Y_i , given $X_i = x$ is *strictly* symmetric if and only if $x = 0$. Further the reduction of the noise level is in favour of the rejection of H_0 . Our simulation shows that the bootstrap test leads to the correct inference (*i.e.*, not to reject H_0 when $x = 0$).

We let $x = 0$ and $n = 500$. Note for $x = 0$, the conditional distribution of Y_i , given $X_i = x$ is normal with mean 2.5 and standard deviation 0.5. To speed up the computation, the normal reference rules are employed to select the bandwidths. Figure 3 plots the empirical distribution of the test statistic $T(x)$ in the simulation with 200 replications, together with three bootstrap approximations. The three bootstrap approximations were selected in such a way that

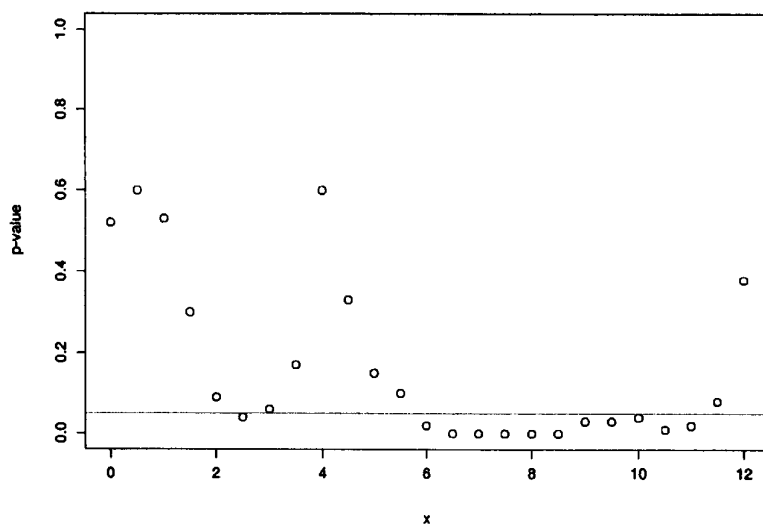


FIGURE 2 The p -values of the bootstrap test for symmetry of the conditional density $g(y|x)$ in Example 1. Here $\hat{g}(y|x)$ is the local parametric estimate with $r = 1$ and bandwidths chosen using the normal reference rules to be $h = 1.1$ and $b = 1.6$. The horizontal line shows the 0.05 level.

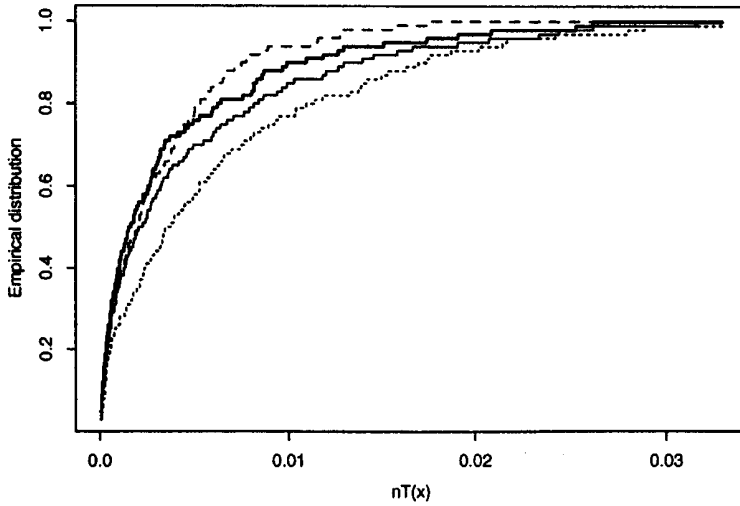


FIGURE 3 The plots of the sampling distribution of $nT(x)$ (thick solid lines) and its bootstrap approximations: first quartile (dotted lines), median (thin solid lines) and third quartile (dashed lines).

the corresponding p -values were equal to the first quartile, the median and the third quartile. Figure 3 shows that the bootstrap approximation is fairly accurate.

Example 2 We next consider a quadratic AR(1) time series model

$$Y_t = 0.23Y_{t-1}(16 - Y_{t-1}) + 0.4\varepsilon_t \quad (5.1)$$

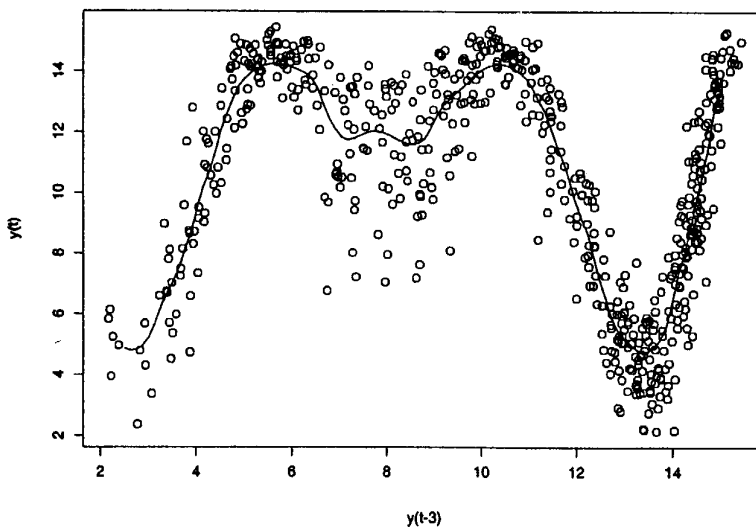


FIGURE 4 Scatterplot of 600 observations from Example 2. The line through the points is $u(x)$, the estimated centre of symmetry, calculated from (4.1).

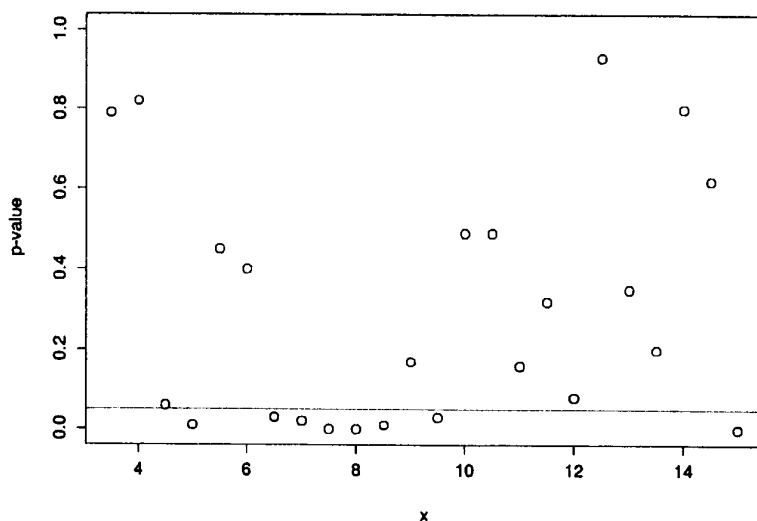


FIGURE 5 The p -values of the bootstrap test for symmetry of the conditional density $g(y|x)$ in Example 2. Here $\hat{g}(y|x)$ is the kernel estimate with bandwidths $h = b = 0.5$. The horizontal line shows the 0.05 level.

where $\{\varepsilon_t\}$ is a sequence of independent random variables each with the standard normal distribution truncated in the interval $[-12, 12]$. The conditional distribution of Y_t given $X_t \equiv Y_{t-m}$ is symmetric for $m = 1$ but not necessarily so for $m > 1$. Figure 4 shows a lagged scatterplot of 600 observations from this model with $m = 3$. The line through the points is $u(x)$ calculated from (4.1) where $\hat{g}(y|x)$ is the local parametric estimate with $r = 1$. Bandwidths were chosen using the algorithm to be $h = 0.4$ and $b = 1.2$. For each of the bootstrap tests, 100 replications were performed. The p -values from the bootstrap test for symmetry are shown in Figure 5. There is a clear evidence that the conditional distribution is not symmetric for x between 6.5 and 8.5.

To demonstrate that our bootstrap approximation works, we conduct simulations with $X_t = Y_{t-1}$. Then the conditional distribution of Y_t given $X_t = x$ is normal with mean $0.23x(16 - x)$ and variance 0.4^2 . For $x = 5$, we simulate 200 data sets for each of $n = 600$ and $n = 1200$. Figure 6 shows that the bootstrap approximations with $n = 600$ tend to be

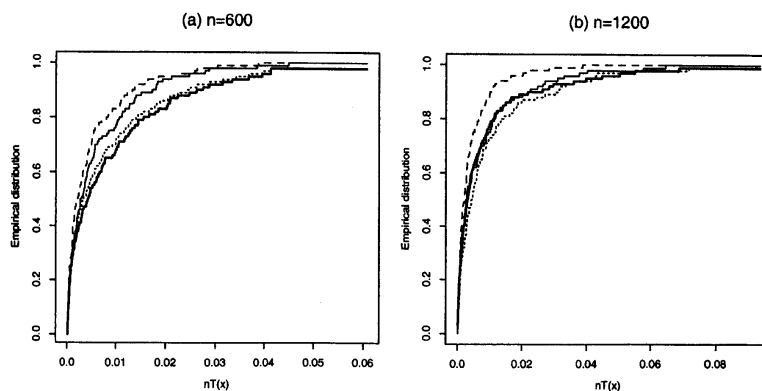


FIGURE 6 The plots of the sampling distribution of $nT(x)$ (thick solid lines) and its bootstrap approximations: first quartile (dotted lines), median (thin solid lines) and third quartile (dashed lines).

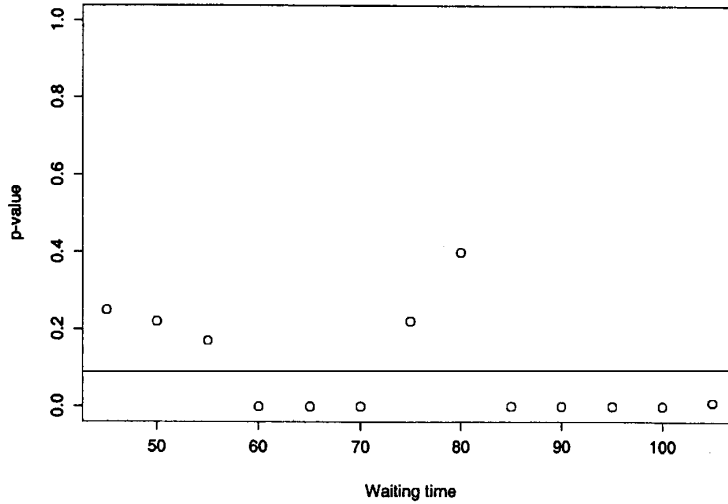


FIGURE 7 The p -values of the bootstrap test for symmetry of the density of the Old Faithful geyser eruption duration conditional on the waiting time between eruptions. Here $\hat{g}(y|x)$ is the kernel estimate with bandwidths $h = 7.2$ and $b = 0.41$ chosen using the bandwidth selection algorithm. The horizontal line shows the 0.05 level.

biased in the sense that the bootstrap distributions seem to have heavier tails on the left. By increasing the sample size to $n = 1200$, the approximation is more satisfactory. This seems to suggest that a large sample size is required to ensure the estimator behaves like the one based on independent data.

5.1 Old Faithful Geyser data

Azzalini and Bowman (1990) give data on the waiting time between the starts of successive eruptions and the duration of the subsequent eruption for the Old Faithful geyser in Yellowstone National Park, Wyoming. The data were collected continuously between 1–15 August

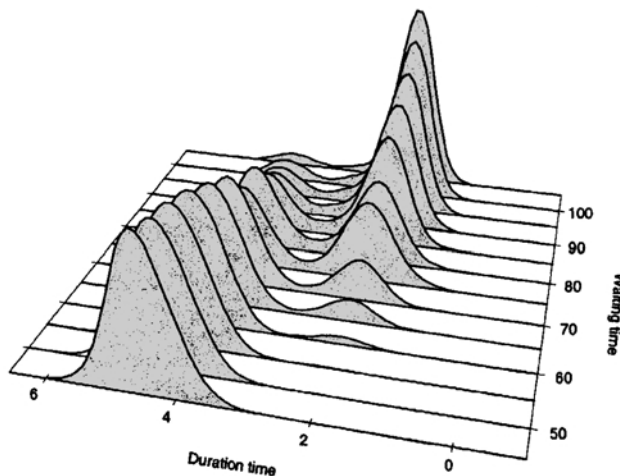


FIGURE 8 Estimated conditional density of eruption duration conditional on waiting time to the eruption. The densities have been symmetrized if the p -values in Figure 7 are greater than 0.05. Bandwidths were chosen using the selection algorithm.

1985. There are a total of 299 observations. The times are measured in minutes. Some duration measurements, taken at night, were originally recorded as S (short), M (medium), and L (long). These values have been coded as 2, 3 and 4 minutes respectively. This data set is also distributed with S-Plus.

We are interested in the distribution of duration time conditional on the previous waiting time. The bandwidth selection algorithm gives bandwidths $h = 8.1$ and $b = 0.33$. Using these values we test the symmetry of the conditional densities (again using the kernel estimator (2.1)) with 100 replications per test. The p -values from the bootstrap test for symmetry are shown in Figure 7. Where the p -value is greater than 0.05, we replace $\hat{g}(y|x)$ by the symmetric estimator (4.2). The resulting estimates are shown in Figure 8 using the stacked density visualization method of Hyndman, Bashtannyk and Grunwald (1996).

Acknowledgements

This work was carried out while Rob Hyndman was a visitor to the Department of Statistics, Colorado State University and Qiwei Yao was a visitor to the Australian National University. Rob Hyndman was supported in part by an Australian Research Council grant. Qiwei Yao was supported partially by EPSRC Grant L16385 and BBSRC/EPSRC Grant 96/MMI09785. The authors would like to thank Clive Loader for making his locfit software available, and Chris Jones for some helpful comments.

References

- Azzalini, A. and Bowman, A.W. (1990) "A look at some data on the Old Faithful geyser", *Applied Statistics* **39**, 357–365.
- Bashtannyk, D. M. and Hyndman, R. J. (2001) "Bandwidth selection for kernel conditional density estimation", *Computat. Statist. and Data Anal.* **36**(3), 279–298.
- Brännäs, K. and De Gooijer, J. G. (1992) "Modelling business cycle data using autoregressive-asymmetric moving average models", *ASA Proceedings of Business and Economic Statistics Section* 331–336.
- Butler, C. (1969) "A test for symmetry using the sample distribution function", *Ann. Math. Statist.* **40**, 2209–2210.
- Copas, J. B. (1995) "Local likelihood based on kernel censoring", *J. Roy. Statist. Soc. Ser. B* **57**, 221–235.
- Csörgö, S. and Heathcote, C. R. (1987) "Testing for symmetry", *Biometrika* **74**, 177–184.
- De Gooijer, J. G. and Gannoun, A. (2000) "Nonparametric conditional predictive regions for time series", *Computational Statistics and Data Analysis* **33**, 259–275.
- Diks, C. and Tong, H. (1999) "A test for symmetries of multivariate probability distributions", *Biometrika* **86**(3), 605–614.
- Doksum, K. A., Fenstad, G. and Aaberge, R. (1977) "Plots and tests for symmetry", *Biometrika* **64**, 473–487.
- Fan, J. and Gijbels, I. (1995) "Data-driven bandwidth selection in local polynomial fitting: variable bandwidth and spatial adaptation", *J. R. Statist. Soc. B* **57**, 371–394.
- Fan, J. and Gijbels, I. (1996) *Local polynomial estimation* (Chapman and Hall, New York).
- Fan, J., Yao, Q. and Tong, H. (1996) Estimation of conditional densities and sensitivity measures in nonlinear dynamical systems. *Biometrika* **83**, 189–206.
- Hall, P. and Presnell, B. (1999) "Intentionally biased bootstrap methods", *J. Royal Statist. Soc. B* **61**, 143–158.
- Hall, P., Wolff, R. and Yao, Q. (1999) "Methods for estimating a conditional distribution function", *J. Amer. Statist. Assoc.* **94**, 154–163.
- Härdle, W. (1991) *Smoothing Techniques with Implementation in S* (Springer-Verlag, New York).
- Hill, D. L. and Rao, P. V. (1977) "Tests for symmetry based on Cramér-von Mises statistics", *Biometrika* **64**, 489–494.
- Hjellvik, V., Yao, Q. and Tjøstheim, D. (1998) "Linearity testing using local polynomial approximation", *J. Statist. Plann. Infer.* **68**(2), 295–321.
- Hjort, N. L. and Jones, M. C. (1996) "Locally parametric nonparametric density estimation", *Ann. Statist.* **24**, 1619–1647.
- Hollander, M. (1971) "A nonparametric test for bivariate symmetry", *Biometrika* **71**, 203–212.
- Hyndman, R. J. (1995) "Highest density forecast regions for non-linear and non-normal time series models", *J. Forecasting* **14**, 431–441.
- Hyndman, R. J., Bashtannyk, D. M. and Grunwald, G. K. (1996) Estimating and visualizing conditional densities, *J. Comp. Graph. Statist.* **5**(4), 315–336.
- Kraft, C. H., Lepage, Y. and Van Eeden, C. (1985) "Estimation of a symmetric density function", *Communications in Statistics: Theory and Methods* **14**, 273–288.

- Kreiss, J. P., Neumann, M. and Yao, Q. (1998) Bootstrap tests for simple structures in non-parametric time series regression, Submitted.
- Loader, C. R. (1996) "Local likelihood density estimation", *Ann. Statist.* **24**, 1602–1618.
- Loader, C. (1997) "Locfit: an introduction", *Statistical Computing and Graphics Newsletter* **8**(1), 11–17.
- Lockhart, R. A. and McLaren, C. G. (1985) "Asymptotic points for a test of symmetry about a specified mean", *Biometrika* **85**, 208–210.
- Meloche, J. (1991) "Estimation of a symmetric density", *Canadian J. Statist.* **19**, 151–164.
- Peligrad, M. (1986). "Recent advances in the central limit theorem and its weak invariance principle for mixing sequences of random variables", In: Eberlein, E. and Taqqu, M. S. (Eds.), *Dependence in Probability and Statistics* (Birkhäuser, Boston), 193–223.
- Polonik, W. and Yao, Q. (2000) "Conditional minimum volume predictive regions for stochastic processes", *J. Amer. Statist. Assoc.* **95**, 509–519.
- Rosenblatt, M. (1969) "Conditional probability density and regression estimators", In: Krishnaiah, P. (Ed.), *Multivariate Analysis II* (Academic Press, New York), pp. 25–31.
- Rothman, E. N. D. and Woodroffe, M. A. (1972) "A Cramér-von Mises type statistic for testing symmetry", *Ann. Math. Statist.* **43**, 2035–2038.
- Sheather, S. J. and Jones, M. C. (1991) "A reliable data-based bandwidth selection method for kernel density estimation", *J. Roy. Statist. Soc. Ser. B* **53**, 683–690.
- Silverman, B. W. (1986) *Density Estimation for Statistics and Data Analysis* (Chapman and Hall, London).
- Simonoff, J. S. (1996) *Smoothing Methods in Statistics* (Springer, New York).
- Srinivasan, R. and Godio, L. B. (1974) A "Cramér-von Mises type statistic for testing symmetry", *Biometrika* **61**, 196–198.
- Tibshirani, R. and Hastie, T. (1987) "Local likelihood estimation", *J. Amer. Statist. Assoc.* **82**, 559–567.
- Yao, Q. and Tong, Q. (2000) "Nonparametric estimation of ratios of noise to signal in stochastic regression", *Statistica Sinica* **10**, 751–770.
- Yu, K. and Jones, M. C. (1998) "Local linear quantile regression", *J. Amer. Statist. Assoc.* **93**, 228–237.
- Zhu, L.-X. (1998) "Assessing elliptical symmetry via a computer-assisted test procedure", *J. Amer. Statist. Assoc.*, to appear.

6 APPENDIX: PROOF OF THEOREM 1

We only prove (2.5); Eq. (2.6) can be proved in a much simpler manner. We use the same notation as in Section 2. We always assume that conditions (C1)–(C4) hold and $r \geq 1$. We first introduce a lemma.

LEMMA 1 As $n \rightarrow \infty$, $\hat{\theta}_{xy} \rightarrow \theta_{xy}$ in probability.

Proof Since $\hat{\theta}_{xy}$ is the minimiser of $R_1(\theta; x, y)$ defined in (2.3), $D_n(x, y, \hat{\theta}_{xy}) = 0$, where

$$D_n(x, y, \theta) = \frac{1}{nh^r} \sum_{i=1}^n \{K_b(Y_i - y) - A(X_i - x, \theta)\} \\ \times A(X_i - x, \theta) W_k(X_i - x) \left(1, \frac{X_i - x}{h}, \dots, \left(\frac{X_i - x}{h}\right)^r\right)^\tau.$$

Define

$$D(x, y, \theta, h) = \frac{f(x)}{h^r} \int (1, t, \dots, t^r)^\tau A(0, \theta) W(t) dt \\ \times \sum_{i=0}^r \frac{(ht)^i}{i!} \{g^{(i)}(y|x) - A^{(i)}(0, \theta)\}.$$

It is easy to see that $D(x, y, \theta_{xy}, h) \equiv 0$. Further, it can be proved that for any compact set G ,

$$\sup_{\theta \in G} \|D_n(x, y, \theta) - D(x, y, \theta, h)\| \xrightarrow{P} 0.$$

Let us assume that $\hat{\theta}_{xy} \not\xrightarrow{P} \theta_{xy}$. Then there exists a sub-sequence of $\{n\}$, still denoted as $\{n\}$ for the simplicity of notation, for which $P\{\|\hat{\theta}_{xy} - \theta_{xy}\| > \varepsilon\} > \varepsilon$ for all sufficiently large n , where $\varepsilon > 0$ is a constant. Consequently, $\inf_{\|\theta - \theta_{xy}\| \leq \varepsilon} \|D_n(x, y, \theta)\| \not\xrightarrow{P} 0$. Hence we have that

$$\begin{aligned} \inf_{\|\theta - \theta_{xy}\| \leq \varepsilon} \|D((x, y, \theta, h))\| &\geq \inf_{\|\theta - \theta_{xy}\| \leq \varepsilon} \|D_n(x, y, \theta)\| \\ &\quad - \sup_{\|\theta - \theta_{xy}\| \leq \varepsilon} \|D_n(x, y, \theta) - D(x, y, \theta, h)\| \\ &= \inf_{\|\theta - \theta_{xy}\| \leq \varepsilon} \|D_n(x, y, \theta)\| + o_p(1) \not\xrightarrow{P} 0, \end{aligned}$$

which contradicts the fact that $D(x, y, \theta_{xy}, h) \equiv 0$. Therefore, $\hat{\theta}_{xy} \xrightarrow{P} \theta_{xy}$. ■

Proof of (2.5) For any $\varepsilon \in (0, 1)$, it follows from Lemma 1 that there exists $\varepsilon_1 \in (0, \infty)$ for which $P\{\|\hat{\theta}_{xy} - \theta_{xy}\| \leq \varepsilon_1\} \geq 1 - \varepsilon$ for all sufficiently large n . Let $G \equiv G(\varepsilon_1)$ be the closed ball centered at θ_{xy} with radius ε_1 . Let $\hat{\theta}_{xy,G}$ be the minimizer of (2.3) with θ restricted on G . Define $\hat{g}_G(y|x) = A(0, \hat{\theta}_{xy,G})$. Then $P\{\hat{g}_G(y|x) \neq \hat{g}(y|x)\} < \varepsilon$ for all sufficiently large n . The above argument indicates that we only need to establish (2.5) for $\hat{g}_G(y|x)$. Therefore we proceed the proof below by assuming θ_{xy} is always within a compact set G .

We consider only the case that r is odd and δ given in condition (C3) is positive. Note that $W(\cdot)$ has a bounded support. By a simple Taylor expansion on A in (2.3), we have that

$$\begin{aligned} R_1(\theta; x, y) &= \sum_{i=1}^n \left(K_b(Y_i - y) - \sum_{j=0}^r \frac{A^{(j)}(0, \theta)}{j!} (X_i - x)^j \right. \\ &\quad \left. - \frac{A^{(r+1)}(c_i(X_i - x), \theta)}{(r+1)!} (X_i - x)^{r+1} \right)^2 W_h(X_i - x), \end{aligned}$$

where $c_i \in [0, 1]$. Define $R_1^*(\theta; x, y)$ as $R_1(\theta; x, y)$ with θ in $A^{(r+1)}(c_i(X_i - x), \theta)$ replaced by $\hat{\theta}_{xy}$. Let $\hat{\theta}_{xy}^*$ be the minimizer of $R_1^*(\theta; x, y)$, and $\hat{g}_1^*(y|x) = A(0, \hat{\theta}_{xy}^*)$. In the sequel, we first prove that (2.5) holds for $\hat{g}_1^*(y|x)$. Then we show that

$$\hat{g}_1(y|x) = \hat{g}_1^*(y|x) + o_p(h^{r+1}). \quad (7.1)$$

It is easy to see that (2.5) follows from the above two statements immediately.

It follows from least squares theory that

$$\begin{aligned} \hat{g}_1^*(y|x) - g(y|x) &= \frac{1}{nh} \sum_{i=1}^n W_n\left(\frac{X_i - x}{h}, x\right) \\ &\quad \times \left\{ \epsilon_i \frac{1}{(r+1)!} \{g^{(r+1)}(y|x + c'_i(X_i - x)) \right. \\ &\quad \left. - A^{(r+1)}(c_i(X_i - x), \hat{\theta}_{xy})\} (X_i - x)^{r+1} \right\}, \end{aligned} \quad (7.2)$$

where $\epsilon_i = K_b(X_i - x) - g(y|x)$, $c'_i \in [0, 1]$,

$$W_n(u, x) = (1, 0, \dots, 0) S_n(x)^{-1} (1, u, \dots, u^r)^\tau W(u),$$

and $S_n(x)$ is an $(r+1) \times (r+1)$ matrix with $s_{i+j-2}(x)$ as its (i, j) th element, and

$$s_j(x) = \frac{1}{nh} \sum_{i=1}^n W_h(X_i - x)(X_i - x)^j.$$

(See, for example, (3.11) of Fan and Gijbels 1996.) It follows from the ergodic theorem that $S_n(x) \xrightarrow{P} f(x)(\kappa_{i+j-2})$. We write

$$\xi_i = \sum_{j=1}^{r+1} \kappa^{(1,j)} \left(\frac{X_i - x}{h} \right)^{j-1},$$

$$\eta_i = [g^{(r+1)}(y|x + c'_i(X_i - x)) - A^{(r+1)}(c_i(X_i - x), \hat{\theta}_{xy})]/(r+1)!.$$

We have that

$$\begin{aligned} \hat{g}_1^*(y|x) - g(y|x) &= \left\{ \frac{1}{nhf(x)} \sum_{i=1}^n \xi_i W\left(\frac{X_i - x}{h}\right) \right. \\ &\quad \left. \times \{\epsilon_i + \eta_i(X_i - x)^{r+1}\} \right\} \{1 + o_p(1)\}. \end{aligned}$$

(See Lemmas 1 and 2 of Yao and Tong, 2000.) Note that we have assumed that $\hat{\theta}_{xy} \in G$. It follows from Theorem 1.7 of Peligrad (1986) and the ergodic theorem that the RHS of the above expression admits the asymptotic expansion in the RHS of (2.5).

To prove (7.1), note that all the $A^{(i)}(0, \hat{\theta}_{xy}^*)$ ($i = 0, 1, \dots, r$) have explicit expressions such as (7.2). Therefore, it is easy to prove that $A^{(i)}(0, \hat{\theta}_{xy}^*) \rightarrow A^{(i)}(0, \theta_{xy})$, where θ_{xy} is determined by (2.4). This implies that $\hat{\theta}_{xy}^* \rightarrow \theta_{xy}$ (see Lemma 1 above). Consequently, $|\hat{\theta}_{xy}^* - \theta_{xy}| \xrightarrow{P} 0$, which implies that $R_1(\hat{\theta}_{xy}^*; x, y) = R_1^*(\hat{\theta}_{xy}^*; x, y) + o_p(nh^{2(r+1)})$, because $\partial R_1^*(\theta; x, y)/\partial \theta = 0$ at

$\theta = \hat{\theta}_{xy}^*$. Note that $R_1(\hat{\theta}_{xy}; x, y) = R_1^*(\hat{\theta}_{xy}; x, y)$ and $\hat{\theta}_{xy}$ and $\hat{\theta}_{xy}^*$ are the minimizers of R_1 and R_1^* . From

$$0 < R_1(\hat{\theta}_{xy}^*; x, y) - R_1(\hat{\theta}_{xy}; x, y) = R_1^*(\hat{\theta}_{xy}^*; x, y) - R_1^*(\hat{\theta}_{xy}; x, y) + o_p(nh^{2(r+1)}),$$

we have that

$$\frac{1}{n}R_1(\hat{\theta}_{xy}; x, y) = \frac{1}{n}R_1(\hat{\theta}_{xy}^*; x, y) + o_p(h^{2(r+1)}).$$

Since $\partial R_1(\theta; x, y)/\partial \theta = 0$ at $\theta = \hat{\theta}_{xy}$, the above expression implies that

$$h^{-2(r+1)}(\hat{\theta}_{xy} - \hat{\theta}_{xy}^*)^r \tilde{R}(\hat{\theta}_{xy})(\hat{\theta}_{xy} - \hat{\theta}_{xy}^*) = \left(\frac{\hat{\theta}_{xy,0} - \hat{\theta}_{xy,0}^*}{h^{(r+1)}}, \frac{\hat{\theta}_{xy,1} - \hat{\theta}_{xy,1}^*}{h^r}, \dots, \frac{\hat{\theta}_{xy,r} - \hat{\theta}_{xy,r}^*}{h} \right) R^* \begin{pmatrix} \frac{\hat{\theta}_{xy,0} - \hat{\theta}_{xy,0}^*}{h^{(r+1)}} \\ \frac{\hat{\theta}_{xy,1} - \hat{\theta}_{xy,1}^*}{h^r} \\ \vdots \\ \frac{\hat{\theta}_{xy,r} - \hat{\theta}_{xy,r}^*}{h} \end{pmatrix} \xrightarrow{P} 0,$$

where $\tilde{R}(\theta) = (1/2n)(\partial^2 R_1(\theta; x, y)/\partial \theta \partial \theta^r)$, and

$$R^* = \text{diag}(1, h^{-1}, \dots, h^{-r}) \tilde{R}(\hat{\theta}_{xy}) \text{diag}(1, h^{-1}, \dots, h^{-r}).$$

It can be proved that $R^* \xrightarrow{P} f(x)g(y|x)\{1 - g(y|x)\}S$, where $S = (\kappa_{i+j-2})$ is a positive definite matrix. Therefore we have that

$$\hat{\theta}_{xy,i} = \hat{\theta}_{xy,i}^* + o_p(h^{r-i+1})$$

for $i = 0, 1, \dots, r$. Now (7.2) follows from the fact that $\hat{g}(y|x) = \exp(\hat{\theta}_{xy,0})$. We have completed the proof. ■