# MINIMALLY BIASED NONPARAMETRIC REGRESSION AND AUTOREGRESSION

Authors:     Timothy L. McMurry
             – Department of Mathematical Sciences, DePaul University,
               Chicago, IL, USA
               tmcmurry@depaul.edu

             Dimitris N. Politis
             – Department of Mathematics, University of California,
               San Diego, La Jolla, CA, USA
               politis@math.ucsd.edu

Abstract:

• A nonparametric regression estimator is introduced which adapts to the smoothness
  of the unknown function being estimated. This property allows the new estimator
  to automatically achieve minimal bias over a large class of locally smooth functions
  without changing the rate at which the variance converges. Optimal convergence rates
  are shown to hold for both i.i.d. data and autoregressive processes satisfying strong
  mixing conditions.

Key-Words:

• *nonparametric regression; autoregression; Fourier transform.*

AMS Subject Classification:

• 62G08.

## 1.  INTRODUCTION

Suppose the data $(X_1, Y_1), ..., (X_n, Y_n)$ are observations from a general real valued bivariate random process. The simplest example is when the data are generated by a model of the form $Y_i = r(X_i) + \epsilon_i$ where the $\epsilon_i$ are mean zero random errors satisfying some conditions; in general, the $\epsilon_i$ will not be independent. A second example of interest is nonparametric autoregression, where $Y_t = X_{t+1}$. The function $r$, the conditional mean of $Y$ given $X$, is unknown and will be estimated from the data. There are many nonparametric approaches to estimating $r$, including various kernel methods proposed by Nadaraya [15] and Watson [25], Gasser and Müller [7, 8], and local polynomial estimators, Fan [6]. In each of these techniques $r(x)$ is in essence estimated through weighted local averaging on the data near $x$. The smoothness of the function $r$ and properties of the weights used in this averaging determine the performance of the estimator. In this paper we propose a new class of kernels which allow the Nadaraya–Watson estimator to automatically achieve asymptotically optimal performance no matter how smooth $r$ happens to be.

The Nadaraya–Watson estimator is defined to be

$$(1.1) \qquad \hat{r}(x) := \frac{\sum_{i=1}^{n} Y_i\, K\big((X_i - x)/h\big)}{\sum_{i=1}^{n} K\big((X_i - x)/h\big)} \ .$$

The function $K(x)$ is the kernel; it is used to weight the observations. The denominator ensures the weights sum to 1. The parameter $h$ is the bandwidth, or smoothing parameter. It balances a tradeoff between bias and variance. Small values of $h$ concentrate the mass of the kernel near $x$, giving heavy weight to nearby observations and relatively little or no weight to more distant observations, resulting in a relatively unbiased but highly variable estimate. By contrast, large values of $h$ average over many data points, resulting in an estimate with relatively low variance, but potentially large bias, as observations which are quite distant from $x$ are included in the average. Since the number of data points included in the average is proportional to $nh$, each of these estimators has pointwise variance proportional to $1/(nh)$. For these reasons, we require that as $n \to \infty$, $h \to 0$ in such a way that $nh \to \infty$.

It is well known that the asymptotic bias of such nonparametric regression estimators is proportional to $h^p$, where $p$ depends on the smoothness of $r$, the smoothness of the marginal density of the $X_i$, and the properties of the kernel, or in the case of local polynomials, the polynomial degree of the local fit. In this paper we show that through appropriate choice of kernel, the rate at which the bias converges to zero will only be limited by properties of the unknown function, and not the kernel.

Sections 2 and 3 contain some important definitions and background. The case where the pairs of data $(X_1, Y_1), ..., (X_n, Y_n)$ are i.i.d. will be studied in Section 4; the case where the data satisfy strong mixing conditions will be presented in Section 5; a small simulation study is presented in Section 6. Technical proofs have been placed in Section 7.

## 2.    INFINITE ORDER KERNELS

If the kernel $K$ has finite moments up to order $q$ and its first $q-1$ moments are 0, then $K$ is said to be of *order q*. The most frequently used kernels are second order; common examples include the Epanechnikov kernel, $K_e(x) := (3/4)\,(1-x^2)\,1_{[-1,1]}(x)$, and the scaled normal density.

In general, if $r$ is $k$ times differentiable, with $k \geq 2$, the bias of a second order kernel estimate is $O(h^2)$. This rate of convergence can be improved up to $O(h^k)$ by choosing a kernel of order greater than or equal to $k$. However, the degree of smoothness in the underlying function is unknown and difficult to estimate, so it is difficult to know what order kernel to use.

In order to alleviate this difficulty we focus on a class of kernels that effectively have infinite order. These kernels automatically reduce the bias to $o(h^k)$ no matter how large $k$ happens to be. As in Politis and Romano [18, 19, 20] and Politis [16, 17], we now state the following general definition.

**Definition 2.1.**    A general flat-top kernel $K$ is defined in terms of its Fourier transform $\lambda$, which in turn is defined as follows. Fix a constant $c > 0$. Let

(2.1) $$\lambda(s) = \begin{cases} 1 & \text{if } |s| \leq c, \\ g\big(|s|\big) & \text{if } |s| > c, \end{cases}$$

where the function $g$ is chosen to make $\lambda(s)$, $\lambda^2(s)$, and $s\,\lambda(s)$ integrable. The flat top kernel is now given by

(2.2) $$K(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \lambda(s)\, e^{-isx}\, ds,$$

i.e., the inverse Fourier transform of $\lambda(s)$.

Note that in the preceding definition, the choice of $g$ is not unique. The function $\lambda$, and hence the kernel $K$, depend on the function $g$ and the parameter $c$ although this dependence will not be explicitly denoted.

Kernels satisfying this definition do not necessarily satisfy the moment conditions $\int z^k K(z)\, dz = 0$ for all integers $k$, as some of these integrals may not

be defined in either the Riemann or Lebesgue sense. However, the Cauchy principal value of each of these integrals is zero, and in many cases this is sufficient for optimal asymptotic performance.

The simplest kernel satisfying Definition 2.1 is determined by

$$\lambda_D(s) = \begin{cases} 1 & \text{if } |s| \leq 1, \\ 0 & \text{if } |s| > 1. \end{cases}$$

This is the example studied in the case of density estimation by Davis [2, 3], Devroye [4], and Ibragimov and Hasminksii [12], and it generates the Dirichlet kernel, $K(x) := \sin(x)/(\pi x)$. Both $\lambda(s)$ and the resulting kernel are shown in Figure 1. We can see that the tails of this kernel are very wiggly. This is problematic in two ways. First, the slow decay in the tails and the large negative oscillations increase $\int K^2(z)\,dz$, which will be shown to increase the variance of the estimate. Secondly, the large wiggles distant from 0 generate a finite sample bias because they allow observations which are relatively distant from $x$ to have a substantial influence on the estimate at $x$. These difficulties make density estimators using this kernel relatively uncompetitive for all but extremely large sample sizes.
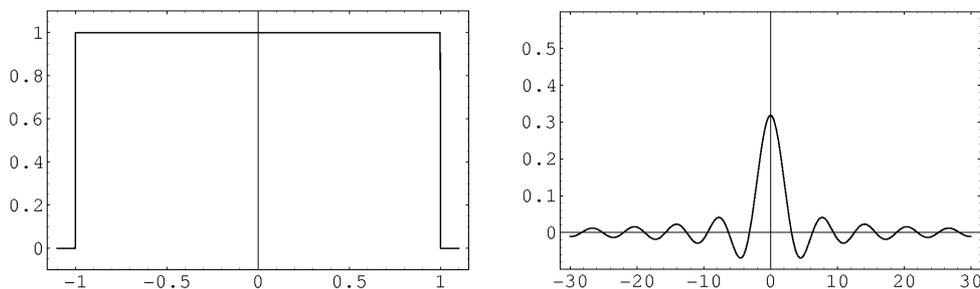


**Figure 1**: $\lambda(s)$ and the resulting Dirichlet kernel.

These problems can be substantially remedied by making the transition from 0 to 1 in the Fourier domain less abrupt. For example, Devroye and Gyorfi [5], Hall and Marron [11], and in the case of spectral density estimation, Politis and Romano [19], studied the kernel whose Fourier transform is given by

$$(2.3) \qquad \lambda_{T,1/2}(s) = \begin{cases} 1 & \text{if } |s| \leq 1/2, \\ 2\left(1 - |s|\right) & \text{if } 1/2 < s \leq 1, \\ 0 & \text{if } |s| > 1. \end{cases}$$

The corresponding kernel is

$$K(x) = \frac{2\left(\cos(x/2) - \cos(x)\right)}{\pi x^2}.$$

These are shown in Figure 2. Note the substantial improvement in the tails of the kernel.
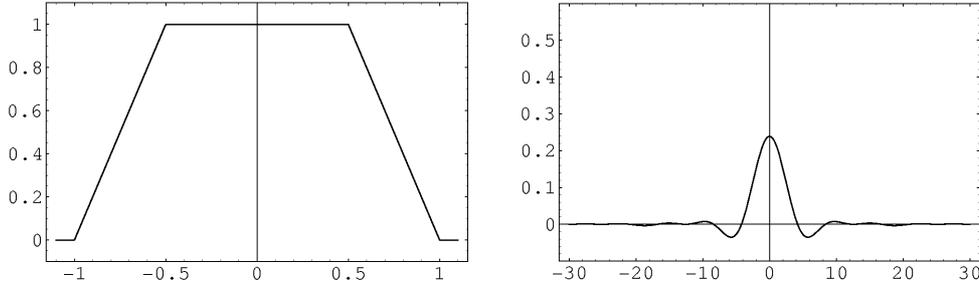


**Figure 2**:  $\lambda(s)$ and the resulting improved kernel.

Unfortunately, in the case of regression, it is often only reasonable to assume that the function being estimated is smooth over some interval rather than over its entire domain; if the marginal density of the $X_i$ has compact support, then the endpoints often generate discontinuities. Since infinite order kernels do not have compact support, the effects caused by these breaks get spread across the whole region of interest, potentially worsening the rate of convergence. For this reason, as discussed in the case of discontinuous density estimation in Politis [16], it is important that the tails of $K$ decay as quickly as possible, to minimize the effect on the interior of the interval. This can be ensured by requiring the Fourier transform of the kernel to be very smooth. If $\lambda$ is infinitely differentiable, then the tails of $K(x)$ decay faster than $x^{-m}$ for any positive $m$. In addition, $\lambda(s)$ as defined in equation (2.1) clearly has an infinite number of zero derivatives at $s = 0$. Together, these two conditions ensure that all moments of $K$ are zero in the Lebesgue sense. For these reasons, for the remainder of this work, we will restrict ourselves to kernels satisfying the following stronger definition.

**Definition 2.2.** An infinitely differentiable flat-top kernel $K$ is a flat-top kernel (as in Definition 2.1) with the added caveat that the function $g$ is chosen to make $\lambda(s)$ infinitely differentiable for all $s$.

We now provide an example of such a kernel, which was first introduced in McMurry and Politis [14], where the case of fixed design regression was studied. Let $b$ and $c$ be constants satisfying $b > 0$ and $0 < c < 1$. Define $\lambda(s)$ by

$$
(2.4) \quad \lambda_{IO}(s) = \begin{cases} 1 & \text{if } |s| \leq c, \\ \exp\left[-b \exp\left[-b/\big(|s| - c\big)^2\right]\big/\big(|s| - 1\big)^2\right] & \text{if } c < |s| < 1, \\ 0 & \text{if } |s| \geq 1. \end{cases}
$$

The parameter $c$ determines the region over which the kernel is identically 1; the parameter $b$ allows the shape of $\lambda$ to be altered, making the transition from 0 to 1 less abrupt. Figure 3 show plots of $\lambda$ (as defined above) and the resulting kernel $K$ for $c = 0.05$ and $b = 1/4$.
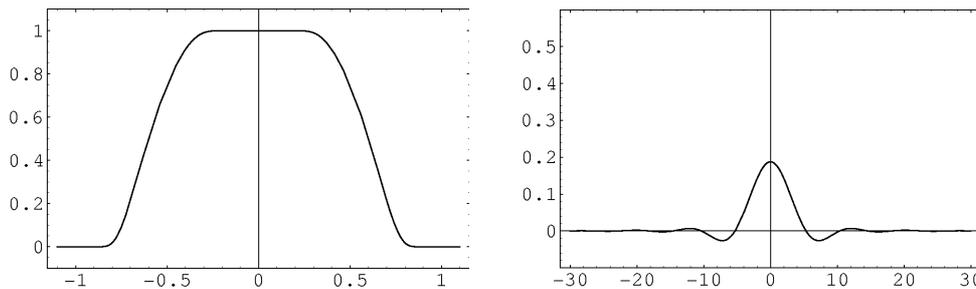


**Figure 3**:   Smooth $\lambda(s)$ and the resulting kernel with $b = 1/4$ and $c = 0.05$.

The function $\exp\!\big[-b\,\exp[-b/(|s| - c)^2]/(|s| - 1)^2\big]$ was chosen because it connects the regions where $\lambda$ is 0 and the region where $\lambda$ is 1 in a manner such that $\lambda(s)$ is infinitely differentiable for all $s$, including where $|s| = c$, and $|s| = 1$.

## 3.    BACKGROUND AND NOTATION

We examine the performance of the Nadaraya–Watson estimator when using infinite order kernels. The observed data is assumed to take the form of identically distributed pairs $(X_1, Y_1), ..., (X_n, Y_n)$, which satisfy $Y_i = r(X_i) + \epsilon_i$. Further restrictions will be necessary, but their discussion will be postponed for the moment. The Nadaraya–Watson estimator introduced in equation (1.1) can be written as

$$\hat{r}(x) := \frac{(1/n) \sum_{i=1}^{n} Y_i \, K_h(X_i - x)}{(1/n) \sum_{i=1}^{n} K_h(X_i - x)} \ ,$$

where

$$K_h(x) := (1/h) \, K(x/h) \ .$$

This estimator should be viewed in two different ways. As mentioned before, it is a weighted local average of the $Y_i$'s, where the denominator normalizes the weights so they sum to 1. It is also an explicit estimator of conditional expectation. The denominator is the standard kernel estimate of the design density, the marginal density of the $X_i$'s. The numerator is an approximation to $\int_{-\infty}^{\infty} y \, f(x, y) \, dy$, where $f(x, y)$ is the joint density of $(X_i, Y_i)$. Put together, this is an approximation to $r(x) = E\big[Y|X = x\big]$.

In order to simplify notation, define

$$\hat{g}(x) := (1/n) \sum_{i=1}^{n} Y_i \, K_h(X_i - x)$$

and

$$\hat{f}(x) := (1/n) \sum_{i=1}^{n} K_h(X_i - x) \, ,$$

which are the finite sample approximations to

$$g(x) := \int_{-\infty}^{\infty} y \, f(x, y) \, dy$$

and

$$f(x) := \int_{-\infty}^{\infty} f(x, y) \, dy \, .$$

We note that $r(x) = g(x)/f(x)$.

---

## 4.    THE INFINITE ORDER NADARAYA-WATSON ESTIMATOR FOR I.I.D. DATA

We first examine the behavior of the Nadaraya–Watson estimator when the observed pairs of data, $(X_1, Y_1), ..., (X_n, Y_n)$, are i.i.d. In order to understand the estimator as a whole, we begin with lemmas quantifying the asymptotic performance of the numerator and denominator, $\hat{g}(x)$ and $\hat{f}(x)$, as they approximate $g(x)$ and $f(x)$. In the process, it will be necessary to impose some assumptions, which will be introduced and discussed as needed. We first place some reasonable restrictions on the behavior of the bandwidth $h$ as the sample size grows large and on the conditional distribution of the errors.

**Assumption 1.**    As the sample size $n \to \infty$, the bandwidth $h \to 0$ in such a way that $nh \to \infty$.

**Assumption 2.**    $E\big[\epsilon_i | X_i = x\big] = 0$,  and  $E\big[\epsilon_i^2 | X_i = x\big] := \sigma^2(x) < \infty$.

Under this assumption, $\hat{f}(x)$ and $\hat{g}(x)$ are infinite order estimators of $f(x)$ and $g(x)$; this is quantified in the following lemma.

**Lemma 4.1.**    *If $x$ is contained in an open interval on which $f(x)$ has $p$ bounded continuous derivatives and $r(x)$ has $q$ bounded continuous derivatives, then under Assumptions 1 and 2,*

(4.1)                                              $E\big[\hat{f}(x)\big] - f(x) \, = \, o(h^p)$

*and*

(4.2) $$E\big[\hat{g}(x)\big] - g(x) = o(h^k) \, ,$$

*where $k = \min\{p, q\}$. If both $f(x)$ and $g(x)$ are infinitely differentiable, then each of these biases become $o(h^m)$ for all positive real $m$.*

      If we impose the additional assumptions that the observed pairs of data are i.i.d., and that $f$, $g$, and $\sigma^2(x)$ are reasonably well behaved, then the variance of $\hat{f}$ and $\hat{g}$ also behaves as expected.

      **Assumption 3.** $(X_1, Y_1), ..., (X_n, Y_n)$ are i.i.d.

      The next assumption is necessary to ensure that the asymptotic approximations we use are valid, and to avoid division by zero.

      **Assumption 4.** The point $x$ is a continuity point of $\sigma^2(x)$, $f(x) > C$ for some $C > 0$, and $r$ and $f$ are each differentiable in a neighborhood of $x$.

      **Lemma 4.2.** *Under Assumptions 1–4,*

(4.3) $$\operatorname{var}\big[\hat{f}(x)\big] = \frac{f(x)}{nh} \int_{-\infty}^{\infty} K^2(z)\, dz + o\left(\frac{1}{nh}\right) + O\left(\frac{1}{n}\right),$$

(4.4) $$\operatorname{var}\big[\hat{g}(x)\big] = \frac{\big(r^2(x) + \sigma^2(x)\big) f(x)}{nh} \int_{-\infty}^{\infty} K^2(z)\, dz + o\left(\frac{1}{nh}\right) + O\left(\frac{1}{n}\right)$$

*and*

(4.5) $$\operatorname{cov}\big[\hat{f}(x), \hat{g}(x)\big] = \frac{r(x)\, f(x)}{nh} \int_{-\infty}^{\infty} K^2(z)\, dz + o\left(\frac{1}{nh}\right) + O\left(\frac{1}{n}\right).$$

      Now that the behaviors of $\hat{f}$ and $\hat{g}$ are understood independently, the analysis will proceed by establishing their joint asymptotic normality. Once this has been shown, a Taylor series argument can be employed to show that $\hat{r}$ also has an asymptotic normal distribution with optimal bias and the standard variance. The joint asymptotic normality of $\hat{f}$ and $\hat{g}$ will be established via the Liapunov condition, which implies the Lindeberg–Feller central limit theorem. This requires a uniform bound on the $2 + \delta$'th moments of the $Y_i$, for some $\delta > 0$.

      **Assumption 5.** There exists a positive constants $M$ and $\delta$ such that

$$E\Big[|Y_i|^{2+\delta}\big|X_i = x\Big] < M \, ,$$

for all $x$.

The final assumption forces the conditional variance of the errors to be bounded above and below for all $x$. The bound from below is assumed for technical simplicity.

**Assumption 6.**   There exist strictly positive constants $b$ and $B$ such that $b < \sigma^2(x) < B$ for all $x$.

**Lemma 4.3.**   *Under Assumptions 1–6, for all real $c_1$ and $c_2$ (not both zero),*

$$(4.6) \quad \sqrt{nh}\left[c_1\left(\hat{f}(x) - E\big[\hat{f}(x)\big]\right) + c_2\left(\hat{g}(x) - E\big[\hat{g}(x)\big]\right)\right] \xrightarrow{\mathcal{D}} N\big(0, \theta(x)\big).$$

*where $\theta(x) := \left(c_1^2 + 2\,c_1\,c_2\,r(x) + c_2^2\big[r^2(x) + \sigma^2(x)\big]\right)f(x)\int_{-\infty}^{\infty}K^2(z)\,dz$. This implies the joint asymptotic normality of $\hat{f}$ and $\hat{g}$.*

The consequence of the preceding Lemma is the asymptotic normality of our estimator.

**Theorem 4.1.**   *If $x$ is contained in an open interval on which $f(x)$ has $p$ bounded continuous derivatives and $r(x)$ has $q$ bounded continuous derivatives, then under Assumptions 1–6,*

$$(4.7) \quad \sqrt{nh}\left(\hat{r}(x) - r(x) + o(h^k)\right) \xrightarrow{\mathcal{D}} N\left(0,\, \frac{\sigma^2(x)}{f(x)}\int_{-\infty}^{\infty}K^2(z)\,dz\right),$$

*where $k = \min\{p, q\}$.*

**Remark 4.1.**   Letting $h$ proportional to $n^{-1/(2k+1)}$, the mean square optimal rate, we get $\hat{r}(x) = r(x) + O_p\big(n^{-k/(2k+1)}\big)$, and $\sqrt{nh}\big(\hat{r}(x) - r(x)\big) \xrightarrow{\mathcal{D}} N\big(0, \frac{\sigma^2(x)}{f(x)}\int_{-\infty}^{\infty}K^2(z)\,dz\big)$, which demonstrates the higher order accuracy provided by infinite order kernels.

## 5.   DEPENDENT DATA AND NONPARAMETRIC AUTOREGRESSION

It is desirable to weaken the condition that $(X_1, Y_1), ..., (X_n, Y_n)$ are i.i.d. In particular we wish to be able to estimate nonparametric autoregression, where $X_t$ may be an unknown function of $X_{t-1}$. Mathematically, autoregressive processes are assumed to satisfy a model of the form, $X_t = r(X_{t-1}) + \sigma(X_{t-1})\,\epsilon_t$, where $r$ and $\sigma$ are unknown, and the $\epsilon_t$ are mean zero errors; further restrictions

similar to those in Section 4 will be imposed as necessary. The problem of interest is the estimation of $r(x) := E\big[X_t | X_{t-1} = x\big]$; this can be done by pairing consecutive observations, $(X_1, X_2), (X_2, X_3), ..., (X_{n-1}, X_n)$, and then performing a standard nonparametric regression.

Since nothing is gained by restricting ourselves to the case of autoregression, we will study the infinite order kernel estimator in the case where $(X_1, Y_1), ..., (X_n, Y_n)$ satisfy the same type of asymptotic dependence conditions that we wish the autoregressive process to satisfy. The results in this more general situation will then imply the desired result in the specific case of interest.

Any meaningful analysis will require conditions that ensure some sort of asymptotic independence; that is, random samples at times which are very distant from each other should behave as if they are independent. We will focus on the case of $\alpha$-mixing because it is the weakest of the most commonly studied conditions.

**Definition 5.1.** Let $\mathcal{F}_l^m$ be the $\sigma$-field generated by $U_l, U_{l+1}, ..., U_m$. A stationary time series $\{U_n\}_{n \in \mathbb{Z}}$ is said to be $\alpha$-mixing (or strong-mixing), if

$$\sup_{k \in \mathbb{Z}} \sup_{A \in \mathcal{F}_{-\infty}^k, B \in \mathcal{F}_{k+i}^{\infty}} \left| P(A)\, P(B) - P(AB) \right| := \alpha(i) \ \rightarrow \ 0 \ ,$$

as $i \to \infty$. The $\alpha(i)$'s are called the $\alpha$-mixing coefficients.

The analysis proceeds as in the i.i.d. case. We begin by proving, under some conditions, the joint asymptotic normality of $\hat{f}(x)$ and $\hat{g}(x)$. Once this has been established, we will be able to use the same argument used in the proof of Theorem 4.1 to show the asymptotic normality of $\hat{g}/\hat{f}$. Similar results for local polynomial regression and for Nadaraya–Watson estimators with finite order kernels have been obtained by Masry and Fan [13] and Robinson [21] respectively. Although their arguments are similar, the central limit theorem we prove here will be more closely related to that of Masry and Fan [13]. Let $c_1$ and $c_2$ be real numbers (not both zero). Define

$$Z_i := c_1 \Big[ K_h(X_i - x) - E\big[ K_h(X_i - x) \big] \Big] + c_2 \Big[ Y_i\, K_h(X_i - x) - E\big[ Y_i\, K_h(X_i - x) \big] \Big]$$

and

$$Q_n := \frac{1}{n} \sum_{i=1}^{n} Z_i \ .$$

We will establish asymptotic normality for $\sqrt{nh}\, Q_n$. In order to do so, we impose further assumptions on the marginal distributions of $(X_1, Y_1)$, and on the $\alpha$-mixing coefficients associated with the time series defined by $U_i = (X_i, Y_i)$.

**Assumption 7.** There exist finite positive bounds $M_1$, $M_2$, and $M_3$, such that

(i)  $f_i(u, v) \leq M_1$, where $f_i(u, v)$ is the joint density of $(X_1, X_i)$;

(ii) $E\big[Y_1^2 + Y_i^2 | X_1, X_i\big] \leq M_2$;

(iii) There exists $\delta > 2$ and $\beta > 1 - 2/\delta$ such that $E\big[|Y_1|^\delta | X\big] \leq M_3$ and $\sum_{i=1}^\infty i^\beta [\alpha(i)]^{1-2/\delta} < \infty$.

**Lemma 5.1.** *Let $x$ be a continuity point of conditional mean and variance functions, $r(\cdot)$ and $\sigma^2(\cdot)$. In addition, suppose that the marginal density of the $X_i$, $f(\cdot)$, and the product of $r(\cdot)$ and $f(\cdot)$, $g(\cdot)$, have $k$ bounded continuous derivatives in a neighborhood of $x$, where $k \geq 1$. Under Assumptions 2, 4, 6, and 7, we have the following convergences as $n \to \infty$, $h \to 0$, and $nh \to \infty$:*

(a) $h \operatorname{var}[Z_1] \to \theta(x)$,

(b) $h \sum_{i=1}^{n-1} \big| \operatorname{cov}[Z_1, Z_{i+1}] \big| \to 0$,

(c) $nh \operatorname{var}[Q_n] \to \theta(x)$,

*where, as before, $\theta(x) := \Big( c_1^2 + 2\, c_1 c_2\, r(x) + c_2^2 \big[ r^2(x) + \sigma^2(x) \big] \Big) f(x) \int_{-\infty}^\infty K^2(z)\, dz$.*

In order to establish the central limit theorem, we need one final condition on the $\alpha$-mixing coefficients.

**Assumption 8.** There exists a sequence of positive integers satisfying $s_n \to \infty$ and $s_n = o(\sqrt{nh})$ such that $\sqrt{n/h}\, \alpha(s_n) \to 0$.

**Remark 5.1.** Assumption 8 is a technical assumption that may dictate some particular rates at which $h \to 0$; nevertheless, Assumption 8 is weak enough to allow for a wide range of useful rates. To elaborate, Assumption 7 (iii) requires that the mixing coefficients decay at a polynomial rate depending on $\delta$. In particular, it requires that there exist $C > 0$, $\epsilon > 0$, and $n_0 > 0$ such that for all $n > n_0$, $\alpha(n) < C n^{-(\delta/(\delta-2)+1+\epsilon)}$. For example, if $\delta = 3$, then the mixing coefficients need to decay slightly faster than $n^{-4}$. Assuming that the function being estimated is at least twice differentiable, $h$ will optimally decrease at a rate equal to or slower than $n^{-1/5}$. This means that $\sqrt{n/h} \leq C n^{3/5}$ for some constant $C$. Similarly, $\sqrt{nh} \geq C n^{2/5}$. Under the strongest moment assumptions, $\alpha(n)$ is required to decay faster than $n^{-2}$. If we put these together, $\sqrt{n/h}\, \alpha(s_n) \leq C n^{3/5} s_n^{-2}$. From this expression it is easily seen that as long as $s_n$ grows faster than $n^{3/10}$, the second requirement of Assumption 8 will be satisfied. By the preceding argument the first requirement is satisfied if $s_n = o(n^{2/5})$. Since these conditions can be met simultaneously, Assumption 8 generally imposes no additional restrictions.

**Lemma 5.2.**  *Under Assumptions 2, 4, and 6–8, we have as $n \to \infty$, $h \to 0$, and $nh \to \infty$,*

$$\sqrt{nh}\, Q_n \xrightarrow{\mathcal{D}} N\big(0, \theta(x)\big).$$

The immediate consequence of this result is the following asymptotic normality for $\hat{r}$.

**Theorem 5.1.**  *If $x$ is contained in an open interval on which $f(x)$ has $p$ bounded continuous derivatives and $r(x)$ has $q$ bounded continuous derivatives, then under Assumptions 1, 2, 4, and 6–8,*

$$(5.1) \qquad \sqrt{nh}\left(\hat{r}(x) - r(x) + o(h^k)\right) \xrightarrow{\mathcal{D}} N\left(0, \frac{\sigma^2(x)}{f(x)} \int_{-\infty}^{\infty} K^2(z)\, dz\right),$$

*where  $k = \min\{p, q\}$.*

---

## 6.    SIMULATIONS

---

An extensive simulation study was undertaken to investigate the performance of the proposed estimator. For each combination of regression function, design density, error variance, and sample size, 100 data sets were created and smoothed using the infinite order and local linear estimators. Finally the integrated square error was estimated using Simpson's rule. Bandwidths for the infinite order estimator were selected by the rule of thumb suggested in [14] and developed further in [17]. Bandwidths for the local linear estimator were selected using the direct plug-in method suggested by Ruppert, Sheather, and Wand [22] and implemented in the R package KernSmooth [24].

The first regression function was taken to be $r(x) = x + 4\exp(-2x^2)/\sqrt{2\pi}$, which includes sections of almost linear behavior and an exponential bump with more curvature. Design densities were uniform on $[-2, 2]$ and $N(0, 1)$, and the integrated square error is over the interval $[-2, 2]$. The resulting integrated square errors for one simulation are shown in Figure 4. A scatterplot along with the two smoothings is shown in Figure 5.

The second regression function was taken to be $r(x) = \sin(4\pi x)$ with uniform design density on $[0, 1]$. The integrated square error was calculated on both the entire interval $[0, 1]$ and over the interval $[0.15, 0.85]$ to exclude edge effects.

It is clear from the simulations that the two estimators have different strengths. The infinite order estimator is clearly superior in the interior of the data set, when the error variance is large, and when the sample size is moderate to large. Since the local linear estimator automatically adapts to the edges of the design and the infinite order estimator does not, it is unsurprising that the local linear estimator is superior in these regions.

**Table 1**:   Comparison of infinite order and local linear estimators.

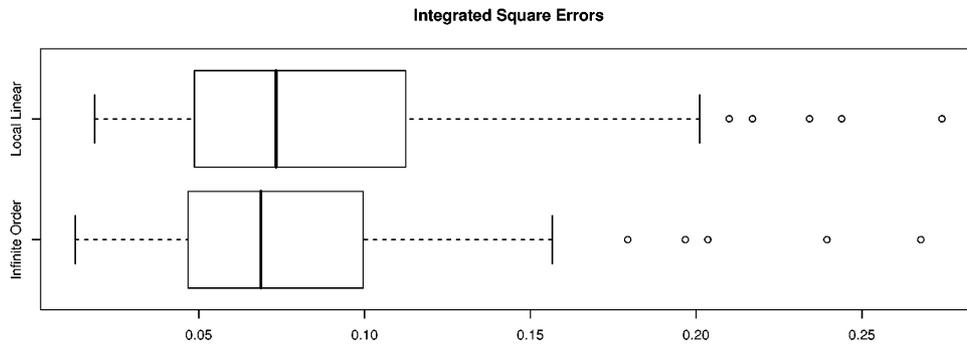| Function | Design | $n$ | $\sigma$ | Median Integrated Square Error | |
|---|---|---|---|---|---|
| | | | | Infinite Order | Local Linear |
| Exponential | Normal | 100 | 0.3 | 0.0870 | 0.0631 |
| | | | 0.5 | 0.1312 | 0.1627 |
| | | | 0.7 | 0.2101 | 0.2521 |
| Exponential | Normal | 200 | 0.3 | 0.0337 | 0.0324 |
| | | | 0.5 | 0.0633 | 0.0686 |
| | | | 0.7 | 0.1091 | 0.1375 |
| Exponential | Normal | 1000 | 0.3 | 0.0065 | 0.0077 |
| | | | 0.5 | 0.0132 | 0.0166 |
| | | | 0.7 | 0.0204 | 0.0274 |
| Exponential | Uniform | 100 | 0.3 | 0.0520 | 0.0384 |
| | | | 0.5 | 0.0954 | 0.1613 |
| | | | 0.7 | 0.0813 | 0.1481 |
| Exponential | Uniform | 200 | 0.3 | 0.0251 | 0.0190 |
| | | | 0.5 | 0.0481 | 0.0474 |
| | | | 0.7 | 0.0731 | 0.0823 |
| Exponential | Uniform | 1000 | 0.3 | 0.0066 | 0.0051 |
| | | | 0.5 | 0.0110 | 0.0112 |
| | | | 0.7 | 0.0175 | 0.0189 |
| Sin (edges excluded) | Uniform | 100 | 0.3 | 0.0065 | 0.0052 |
| | | | 0.5 | 0.0135 | 0.0139 |
| | | | 0.7 | 0.0232 | 0.0228 |
| Sin (edges included) | Uniform | 100 | 0.3 | 0.0120 | 0.0077 |
| | | | 0.5 | 0.0221 | 0.0191 |
| | | | 0.7 | 0.0412 | 0.0333 |
| Sin (edges excluded) | Uniform | 200 | 0.3 | 0.0032 | 0.0031 |
| | | | 0.5 | 0.0069 | 0.0068 |
| | | | 0.7 | 0.0108 | 0.0126 |
| Sin (edges included) | Uniform | 200 | 0.3 | 0.0066 | 0.0042 |
| | | | 0.5 | 0.0115 | 0.0091 |
| | | | 0.7 | 0.0191 | 0.0183 |
| Sin (edges excluded) | Uniform | 1000 | 0.3 | 0.0007 | 0.0008 |
| | | | 0.5 | 0.0013 | 0.0019 |
| | | | 0.7 | 0.0021 | 0.0030 |
| Sin (edges included) | Uniform | 1000 | 0.3 | 0.0029 | 0.0011 |
| | | | 0.5 | 0.0037 | 0.0026 |
| | | | 0.7 | 0.0050 | 0.0040 |

**Integrated Square Errors**



**Figure 4**:   Comparison of integrated square errors for 100 simulations
of the exponential function with $n = 200$ and $\sigma = 0.5$.
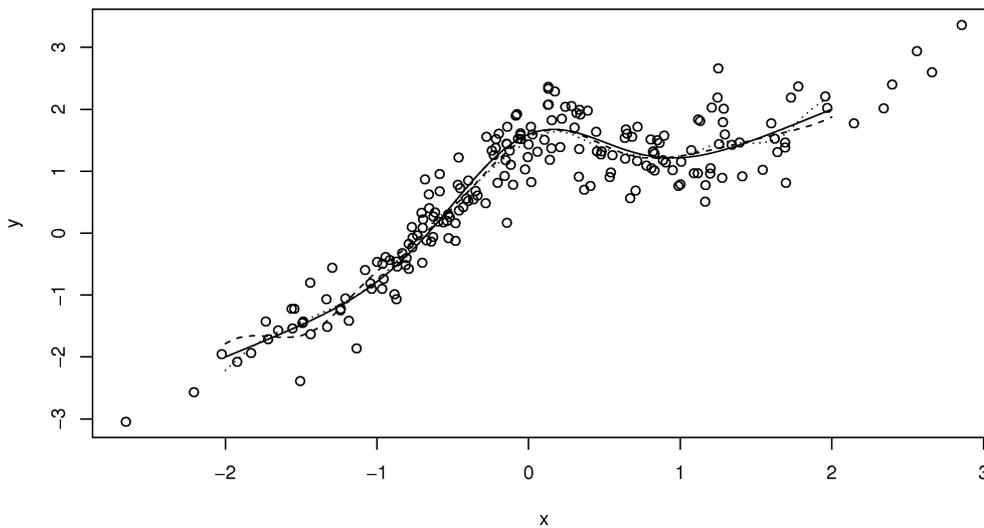


**Figure 5**:   A sample regression. The solid line is the true function,
the dashed line is the infinite order estimate, and the
dotted line is the local linear estimate.

## 7.    TECHNICAL PROOFS

**Proof of Lemma 4.1:**   The proof of (4.2) is almost identical to, but slight-
ly more complicated than the proof of (4.1). For this reason, only (4.2) will be
shown. The proof is similar in spirit to the proof of Theorem 2 in McMurry and
Politis [14], except this time it requires both $r$ and $f$ to be smooth. It will be
proved using a different, although more standard, technique.  This method of

proof is somewhat less elegant than the Fourier transform method used previously, but it has the advantage of showing that the same convergence rates hold even if $f$ and $g$ are only smooth on an open interval containing $x$. This proof technique does have a slight disadvantage. If $f$ and $g$ are smooth over all $\mathbb{R}$, then the Fourier transform technique can be employed to show the same convergence rates hold even if $\lambda(s)$, the Fourier transform of $K$, is not smooth. By conditioning on $X_i$,

$$
\begin{aligned}
E\big[\hat{g}(x)\big] - g(x) &= E\Big[E\big[Y_i\,K_h(X_i - x)\,\big|\,X_i\big]\Big] - g(x) \\
&= E\Big[r(X_i)\,K_h(X_i - x)\Big] - g(x) \\
&= \int_{-\infty}^{\infty} r(u)\,f(u)\,K_h(u - x)\,du \;-\; g(x)\ .
\end{aligned}
$$

Suppose $(rf)$ has $k$ bounded continuous derivatives on an interval $(a,b)$ containing $x$. Should $(rf)$ be smooth over all $\mathbb{R}$, then the proof can be simplified by taking $(a,b) = (-\infty, \infty)$:

$$
\begin{aligned}
E\big[\hat{g}(x)\big] - g(x) &= \int_{a}^{b} r(u)\,f(u)\,K_h(u - x)\,du \;+\; \int_{-\infty}^{a} r(u)\,f(u)\,K_h(u - x)\,du \\
&\quad + \int_{b}^{\infty} r(u)\,f(u)\,K_h(u - x)\,du \;-\; g(x)\ .
\end{aligned}
$$

Since the tails of $K(x)$ decay faster than $x^{-m}$ for all positive $m$, the two error terms are $o(h^m)$ for all positive $m$. At this point we perform a Taylor series expansion of the product $(rf)(z)$ around $x$:

$$
\begin{aligned}
E\big[\hat{g}(x)\big] - g(x) &= \\
&= \int_{a}^{b} r(u)\,f(u)\,K_h(u - x)\,du \;-\; g(x)\;+\; o(h^m) \\
&= \int_{(a-x)/h}^{(b-x)/h} r(x + hv)\,f(x + hv)\,K(v)\,dv \;-\; g(x)\;+\; o(h^m) \\
&= \int_{(a-x)/h}^{(b-x)/h} \left[(rf)(x) + hv\,(rf)'(x) + \cdots + \frac{(hv)^k}{k!}\,(rf)^{(k)}(x + \xi)\right] K(v)\,dv \\
&\quad - g(x)\;+\; o(h^m)\ ,
\end{aligned}
$$

where $\xi$ is between $x$ and $x + hv$. Since $K$ integrates to one, its moments are zero, and since $g(x) = r(x)\,f(x)$,

$$
\begin{aligned}
E\big[\hat{g}(x)\big] - g(x) &= \\
&= \int_{-\infty}^{\infty} \left[(rf)(x) + hv\,(rf)'(x) + \cdots + \frac{(hv)^{k-1}}{(k-1)!}\,(rf)^{(k-1)}(x)\right] K(v)\,dv \\
&\quad + \int_{(a-x)/h}^{(b-x)/h} \frac{(hv)^k}{k!}\,(rf)^{(k)}(x + \xi)\,K(v)\,dv \;-\; g(x)\;+\; o(h^m) \\
&= \int_{(a-x)/h}^{(b-x)/h} \frac{(hv)^k}{k!}\,(rf)^{(k)}(x + \xi)\,K(v)\,dv \;+\; o(h^m)\ .
\end{aligned}
$$

Since $(rf)$ has bounded continuous derivatives on $(a, b)$, we can apply the dominated convergence theorem, yielding

$$\lim_{h \to 0} \int_{(a-x)/h}^{(b-x)/h} (rf)^{(k)}(x+\xi)\, v^k\, K(v)\, dv \;=\; \int_{-\infty}^{\infty} (rf)^{(k)}(x)\, v^k\, K(v)\, dv \;=\; 0\ .$$

Therefore,

$$E\big[\hat{g}(x)\big] - g(x) \;=\; o(h^k)\ . \qquad \square$$

**Proof of Lemma 4.2:**  As $f(x)$ can be viewed as a special case of $g(x)$ with $Y_i = 1$ for all $i$, (4.3) will follow immediately from the proof of (4.4). The proof of (4.5) is almost identical to the proof of (4.4), so it is omitted.

By conditioning on $X_1$, and by Lemma 4.1,

$$
\begin{aligned}
\text{var}\big[\hat{g}(x)\big] \;&=\; \frac{1}{n}\, \text{var}\big[Y_1\, K_h(X_1 - x)\big] \\
&=\; \frac{1}{n}\left[\int_{-\infty}^{\infty}\Big(r^2(u) + \sigma^2(u)\Big) f(u)\, K_h^2(u - x)\, du \;-\; r^2(x) \;+\; o(h^k)\right] \\
&=\; \frac{1}{nh}\int_{-\infty}^{\infty}\Big(r^2(x+hz) + \sigma^2(x+hz)\Big) f(x+hz)\, K^2(z)\, dz \;+\; O\Big(\frac{1}{n}\Big)\ ,
\end{aligned}
$$

since $x$ is a continuity point of $r$, $f$, and $\sigma^2$, the dominated convergence theorem yields

$$\text{var}\big[\hat{g}(x)\big] \;=\; \frac{\Big(r^2(x) + \sigma^2(x)\Big) f(x)}{nh}\int_{-\infty}^{\infty} K^2(z)\, dz \;+\; o\Big(\frac{1}{nh}\Big) \;+\; O\Big(\frac{1}{n}\Big)\ . \qquad \square$$

**Proof of Lemma 4.3:**  The proof proceeds by verifying that the Liapunov condition holds, which is sufficient for the Lindeberg–Feller central limit theorem. The result is trivial if $c_1 = c_2 = 0$, so assume that at least one of these constants is nonzero. Let $C$ denote a positive constant.

$$
\frac{\sum_{i=1}^{n} E\Big[\big|(c_1 + c_2 Y_i)\, K_h(X_i - x)\big|^{2+\delta}\Big]}{\text{var}\Big[c_1 \sum_{i=1}^{n} K_h(X_i - x) + c_2 \sum_{i=1}^{n} Y_i\, K_h(X_i - x)\Big]^{(2+\delta)/2}} \;\leq\;
$$

$$
\leq\; \frac{C\, n\, h^{1+\delta}}{n^{1+\delta/2}\Big[\theta(x) + o(1/h) + O(1)\Big]^{(2+\delta)/2}}\ ,
$$

where the inequality follows from the proof of Lemma 4.2. After multiplying the numerator and denominator by $h^{(2+\delta)/2}$, it is clear that this quantity goes to zero as $n$ goes to infinity. Therefore, the Liapunov condition is satisfied, and the Lemma follows immediately. $\qquad \square$

**Proof of Theorem 4.1:** We begin with a lemma which ensures that for large enough $n$, $\hat{f}(x) \geq c > 0$ for some constant $c$, as long as $f(x) > 0$.

**Lemma 7.1.** *Suppose $f(x) > c$ for some $c > 0$. Also suppose $h$ decreases slowly enough that $n^{-2/7} + \delta = o(h)$ for some $\delta > 0$. Then for all $\epsilon > 0$,*
$P\big[|f(x) - \hat{f}(x)| > \epsilon \ i.o.\big] = 0.$

It should also be noted that much stronger results hold. Under additional conditions, rates of uniform almost sure convergence over compact sets can be established. See Bosq [1] or Györfi *et al.* [9].

**Proof of Lemma 7.1:** We make use of the following Bernstein inequality, which is Theorem 1.3 part (2) in Bosq [1].

**Lemma 7.2.** *Let $\{W_t\}_{t \in \mathbb{Z}}$ be a mean zero real valued random process such that $\sup_{1 \leq t \leq n} \|W_t\|_\infty \leq b$, and let $S_n = \sum_{t=1}^n W_t$. Then for each integer $q \in \lfloor 1, n/2 \rfloor$ and each $\epsilon > 0$,*

$$\mathrm{P}\big[|S_n| > n\epsilon\big] \ \leq \ 4 \exp\left(-\frac{\epsilon^2}{8\, v^2(q)}\, q\right) + 22\left(1 + \frac{4\,b}{\epsilon}\right)^{1/2} q\, \alpha\left(\left\lfloor \frac{n}{2\,q} \right\rfloor\right),$$

*where $v^2(q) = \frac{2}{p^2}\, \sigma^2(q) + \frac{b\epsilon}{2}$, $p = \frac{n}{2q}$, and*

$$\sigma^2(q) \ = \ \max_{0 \leq j \leq 2q-1} E\bigg[\Big(\lfloor jp \rfloor + 1 - jp\Big) X_{\lfloor jp \rfloor + 1} + X_{\lfloor jp \rfloor + 2} + \cdots$$
$$+ X_{\lfloor (j+1)p \rfloor} + \Big[(j+1)p - \lfloor (j+1)p \rfloor\Big] X_{\lfloor (j+1)p+1 \rfloor}\bigg].$$

The Borel–Cantelli lemma will be used to show

$$\mathrm{P}\Big[\big|\hat{f}_n(x) - E\,\hat{f}_n(x)\big| > \epsilon \ \text{i.o.}\Big] \ = \ 0 \,.$$

Since $E\,\hat{f}_n(x) \to f(x)$, this will establish the desired result.

Let $W_i = K_h(X_i - x) - E\big[K_h(X_i - x)\big]$. Then $\|W_i\|_\infty \leq \hat{K}/h$ for all $i$, where $\hat{K} = 2 \sup_{x \in \mathbb{R}} K(x)$. In addition, it can easily be seen that $\sigma^2(q) \leq \frac{(p+1)^2}{n^2 h^2}\, \hat{K}^2$.

Therefore,

$$p_n \ := \ \mathrm{P}\left[\frac{1}{n}\left|\sum_{i=1}^n W_i\right| > \epsilon\right]$$

$$\leq \ 4 \exp\left(-\frac{C_1\, \epsilon^2}{\frac{(p+1)^2}{n^2\, h^2\, p^2} + \frac{1}{h}}\, q\right) + C_2\left(1 + \frac{C_3}{h\epsilon}\right)^{1/2} q\, \alpha\left(\left\lfloor \frac{n}{2\,q} \right\rfloor\right) \,.$$

We need $\sum_n p_n < \infty$. In order for the first term in the sum to be fininte, it is necessary that $n^2 h^2 q/(n^2 h) \to \infty$ at a rate equal to or faster than $n^{\delta_1}$ for some $\delta_1 > 0$; this requires $q$ grow at least as fast as $n^{\delta_1} h^{-1}$. On the other hand, the second term requires $\sum_n (q/\sqrt{h}) \, \alpha(\lfloor n/2q \rfloor) < \infty$. As noted in the Remark 5.1, it is sufficient to choose $q$ such that $\sum_n (q/\sqrt{h})(q/n)^2 < \infty$. The latter condition can be satisfied if $q^3 h^{-1/2}$ grows at a rate $n^{1-\delta_2}$, for some $\delta_2 > 0$. Equivalently, it suffices for $q$ to grow at a rate equal to or slower than $n^{(1/3)-\delta_2} h^{1/6}$. It can easily be seen that these requirements can be met simultaneously as long as $h$ satisfies $n^{-2/7+\delta} = o(h)$, for some $\delta > 0$; this includes all optimal rates. This result could be further strengthened by imposing additional assumptions. For example, in the case where $X_1, ..., X_n$ are i.i.d., the mixing coefficients are 0, and hence summable. In this situation, the only restriction on $h$ is that it decrease slightly slower than $1/n$. In the case of a mixing process, the possible range of rates for $h$ could be expanded if one were to assume that the joint density of $(X_1, X_i)$ is differentiable with partial derivatives uniformly bounded in $i$. $\qquad\square$

We now return to the proof of the main result. By the preceding lemma, for large enough $n$, we can assume that $\hat{f}_n(x) > c/2$. So, we can apply the intermediate value theorem to see,

$$
\hat{r}(x) - r(x) \;=\; \frac{\hat{g}(x)}{\hat{f}(x)} - \frac{g(x)}{f(x)}
$$

$$
\;=\; \hat{g}(x)\left(\frac{1}{f(x)} - \frac{1}{\xi_n^2}\Big(\hat{f}(x) - f(x)\Big)\right) - \frac{g(x)}{f(x)}\,,
$$

where $|\xi_n - f(x)| \le |\hat{f}(x) - f(x)|$. This can be further simplified to

$$
\hat{r}(x) - r(x) \;=\; \frac{1}{f(x)}\Big(\hat{g}(x) - g(x)\Big) - \frac{\hat{g}(x)}{\xi_n^2}\Big(\hat{f}(x) - f(x)\Big)
$$

$$
\;=\; \frac{1}{f(x)}\Big(\hat{g}(x) - E\big[\hat{g}(x)\big]\Big) - \frac{\hat{g}(x)}{\xi_n^2}\Big(\hat{f}(x) - E\big[\hat{f}(x)\big]\Big) + o(h^k)\,.
$$

By Lemmas 4.1 and 4.2, $\hat{g}(x)$ and $\hat{f}(x)$ converge in probability to $g(x)$ and $f(x)$ respectively. Therefore, $\xi_n$ also converges to $f(x)$ in probability. By Slutsky's theorem, and Lemma 4.3,

$$
\sqrt{nh}\,\Big(\hat{r}(x) - r(x) + o(h^k)\Big) \;\xrightarrow{\;\mathcal{D}\;}\; N\!\left(0,\; \frac{\sigma^2(x)}{f(x)}\int_{-\infty}^{\infty} K^2(z)\,dz\right),
$$

the desired result. $\qquad\square$

**Proof of Lemma 5.1:** The proof of part (a) follows from similar results for the i.i.d. case:

$$
\mathrm{var}\big[Z_1\big] \;=\; E\Big[(c_1 + c_2 Y_1)^2 \, K_h^2(X_1 - x)\Big] - \Big(c_1 f(x) + c_2\, r(x)\, f(x) + O(h^k)\Big)^2
$$

$$
\;=\; E\Big[\Big(c_1^2 + 2\,c_1 c_2\, r(X_1) + c_2^2\big[r^2(X_1) + \sigma^2(X_1)\big]\Big) K_h^2(X_1 - x)\Big] + O(1)
$$

$$
\;=\; \theta(x)/h + o(1/h) + O(1)\,.
$$

The proof of part (b) is more challenging. Let $d_n$ be a sequence of integers such that $d_n \to \infty$ and $d_n h \to 0$. Define

$$J_1 := \sum_{i=1}^{d_n-1} \left| \mathrm{cov}[Z_1, Z_{i+1}] \right|,$$

and

$$J_2 := \sum_{i=d_n}^{n-1} \left| \mathrm{cov}[Z_1, Z_{i+1}] \right|.$$

We wish to show $J_1 = o(1/h)$ and $J_2 = o(1/h)$. We begin with $J_1$. By conditioning on $(X_1, X_i)$,

$$\left| \mathrm{cov}[Z_1, Z_i] \right| \le$$

$$\le \left| E\Big[ (c_1 + c_2 Y_1)\, K_h(X_1 - x)\,(c_1 + c_2 Y_i)\, K_h(X_i - x) \Big] \right| + O(1)$$

$$\le E\Big[ \left| K_h(X_1 - x)\, K_h(X_i - x) \right| E\Big[ \left| (c_1 + c_2 Y_1)(c_1 + c_2 Y_i) \right| \Big| X_1, X_i \Big] \Big] + O(1)$$

$$\le E\Big[ \left| K_h(X_1 - x)\, K_h(X_i - x) \right|$$
$$\times \left( E\Big[ (c_1 + c_2 Y_1)^2 \big| X_1, X_i \Big] E\Big[ (c_1 + c_2 Y_i)^2 \big| X_1, X_i \Big] \right)^{1/2} \Big] + O(1)$$

$$\le C\, E\Big[ \left| K_h(X_1 - x)\, K_h(X_i - x) \right| \Big] + O(1)$$

$$\le C \left( \int_{-\infty}^{\infty} \left| K_h(u - x) \right| du \right)^2 + O(1).$$

Since the $O(1)$ term is the same for all $i$, $\left| \mathrm{cov}[Z_1, Z_i] \right| < C$ for some positive $C$. Therefore $J_1 = o(1/h)$. For the second term, $J_2$, we employ Davydov's Lemma (see Hall and Heyde [10]), which tells us

$$\left| \mathrm{cov}[Z_1, Z_{i+1}] \right| \le 8 \left[ \alpha(i) \right]^{1-2/\delta} \left[ E|Z_1|^\delta \right]^{2/\delta}.$$

We now need to put a bound on $E|Z_i|^\delta$.

(7.1)
$$E|Z_i|^\delta = E\left[ \left| c_1 \Big[ K_h(X_i - x) - E\big[ K_h(X_i - x) \big] \Big] \right. \right.$$
$$\left. \left. + c_2 \Big[ Y_i\, K_h(X_i - x) - E\big[ Y_i\, K_h(X_i - x) \big] \Big] \right|^\delta \right]$$

$$\le 2\, E\left[ \left| c_1 \Big[ K_h(X_i - x) - E\big[ K_h(X_i - x) \big] \Big] \right|^\delta \right]$$
$$+ 2\, E\left[ \left| c_2 \Big[ Y_i\, K_h(X_i - x) - E\big[ Y_i\, K_h(X_i - x) \big] \Big] \right|^\delta \right].$$

These two terms behave similarly, so it is sufficient to examine only the second. Let $C$ denote a generic positive constant which may take on different values:

$$2\,E\left[\left|\,c_2\Big[Y_i\,K_h(X_i-x)-E\big[Y_i\,K_h(X_i-x)\big]\Big]\,\right|^{\delta}\right] \le$$

$$\le\ 4\,E\left[\big|\,c_2\,Y_i\,K_h(X_i-x)\,\big|^{\delta}\right]+4\,\big|E\big[Y_i\,K_h(X_i-x)\big]\big|^{\delta}$$

$$\le\ C\,E\left[\big|K_h(X_i-x)\big|^{\delta}\,E\big[|Y_i|^{\delta}\,\big|\,X_i\big]\right]+C$$

$$\le\ C\,h^{1-\delta}+C\ .$$

An identical result holds for the first term in equation (7.1). Putting these two terms together yields

$$\big[E\,|Z_1|^{\delta}\big]^{2/\delta}\ \le\ \big[C\,h^{1-\delta}+C\big]^{2/\delta}$$

$$\le\ C\,h^{2/\delta-2}+C\ .$$

Returning to $J_2$,

$$J_2\ \le\ \sum_{i=d_n}^{\infty}8\,\big[\alpha(i)\big]^{1-2/\delta}\,\big[E\,|Z_1|^{\delta}\big]^{2/\delta}$$

$$\le\ \sum_{i=d_n}^{\infty}\big[\alpha(i)\big]^{1-2/\delta}\big(C\,h^{2/\delta-2}+C\big)\ .$$

By Assumption 7, $\sum_{i=d_n}^{\infty}[\alpha(i)]^{1-2/\delta}\to 0$ as $d_n\to\infty$. So,

$$J_2\ \le\ C\,h^{2/\delta-2}\sum_{i=d_n}^{\infty}\big[\alpha(i)\big]^{1-2/\delta}+o(1)$$

$$\le\ C\,h^{2/\delta-2}\,d_n^{-\beta}\sum_{i=d_n}^{\infty}i^{\beta}\big[\alpha(i)\big]^{1-2/\delta}+o(1)\ .$$

By choosing $d_n$ such that $h^{2/\delta-1}d_n^{-\beta}\to 1$, we see $J_2=o(1/h)$ and $hd_n\to 0$, which ensures the convergence of $J_1$. This completes the proof of (b). The proof of (c) is an immediate consequence of (a) and (b). $\qquad\square$

**Proof of Lemma 5.2:** The proof employs a small-block large-block argument. The set $\{1,...,n\}$ is partitioned into $2k+1$ alternating large and small subsets. Let $r_n$ be the size of the large blocks and $s_n$ be the size of the small blocks. Then $k_n=\lfloor n/(r_n+s_n)\rfloor$. For $0\le j\le k-1$, define

$$U_j\ :=\ \sqrt{h}\sum_{i=j(r+s)+1}^{j(r+s)+r}Z_i\ ,$$

$$V_j\ :=\ \sqrt{h}\sum_{i=j(r+s)+r+1}^{(j+1)(r+s)}Z_i$$

and

$$W_j := \sqrt{h} \sum_{i=k(r+s)+1}^{n} Z_i .$$

We see immediately that $U_j$ sums the $Z_i$ over the blocks of size $r$, $V_j$ sums the $Z_i$ over the blocks of size $s$, and $W_j$ accounts for the remaining terms that do not fit evenly into the first $2k$ blocks.

The idea of the proof is to show that the small blocks separate the large blocks by enough to make them asymptotically independent while being small enough that they don't make a substantial contribution to the limiting distribution. The Lindeberg condition can then be checked for the separated large blocks.

To formalize this, we write

$$\sqrt{nh}\, Q_n = \frac{1}{\sqrt{n}} \left[ \sum_{j=0}^{k-1} U_j + \sum_{j=0}^{k-1} V_j + W_j \right]$$

$$:= \frac{1}{\sqrt{n}} \left[ Q'_n + Q''_n + Q'''_n \right] .$$

We will establish the following identities:

$$(7.2) \qquad \frac{1}{n} E\big[(Q''_n)^2\big] \to 0 ,$$

$$(7.3) \qquad \frac{1}{n} E\big[(Q'''_n)^2\big] \to 0 ,$$

$$(7.4) \qquad \left| E \exp(it\, Q'_n) - \prod_{j=1}^{k} E\big[\exp(it\, U_j)\big] \right| \to 0 ,$$

$$(7.5) \qquad \frac{1}{n} \sum_{i=1}^{k} E\big[U_i^2\big] \to \theta^2(x) ,$$

and

$$(7.6) \qquad \frac{1}{n} \sum_{j=0}^{k-1} E\Big[U_j^2\, 1_{[|U_j| \geq \epsilon\, \theta(x)\sqrt{n}]}\Big] \to 0 ,$$

for all $\epsilon > 0$.

Once these have been established, by a Taylor Series expansion, we will have

$$\left| E\Big[\exp\big(it\,\sqrt{nh}\, Q_n\big)\Big] - \exp\big(-t^2\theta^2(x)/2\big) \right| =$$

$$(7.7) \qquad = \left| E\left[\exp\Big(it\, \frac{1}{\sqrt{n}}\, \big[Q'_n + Q''_n + Q'''_n\big]\Big)\right] - \exp\big(-t^2\theta^2(x)/2\big) \right|$$

$$\leq \left| E\left[\exp\Big(it\, \frac{1}{\sqrt{n}}\, Q'_n\Big)\right] - \exp\big(-t^2\theta^2(x)/2\big) \right|$$

$$+ E\left[\left|\frac{2t}{\sqrt{n}}\, Q''_n\right|\right] + E\left[\left|\frac{2t}{\sqrt{n}}\, Q'''_n\right|\right] .$$

The final two terms will converge to 0 by the Cauchy–Schwarz inequality and equations (7.2) and (7.3). Equations (7.4), (7.5), and (7.6) are enough to verify the conditions of the Lindeberg–Feller central limit theorem, which will establish the desired result. The proof will be complicated somewhat because (7.6) will be established first for bounded random variables, and then the bound will be allowed to tend to infinity.

We begin by choosing block sizes. By Assumption 8, there exists a sequence $q_n$ such that $q_n \to \infty$ and $q_n s_n = o(\sqrt{nh})$, and $q_n \sqrt{n/h}\, \alpha(s_n) \to 0$. Define the large block size $r_n$ by

$$r_n := \lfloor \sqrt{nh}/q_n \rfloor \ .$$

From this definition, we see

$$(7.8) \qquad \frac{s_n}{r_n} \leq \frac{s_n}{\sqrt{nh}/q_n - 1} = \frac{q_n s_n/\sqrt{nh}}{1 - q_n/\sqrt{nh}} \ \to \ 0 \ ,$$

$$(7.9) \qquad \begin{aligned} \frac{n}{r_n}\, \alpha(s_n) &= \frac{n}{\lfloor \sqrt{nh}/q_n \rfloor}\, \alpha(s_n) \leq \frac{n}{\sqrt{nh}/q_n - 1}\, \alpha(s_n) \\ &= \left( \sqrt{n/h}\, q_n + o(1) \right) \alpha(s_n) \ \to \ 0 \end{aligned}$$

and

$$(7.10) \qquad \frac{r_n}{\sqrt{nh}} = \frac{\lfloor \sqrt{nh}/q_n \rfloor}{\sqrt{nh}} \leq \frac{\sqrt{nh}/q_n + 1}{\sqrt{nh}} = \frac{1}{q_n} + \frac{1}{\sqrt{nh}} \ \to \ 0 \ .$$

We begin by verifying equations (7.2) and (7.3):

$$(7.11) \qquad E\big[(Q_n'')^2\big] = \sum_{i=0}^{k-1} \mathrm{var}[V_j] + \sum_{i \neq j} \mathrm{cov}[V_i, V_j] \ ,$$

where, by Lemma 5.1,

$$\begin{aligned} \mathrm{var}[V_i] &= sh\, \mathrm{var}[Z_1] + 2\, sh \sum_{i=1}^{s-1} (1 - j/s)\, \mathrm{cov}[Z_1, Z_{1+i}] \\ &= s\big[\theta(x) + o(1)\big] \ . \end{aligned}$$

By equation (7.8),

$$\begin{aligned} \sum_{i=0}^{k-1} \mathrm{var}[V_j] &\leq k_n s_n \big[\theta(x) + o(1)\big] \\ &\leq \frac{n\, s_n}{s_n + r_n} \big[\theta(x) + o(1)\big] \\ &= o(n) \ . \end{aligned}$$

The second term of (7.11) can be treated as follows:

$$\sum_{i \neq j} \mathrm{cov}[V_i, V_j] = h \sum_{i \neq j}^{k-1} \sum_{l=1}^{s} \sum_{m=1}^{s} \mathrm{cov}\big[ Z_{i(r+s)+r+m},\, Z_{j(r+s)+r+l} \big] \ .$$

Since $i \neq j$, the difference between the indices, $\left|i(r+s)+r+m-(j(r+s)+r+l)\right| \geq r$, so

$$\left|\sum_{i \neq j} \mathrm{cov}[V_i, V_j]\right| \leq 2\,h \sum_{l=1}^{n-r} \sum_{m=l+r}^{n} \left|\mathrm{cov}[Z_l, Z_m]\right|$$

$$= 2\,h \sum_{l=1}^{n-r} \sum_{j=r}^{n-l} \left|\mathrm{cov}[Z_l, Z_{l+j}]\right|$$

$$\leq 2\,n\,h \sum_{j=r}^{n-1} \left|\mathrm{cov}[Z_1, Z_{j+1}]\right| = o(n)\,.$$

This establishes (7.2). We now turn our attention to (7.3):

$$\frac{1}{n}\,E[(Q_n''')^2] = \frac{h}{n}\Big(n - k(r+s)\Big)\,\mathrm{var}[Z_1] + 2\,h \sum_{i=2}^{n-k(r+s)} \mathrm{cov}[Z_1, Z_i]$$

$$\leq \frac{r_n + s_n}{n}\,\theta^2(x) + o(1) \;\to\; 0\,.$$

To prove (7.4), we use a lemma of Volkonskii and Rozanov [23], which is stated in Lemma 7.3, following this proof:

$$\left|E \exp(it\,Q_n') - \prod_{j=1}^{k} E\big[\exp(it\,U_j)\big]\right| \leq 16\,(k-1)\,\alpha(s_n+1)$$

$$= 16\,\frac{n}{r_n}\,\alpha(s_n+1) + o(1) \;\to\; 0\,,$$

by equation (7.9).

We now turn our attention to equation (7.5):

$$\frac{1}{n} \sum_{i=1}^{k} E[U_i^2] = \frac{k_n}{n}\,\mathrm{var}[U_1]$$

$$= \frac{k_n r_n}{n}\,\big(\theta^2(x) + o(1)\big)$$

$$= \frac{r_n}{r_n + s_n}\,\theta^2(x) + o(1) \;\to\; \theta^2(x)\,.$$

Finally, we verify equation (7.6). We begin establishing the result for truncated random variables, and then subsequently letting the truncation point go to infinity. Define

$$Z_i^L := (c_1 + c_2 Y_i)\,\mathbf{1}_{[|Y_i| \leq L]}\,K_h(X_i - x) - E\Big[(c_1 + c_2 Y_i)\,\mathbf{1}_{[|Y_i| \leq L]}\,K_h(X_i - x)\Big]\,,$$

$$Q_n^L := \frac{1}{n} \sum_{i=1}^{n} Z_i^L\,,$$

$$\tilde{Q}_n^L := \frac{1}{n} \sum_{i=1}^{n} (Z_i - Z_i^L)\,,$$

and

$$U_j^L := \sqrt{h} \sum_{i=j(r+s)+1}^{j(r+s)+r} Z_i^L.$$

We first need to estimate the asymptotic variance of $Z_1^L$. Assume conditions strong enough that for all $c_1$ and $c_2$, and for all $L > L_0$,

$$E\left[(c_1 + c_2 Y_i)^2 \, 1_{[|Y_i| \leq L]} \,\middle|\, X_1 = x\right]$$

is continuous as a function of $x$:

$$\left| E\left[(c_1 + c_2 Y_i) \, 1_{[|Y_i| \leq L]} \, K_h(X_i - x)\right] \right| \leq \left(|c_1| + |c_2 L|\right) E\left[\left|K_h(X_1 - x)\right|\right] \leq C \, ,$$

and

$$h \, E\left[(c_1 + c_2 Y_i)^2 \, 1_{[|Y_i| \leq L]} \, K_h^2(X_i - x)\right] =$$

$$= h \, E\left[K_h^2(X_1 - x) \, E\left[(c_1 + c_2 Y_i)^2 \, 1_{[|Y_i| \leq L]} \,\middle|\, X_1\right]\right]$$

$$= h \int_{-\infty}^{\infty} K_h^2(u - x) \, E\left[(c_1 + c_2 Y_i)^2 \, 1_{[|Y_i| \leq L]} \,\middle|\, X_1 = u\right] f(u) \, du$$

$$= \int_{-\infty}^{\infty} K^2(v) \, E\left[(c_1 + c_2 Y_i)^2 \, 1_{[|Y_i| \leq L]} \,\middle|\, X_1 = x + hv\right] f(x + hv) \, dv$$

$$= E\left[(c_1 + c_2 Y_i)^2 \, 1_{[|Y_i| \leq L]} \,\middle|\, X_1 = x\right] f(x) \int_{-\infty}^{\infty} K^2(v) \, dv \, + \, o(1) \, .$$

Putting these two together,

$$h \, \mathrm{var}[Z_1^L] = E\left[(c_1 + c_2 Y_i)^2 \, 1_{[|Y_i| \leq L]} \,\middle|\, X_1 = x\right] f(x) \int_{-\infty}^{\infty} K^2(v) \, dv \, + \, o(1) \, .$$

For the sake of notational simplicity, we now define

$$(\theta^L)^2(x) := E\left[(c_1 + c_2 Y_i)^2 \, 1_{[|Y_i| \leq L]} \,\middle|\, X_1 = x\right] f(x) \int_{-\infty}^{\infty} K^2(v) \, dv \, .$$

Returning now to the proof of (7.6), since $K$ and $Y_i^L$ are bounded, $hZ_i^L$ is also bounded. Equivalently, for some $D$,

$$\sqrt{h} \, Z_i^L \leq D/\sqrt{h} \, .$$

Therefore, by equation (7.10), $\max_{0 \leq j \leq k-1} U_j^L / \sqrt{n} \leq (Dr_n)/\sqrt{nh} \to 0$. For large enough $n$, the set $\left\{|U_j^L| \geq \theta^L(x) \, \epsilon \, \sqrt{n}\right\}$ becomes empty. Therefore, by the same argument as used to establish (7.7),

$$(7.12) \qquad\qquad \sqrt{nh} \, Q_n^L \xrightarrow{\mathcal{D}} N\left(0, (\theta^L)^2(x)\right) \, .$$

We are now prepared to put these pieces together to finish the proof of the theorem:

$$\left| E\left[\exp\left(it\sqrt{nh}\,Q_n\right)\right] - \exp\left(-t^2\theta^2(x)/2\right) \right| =$$

$$= \left| E\left[\exp\left(it\sqrt{nh}\,\left[Q_n^L + \tilde{Q}_n^L\right]\right)\right] - \exp\left(-t^2(\theta^L)^2(x)/2\right) \right.$$
$$\left. + \exp\left(-t^2(\theta^L)^2(x)/2\right) - \exp\left(-t^2\theta^2(x)/2\right) \right|$$

$$= \left| E\left[\exp\left(it\sqrt{nh}\,Q_n^L\right) + \exp\left(it\sqrt{nh}\,Q_n^L\right)\left(\exp\left[it\sqrt{nh}\,\tilde{Q}_n^L\right] - 1\right)\right] \right.$$
$$\left. - \exp\left(-t^2(\theta^L)^2(x)/2\right) + \exp\left(-t^2(\theta^L)^2(x)/2\right) - \exp\left(-t^2\theta^2(x)/2\right) \right|$$

$$\leq \left| E\left[\exp\left(it\sqrt{nh}\,Q_n^L\right)\right] - \exp\left(-t^2(\theta^L)^2(x)/2\right) \right|$$
$$+ \left| E\left[\exp\left(it\sqrt{nh}\,Q_n^L\right)\left(\exp\left[it\sqrt{nh}\,\tilde{Q}_n^L\right] - 1\right)\right] \right|$$
$$+ \left| \exp\left(-t^2(\theta^L)^2(x)/2\right) - \exp\left(-t^2\theta^2(x)/2\right) \right|$$

$$\leq \left| E\left[\exp\left(it\sqrt{nh}\,Q_n^L\right)\right] - \exp\left(-t^2(\theta^L)^2(x)/2\right) \right|$$
$$+ E\left[\left|\exp\left[it\sqrt{nh}\,\tilde{Q}_n^L\right] - 1\right|\right] + \left| \exp\left(-t^2(\theta^L)^2(x)/2\right) - \exp\left(-t^2\theta^2(x)/2\right) \right| .$$

We analyze each term separately, first letting $n \to \infty$, and then letting $L \to \infty$. For fixed $t$, the first term goes to zero by equation (7.12). The third term goes to 0 by dominated convergence, since $(\theta^L)^2(x) \to \theta^2(x)$ as $L \to \infty$. Only the second term remains to be analyzed. By a Taylor series expansion,

$$\left| \exp\left[it\sqrt{nh}\,\tilde{Q}_n^L\right] - 1 \right| \leq 2\left| t\sqrt{nh}\,\tilde{Q}_n^L \right| .$$

By the Cauchy–Schwarz inequality, the bound will converge to 0 if it can be shown that $nh\,\mathrm{var}\left[(\tilde{Q}_N^L)\right] \to 0$. As $\tilde{Q}_n^L$ satisfies the same dependence assumptions as $Q_n$, the calculations of Lemma 5.1(c) apply. So, it is sufficient to show that $(\theta^L)^2(x) \to 0$ as $L \to \infty$. This follows immediately by dominated convergence.  $\square$

**Lemma 7.3** (Volkonskii and Rozanov [23])**.** *Let $V_1, ..., V_N$ be strong mixing random variables, which are measurable with respect to the $\sigma$-algebras $\mathcal{F}_{i_1}^{j_1}, ..., \mathcal{F}_{i_N}^{j_N}$ respectively, with $1 \leq i_1 < j_1 < i_2 < ... < j_N \leq n$, $i_{l+1} - j_l \geq w \geq 1$, and $|V_l| \leq 1$, for $l = 1, ..., N$. Then*

$$\left| E\left[\prod_{j=1}^{N} V_j\right] - \prod_{j=1}^{N} E[V_j] \right| \leq 16\,(L-1)\,\alpha(w) ,$$

*where $\alpha(w)$ is the strong mixing coefficient.*

# REFERENCES

[1] BOSQ, D. (1998). *Nonparametric statistics for stochastic processes*, "Lecture Notes in Statistics", vol. 110, Springer-Verlag, New York, second edition.

[2] DAVIS, K. (1975). Mean square error properties of density estimates, *Annals of Statistics*, **3**(4), 1025–1030.

[3] DAVIS, K. (1977). Mean integrated square error properties of density estimates, *Annals of Statistics*, **5**(3), 530–535.

[4] DEVROYE, L. (1992). A note on the usefulness of superkernels in density estimates, *Annals of Statistics*, **20**(4), 2037–2056.

[5] DEVROYE, L. and GYÖRFI, L. (1985). *Nonparametric Density Estimation: The $L^1$ View*, Wiley, New York.

[6] FAN, J. (1992). Design-adaptive nonparametric regression, *Journal of the American Statistical Association*, **87**, 998–1004.

[7] GASSER, TH. and MÜLLER, H.-G. (1979). *Kernel estimation of regression functions.* In "Smoothing Techniques for Curve Estimation" (Th. Gasser and M. Rosenblatt, Eds.), Springer Lecture Notes in Mathematics, number 757, Springer-Verlag, Berlin, 23–68.

[8] GASSER, TH. and MÜLLER, H.-G. (1984). Estimating regression functions and their derivatives by the kernel method, *Scandinavian Journal of Statistics*, **11**, 171–185.

[9] GYÖRFI, L.; HÄRDLE, W.; SARDA, P. and VIEU, P. (1989). *Nonparametric Curve Estimation from Time Series*, "Lecture Notes in Statistics", vol. 60, Springer-Verlag, New York.

[10] HALL, P. and HEYDE, C.C. (1980). *Martingale Limit Theory and its Applications*, Academic Press, New York.

[11] HALL, P. and MARRON, J.S. (1988). Choice of kernel order in density estimation, *Annals of Statistics*, **16**(1), 161–173.

[12] IBRAGIMOV, I.A. and KHASMINKSII, R.Z. (1983). Estimation of distribution density belonging to a class of entire functions, *Theor. Probab. Appl.*, **27**(3), 551–562.

[13] MASRY, E. and FAN, J. (1997). Local polynomial estimation of regression functions for mixing processes, *Scandinavian Journal of Statistics*, **24**, 165–179.

[14] MCMURRY, T. and POLITIS, D.N. (2004). Nonparametric regression with infinite order flat-top kernels, *Journal of Nonparametric Statistics*, **16**(3–4), 549–562.

[15] NADARAYA, E.A. (1964). On estimating regression, *Theory Probab. Appl.*, **9**, 141–142.

[16] POLITIS, D.N. (2001). *On nonparametric function estimation with infinite order flat-top kernels.* In "Probability and Statistical Models with Applications" (Ch.A. Charalambides, Markos V. Koutras, and N. Balakrishnan, Eds.), Chapman & Hall/CRC, Boca Raton, 469–483.

[17] POLITIS, D.N. (2003). Adaptive bandwidth choice, *Journal of Nonparametric Statistics*, **25**(4–5), 517–533.

[18]  POLITIS, D.N. and ROMANO, J.P. (1993). On a family of smoothing kernels of infinite order. In "Computing Science and Statistics, Proceedings of the 25th Symposium on the Interface" (M. Tarter and M. Lock, Eds.), The Interface Foundation of North America, 141–145.

[19]  POLITIS, D.N. and ROMANO, J.P. (1995). Bias-corrected nonparametric spectral estimation, *Journal of Time Series Analysis*, **16**(1), 67–103.

[20]  POLITIS, D.N. and ROMANO, J.P. (1999). Multivariate density estimation with general flat-top kernels of infinite order, *Journal of Multivariate Analysis*, **68**, 1–25.

[21]  ROBINSON, P.M. (1983). Nonparametric estimators for time series, *Journal of Time Series Analysis*, **4**(3), 185–207.

[22]  RUPPERT, D.; SHEATHER, S.J. and WAND, M.P. (1995). An effective bandwidth selector for local least squares regression, *Journal of the American Statistical Association*, **90**, 1257–1270.

[23]  VOLKONSKI, V.A. and ROZANOV, YU.A. (1959). Some limit theorems for random functions, *Theory of Probability and its Applications*, **4**, 178–197.

[24]  WAND, M. and RIPLEY, B. (2005). *KernSmooth: Functions for kernel smoothing for Wand and Jones (1995)*, S original by Matt Wand, R port by Brian Ripley, R package version 2.22-15.

[25]  WATSON, G.S. (1964). Smooth regression analysis, *Sankhya Ser. A*, **26**, 359–372.