

Exact and asymptotically optimal bandwidths for kernel estimation of density functionals*

José E. Chacón[†] and Carlos Tenreiro[‡]

July 20, 2011

Abstract

Given a density f we pose the problem of estimating the density functional $\psi_r = \int f^{(r)} f$ for a non-negative even r making use of kernel methods. This is a well-known problem but some of its features remained unexplored. We focus on the problem of bandwidth selection. Whereas all the previous studies concentrate on an asymptotically optimal bandwidth here we study the properties of exact, non-asymptotic ones, and relate them with the former. Our main conclusion is that, despite being asymptotically equivalent, for realistic sample sizes much is lost by using the asymptotically optimal bandwidth. In contrast, as a target for data-driven selectors we propose another bandwidth which retains the small sample performance of the exact one.

Keywords: density functional, exact optimal bandwidth, kernel estimator, normal mixture densities

*This is an electronic version of an article published in *Methodology and Computing in Applied Probability* (Vol. 14, 2012, 523–548), and available on line at <http://dx.doi.org/10.1007/s11009-011-9243-x>

[†]Departamento de Matemáticas, Universidad de Extremadura, Spain. E-mail: jechacon@unex.es

[‡]CMUC, Department of Mathematics, University of Coimbra, Portugal. E-mail: tenreiro@mat.uc.pt

1 Introduction

Let X_1, \dots, X_n denote independent and identically distributed random variables with common and unknown density f on the real line. In this paper we focus on the problem of estimating the functional

$$\psi_r = \int f^{(r)}(x)f(x)dx \quad (1)$$

for $r = 0, 2, 4, \dots$ whenever it makes sense and is finite, where $f^{(r)}$ denotes the r th derivative of f . Notice that for such a functional to be finite it suffices, for instance, that both f and $f^{(r)}$ be square integrable. Moreover, using integration by parts it is easy to show that $\psi_r = (-1)^{r/2} \int \{f^{(r/2)}(x)\}^2 dx$ under some additional conditions on f .

There exists a wide variety of estimates of these functionals. For instance, van Es (1992) proposes an estimator of ψ_0 based on the spacings of the order statistics and Laurent (1997) and Prakasa Rao (1999), respectively, describe series and wavelet estimates for the problem. However, here we will concentrate on kernel estimators,

$$\hat{\psi}_r(g) = \frac{1}{n^2} \sum_{i,j=1}^n L_g^{(r)}(X_i - X_j), \quad (2)$$

where L is the kernel, that is, a real function such that $\int L(x)dx = 1$, $g > 0$ is the bandwidth and $L_g^{(r)}$ represents the r th derivative of the function $L_g(x) = L(x/g)/g$, that is, $L_g^{(r)}(x) = L^{(r)}(x/g)/g^{r+1}$. The motivation for this precise type of kernel estimator can be found, for instance, in Wand and Jones (1995).

This problem is also addressed in many other papers. For instance, the case $r = 0$ (estimation of the integral of a squared density) is closely related with the study of rank-based nonparametric statistics, since it appears in the asymptotic variance of the Wilcoxon signed-rank statistic and in the Pitman asymptotic efficiency of the Wilcoxon test relative to the t -test (see Hettmansperger, 1984). The first kernel estimators of ψ_0 date back to at least Bhattacharya and Roussas (1969), Dmitriev and Tarasenko (1973, 1975) and Schuster (1974), but see also Prakasa Rao (1983), Sheather, Hettmansperger and Donald (1994), and references therein. A recent paper on the topic is Giné and Nickl (2008).

The quantities ψ_2 , ψ_4 and ψ_6 appear in the expression of the asymptotically optimal bandwidths for histogram, frequency polygon and kernel density estimators (see Scott, 1992). The first papers analyzing the kernel-type estimates of ψ_r for arbitrary r , as a particular case of a more general nonlinear functional, are Dmitriev and Tarasenko (1973) and Levit (1978), but the problem of bandwidth selection for the kernel estimator is considered in Hall and Marron (1987) for the first time, although there the kernel estimate is defined as $n(n-1)^{-1} \{ \hat{\psi}_r(g) - n^{-1} L_g^{(r)}(0) \}$, in order to delete the non-stochastic terms in $\hat{\psi}_r(g)$. However, Jones and Sheather (1991) show that indeed the estimator $\hat{\psi}_r(g)$ has improved rates of convergence over the one proposed by Hall and Marron when the bandwidth g is properly chosen. On the other hand, Bickel and Ritov (1988) discuss the information bounds for

this nonparametric problem and propose an efficient estimator. References dealing with adaptive kernel procedures include Wu (1995) and Giné and Mason (2008), among others. Multistep kernel estimators are investigated in Aldershof (1991) and also more recently in Tenreiro (2003) and Chacón and Tenreiro (2011).

As usual for real-valued parameters, we will measure the accuracy of the estimator $\hat{\psi}_r(g)$ through its mean squared error (MSE), defined as $\text{MSE}(g) = \mathbb{E}[\{\hat{\psi}_r(g) - \psi_r\}^2]$. In this sense, the optimal bandwidth can be defined to be $g_{\text{MSE}} = \text{argmin}_{g>0} \text{MSE}(g)$. However, it is not clear at all from its definition that such a minimizer exists, and well-experienced researchers in the field take good care not to refer to this bandwidth, but to its asymptotic counterpart (see Jones and Sheather, 1991, or Wand and Jones, 1995). In fact, the typical approach to bandwidth selection starts from considering an asymptotic expansion of the MSE function, say $\text{AMSE}(g)$, and considering the asymptotically optimal bandwidth $g_0 = \text{argmin}_{g>0} \text{AMSE}(g)$ as a surrogate for g_{MSE} , which is the exact (i.e., non-asymptotic) one. The study of the asymptotically optimal bandwidth presents no doubts about its existence, and even an explicit formula for it is available. But then another question may be raised: how well does g_0 approximate g_{MSE} ? The study of this question leads to the identification of a new bandwidth g_{BA} that annihilates the exact bias of $\hat{\psi}_r(g)$. How well does this new bandwidth approximate g_{MSE} is another question that arises naturally. Therefore, the main purposes of this paper are to present a set of sufficient conditions for the existence of an exact optimal bandwidth and to examine, from an asymptotic and finite sample size point of view, the quality of g_0 and g_{BA} as approximations of the exact optimal bandwidth.

The rest of the paper is organized as follows. In Section 2 we provide mild conditions on the kernel and the density that ensure the existence of an exact optimal bandwidth g_{MSE} and a bias-annihilating bandwidth g_{BA} . In Section 3 we study the asymptotic properties of these bandwidths. In Section 4 we obtain the relative rates of convergence of g_0 and g_{BA} to g_{MSE} and so we quantify the order of these asymptotic approximations. We also establish the order of convergence for $\text{MSE}(g_0) - \text{MSE}(g_{\text{BA}})$ which enables us to compare g_0 and g_{BA} in the sense of the mean squared error. As the results in Section 4 are asymptotic in nature, to assess the quality of the approximations Section 5 contains the case-study of normal mixture densities, for which small- and moderate-sample-size comparisons are made between the three different bandwidths. We will see that for small and moderate sample sizes $\text{MSE}(g_{\text{BA}})$ seems to be much closer to $\text{MSE}(g_{\text{MSE}})$ than $\text{MSE}(g_0)$. In view of these finite sample size results we conclude that bandwidth selectors oriented to g_{BA} should be preferred to the usual ones, which are designed to estimate g_0 . All the proofs are deferred to Section 6.

2 Existence of an exact optimal bandwidth

Recall the definitions of ψ_r and $\hat{\psi}_r(g)$ from (1) and (2) in Section 1. The mean squared error (MSE) of the estimator $\hat{\psi}_r(g)$ can be decomposed as $\text{MSE}(g) = B^2(g) + V(g)$, where $B(g)$ and $V(g)$ are the bias and variance of $\hat{\psi}_r(g)$. If we denote

$$R_{L,r,g}(f) = \mathbb{E}L_g^{(r)}(X_1 - X_2) = \iint L_g^{(r)}(x - y)f(x)f(y)dxdy = \int (L_g^{(r)} * f)(x)f(x)dx,$$

with $*$ standing for the convolution operator, then it is clear that

$$B(g) = \mathbb{E}\hat{\psi}_r(g) - \psi_r = n^{-1}g^{-r-1}L^{(r)}(0) + (1 - n^{-1})R_{L,r,g}(f) - \psi_r. \quad (3)$$

Moreover, using standard U -statistics theory we get that $V(g) = \text{Var}\hat{\psi}_r(g)$ can be written as

$$V(g) = 4(n - 2)(n - 1)n^{-3}\xi_1 + 2(n - 1)n^{-3}\xi_2 - (4n - 6)(n - 1)n^{-3}\xi_0, \quad (4)$$

where $\xi_0 = \{\mathbb{E}L_g^{(r)}(X_1 - X_2)\}^2$, $\xi_1 = \mathbb{E}\{L_g^{(r)}(X_1 - X_2)L_g^{(r)}(X_1 - X_3)\}$ and $\xi_2 = \mathbb{E}[\{L_g^{(r)}(X_1 - X_2)\}^2]$. If we denote

$$S_{L,r,g}(f) = \iiint L_g^{(r)}(x - y)L_g^{(r)}(x - z)f(x)f(y)f(z)dxdydz = \int \{(L_g^{(r)} * f)(x)\}^2 f(x)dx,$$

we just have $\xi_1 = S_{L,r,g}(f)$. Besides, clearly $\xi_0 = \{R_{L,r,g}(f)\}^2$ and, using the fact that $\{L_g^{(r)}(x)\}^2 = \{L^{(r)}(x/g)\}^2/g^{2r+2} = (\{L^{(r)}\}^2)_g/g^{2r+1}$, we can also express $\xi_2 = g^{-2r-1}R_{\{L^{(r)}\}^2,0,g}(f)$.

Combining (3) and (4) with the former representations for ξ_0, ξ_1, ξ_2 , we obtain an exact formula for the MSE of the estimator $\hat{\psi}_r(g)$,

$$\begin{aligned} \text{MSE}(g) &= \{n^{-1}g^{-r-1}L^{(r)}(0) + (1 - n^{-1})R_{L,r,g}(f) - \psi_r\}^2 \\ &\quad + 4(n - 2)(n - 1)n^{-3}S_{L,r,g}(f) + 2(n - 1)n^{-3}g^{-2r-1}R_{\{L^{(r)}\}^2,0,g}(f) \\ &\quad - (4n - 6)(n - 1)n^{-3}\{R_{L,r,g}(f)\}^2. \end{aligned} \quad (5)$$

This exact error formula is the analogue of formula (2.2) in Marron and Wand (1992) for kernel density estimators, and will be useful to explore the existence and limit behavior of the optimal bandwidth as well as for the results in Section 5.

In the following results we will make the next assumptions on the kernel and the density:

(L1) L is a symmetric kernel with bounded and square integrable derivatives up to order r such that $L^{(r)}$ is continuous at zero with $(-1)^{r/2}L^{(r)}(0) > 0$.

(D1) The density f has bounded and square integrable derivatives up to order r .

The next result shows that under these mild conditions there is always an exact optimal bandwidth, that is, a bandwidth which minimizes the exact MSE of the kernel estimator. In this sense, it can be considered as the analogue of Theorem 1 in Chacón et al. (2007a) for kernel density estimators.

Theorem 1. *Under assumptions (L1) and (D1), there exists $g_{\text{MSE}} = g_{\text{MSE},r,n}(f)$ such that $\text{MSE}(g_{\text{MSE}}) \leq \text{MSE}(g)$, for all $g > 0$.*

Notice that the previous result says nothing about the uniqueness of the optimal bandwidth. Presumably, as in the examples in Marron and Wand (1992) it could be possible to find a situation where the optimal bandwidth is not unique, however we do not pursue this further in this paper.

From an asymptotic point of view, however, it is well known that the choice of g can be made on the basis of annihilation of the dominant part of the bias (see Section 4 below). We show next that, in fact, for every density f there is a choice of $g = g_{\text{BA}}$ that makes the estimator $\hat{\psi}_r(g)$ unbiased, that is, that annihilates the exact bias, rather than its asymptotic counterpart.

Theorem 2. *Under assumptions (L1) and (D1), there exists $g_{\text{BA}} = g_{\text{BA},r,n}(f)$ such that $B(g_{\text{BA}}) = 0$.*

The existence of global bandwidths that make the kernel density estimate unbiased at every point has been shown in Chacón et al. (2007b). In fact, strictly speaking we cannot consider it an unbiased estimator since such bandwidths depend on the unknown f , but at least we could say that there exists an ‘unbiased oracle estimator’. However, only a very special class of density functions allows for this situation, namely the class of densities whose characteristic function has bounded support.

In contrast, in the previous result we show that unbiased oracle kernel estimates of ψ_r (not only asymptotically unbiased) exist under the same mild conditions needed for the existence of the optimal bandwidth. This is a key difference between the problems of estimating the density and the functionals ψ_r .

3 Limit behavior of exact bandwidths

From formula (5) and Lemma 1 in Section 6 below it readily follows that $\text{MSE}(g) \rightarrow 0$ for *any* bandwidth sequence $g = g_n$ such that $g \rightarrow 0$ and $ng^{r+1} \rightarrow \infty$ as $n \rightarrow \infty$. Therefore, conditions $g \rightarrow 0$ and $ng^{r+1} \rightarrow \infty$ are sufficient for $\hat{\psi}_r(g)$ to be consistent. It is natural, then, to wonder if the bandwidths g_{MSE} and g_{BA} also fulfill the previous consistency conditions. We will see that the second condition holds quite generally but the same is not necessarily true for the first one. This is similar to the situation with the optimal bandwidth for kernel density estimation, as shown in Chacón et al. (2007b).

Theorem 3. *Under assumptions (L1) and (D1), both $ng_{\text{MSE}}^{r+1} \rightarrow \infty$ and $ng_{\text{BA}}^{r+1} \rightarrow \infty$ as $n \rightarrow \infty$.*

For the analysis of the limit behavior of the sequences g_{MSE} and g_{BA} we use the notation $\varphi_F(t) = \int e^{itx} F(x) dx$, $t \in \mathbb{R}$, for the characteristic function of an integrable real function

F , and for every density f and every symmetric kernel L , we denote

$$\begin{aligned} C_f &= \sup\{r \geq 0 : \varphi_f(t) \neq 0 \text{ a.e. for } t \in [0, r]\}, \\ D_f &= \sup\{t \geq 0 : \varphi_f(t) \neq 0\}, \\ S_L &= \inf\{t \geq 0 : \varphi_L(t) \neq 1\}, \\ T_L &= \inf\{r \geq 0 : \varphi_L(t) \neq 1 \text{ a.e. for } t \geq r\} \end{aligned}$$

A detailed discussion about these quantities is presented in Chacón et al. (2007b). In particular, we remark that all these exist, with C_f, D_f possibly being infinite, $S_L, T_L \in [0, \infty)$, $C_f \leq D_f$ and $S_L \leq T_L$. Notice that, by definition, $S_L > 0$ for superkernels and $S_L = 0$ if L is a kernel of finite order ν (even), that is, if $m_j(L) = 0$ for $j = 1, 2, \dots, \nu - 1$ and $m_\nu(L) \neq 0$ with $|m_\nu|(L) < \infty$, where $m_j(L) = \int u^j L(u) du$ and $|m_j|(L) = \int |u^j L(u)| du$ (see Chacón et al. 2007a).

In the following result we show that both the exact optimal bandwidth g_{MSE} and the exact bias-annihilating bandwidth g_{BA} converge to zero under very general conditions. In particular, if L is a kernel of finite order the convergence to zero takes place with no additional conditions on f other than (D1). The same property occurs in the superkernel case whenever the characteristic function of f has unbounded support.

Theorem 4. *Assume conditions (L1) and (D1). If $S_L = 0$ or $D_f = \infty$ then both $g_{\text{MSE}} \rightarrow 0$ and $g_{\text{BA}} \rightarrow 0$ as $n \rightarrow \infty$.*

In the remaining case $S_L > 0$ and $D_f < \infty$ non-zero limits may occur. In the next example we show that if we use a superkernel and the characteristic function of the density has finite support then any positive number is a possible limit for g_{MSE} or g_{BA} .

Example 1. As in Chacón et al. (2007b), consider the trapezoidal superkernel given by $L(x) = (\pi x^2)^{-1}[\cos x - \cos(2x)]$ for $x \neq 0$ and $L(0) = 3/(2\pi)$, whose characteristic function is $\varphi_L(t) = I_{[0,1]}(|t|) + (2 - |t|) I_{[1,2]}(|t|)$, with $I_A(t)$ standing for the indicator function of the set A , so that $S_L = T_L = 1$. This kernel is symmetric, differentiable of any order, with bounded square integrable derivatives, and such that $L^{(r)}(0) = (-1)^{r/2} [\pi(r+1)(r+2)]^{-1} (2^{r+2} - 1)$, so that it fulfils condition (L1). This is also an example of the so-called flat-top kernels, in the terminology of Politis and Romano (1999).

Consider also the Fejér-de la Vallée-Poussin density, defined as $f(x) = (\pi x^2)^{-1} (1 - \cos x)$ for $x \neq 0$ and $f(0) = 1/(2\pi)$, and let $f_a(x) = f(x/a)/a$ for any $a > 0$; see Figure 1. This density is differentiable of any order, with bounded square integrable derivatives, so that it fulfils condition (D1). The characteristic function of f_a is $\varphi_{f_a}(t) = (1 - a|t|) I_{[-1/a, 1/a]}(t)$, so that $C_{f_a} = D_{f_a} = 1/a$. Besides, we easily obtain $\psi_r = (-1)^{r/2} 2 [\pi a^{r+1} (r+1)(r+2)(r+3)]^{-1}$.

From (14) in Section 6 below we know that $\limsup g \leq a$ for both $g = g_{\text{MSE}}$ and $g = g_{\text{BA}}$. Also, using the formulas for ψ_r and $R_{L,r,g}(f)$ in the Fourier domain given in Section 6,

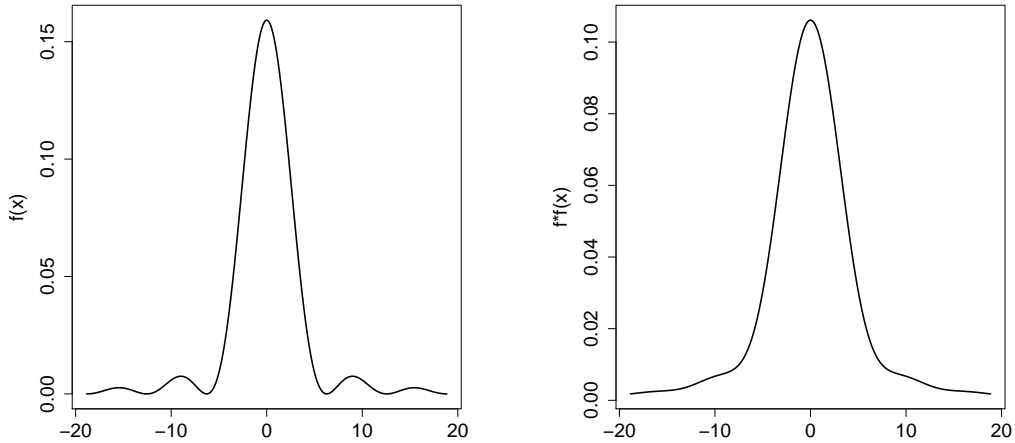


Figure 1: *Fejér-de la Vallée-Poussin density (left) and the convolution with itself (right).*

equation (12), it is not hard to show that, in this case, for $g \in (0, S_L/D_{f_a}] = (0, a]$ we have

$$B(g) = n^{-1}g^{-r-1}L^{(r)}(0) - n^{-1}\psi_r, \quad V(g) = 2(n-1)n^{-3}g^{-2r-1}R_{\{L^{(r)}\}^2,0,g}(f_a) + A,$$

where $A \in \mathbb{R}$ is a constant depending on L, f, r and n , but not on g .

With the formulas for $L^{(r)}(0)$ and ψ_r given above, it is clear that $B(g) \neq 0$ for $g \in (0, a]$, so that it should be that $g_{\text{BA}} \geq a$ for every $n \in \mathbb{N}$ and this, together with the upper bound for the limsup, implies that $g_{\text{BA}} \rightarrow a$ as $n \rightarrow \infty$.

On the other hand, it can be shown that $f * f(x) = 2(\pi x^3)^{-1}(x - \sin x)$ for $x \neq 0$ and $f * f(0) = 1/(3\pi)$, so that $f * f$ is a symmetric density, decreasing for $x > 0$, and the same is true for $f_a * f_a$. Therefore the function $g \mapsto R_{\{L^{(r)}\}^2,0,g}(f_a)$ is decreasing since from (9) in Section 6 we can write $R_{\{L^{(r)}\}^2,0,g}(f_a) = 2 \int_0^\infty \{L^{(r)}(u)\}^2 (f_a * f_a)(gu) du$. This implies that for $g \in (0, a]$ the function $\text{MSE}(g) = B^2(g) + V(g)$ is decreasing and so, that $g_{\text{MSE}} \geq a$ for every $n \in \mathbb{N}$, leading to $g_{\text{MSE}} \rightarrow a$ as $n \rightarrow \infty$.

4 The asymptotically optimal bandwidth

It is well known that the finite sample performance of $\hat{\psi}_r(g)$ depends strongly on the choice of the bandwidth g . In practice, this choice is usually based on the so called asymptotically optimal bandwidth, g_0 , that is, the bandwidth that minimizes the main terms of an asymptotic expansion of $\text{MSE}(g)$ when g tends to zero (see Jones and Sheather, 1991). In order to present such an expansion, some additional conditions on the density f and on the kernel L are needed.

(L2) L is a kernel of finite order ν (even) such that $(-1)^{\nu/2}m_\nu(L) < 0$.

(D2) The density f has bounded and continuous derivatives up to order $r + \nu$.

Under conditions (L1), (L2), (D1) and (D2), if $g \rightarrow 0$ the bias and variance of $\hat{\psi}_r(g)$ given by (3) and (4), respectively, admit the asymptotic expansions

$$B(g) = n^{-1}g^{-r-1}L^{(r)}(0) + g^\nu\psi_{r+\nu}m_\nu(L)/\nu! - n^{-1}\psi_r + o(g^\nu) \quad (6)$$

and

$$V(g) = 4n^{-1}\text{Var}f^{(r)}(X_1) + O(n^{-1}g^\nu + n^{-2}g^{-2r-1}).$$

Therefore,

$$\text{MSE}(g) = 4n^{-1}\text{Var}f^{(r)}(X_1) + B_0^2(g) + o(n^{-2}g^{-2r-2} + n^{-1}g^{\nu-r-1} + g^{2\nu}), \quad (7)$$

where $B_0(g) = n^{-1}g^{-r-1}L^{(r)}(0) + g^\nu\psi_{r+\nu}m_\nu(L)/\nu!$ denotes the asymptotic bias. The asymptotically optimal bandwidth corresponds to the value of g such that $B_0(g) = 0$, that is,

$$g_0 = \left(-\frac{\nu!L^{(r)}(0)}{m_\nu(L)\psi_{r+\nu}n} \right)^{1/(r+\nu+1)}. \quad (8)$$

Notice that the term inside the parenthesis is positive with our assumptions, since we have $(-1)^{r/2}L^{(r)}(0) > 0$, $(-1)^{(r+\nu)/2}\psi_{r+\nu} > 0$ and $(-1)^{\nu/2}m_\nu(L) < 0$.

As the practical choice of g is usually based on this asymptotically optimal bandwidth, g_0 , it is natural to wonder if g_0 is a good approximation of the exact optimal bandwidth, g_{MSE} . In the following theorem we establish the asymptotic equivalence between g_0 , g_{BA} and g_{MSE} , and also the order of convergence to zero of the relative errors $g_0/g_{\text{MSE}} - 1$, $g_{\text{BA}}/g_{\text{MSE}} - 1$ and $g_0/g_{\text{BA}} - 1$.

Theorem 5. *Under assumptions (L1), (L2), (D1) and (D2) we have:*

a) *The bandwidths g_{MSE} , g_{BA} and g_0 are all of the same order; that is,*

$$0 < \liminf n^{1/(r+\nu+1)}g_{\text{MSE}} \leq \limsup n^{1/(r+\nu+1)}g_{\text{MSE}} < \infty,$$

$$0 < \liminf n^{1/(r+\nu+1)}g_{\text{BA}} \leq \limsup n^{1/(r+\nu+1)}g_{\text{BA}} < \infty.$$

b) *Additionally, if $\int |u|\{L^{(r)}(u)\}^2 du < \infty$ then*

$$g_0/g_{\text{MSE}} \rightarrow 1 \quad \text{and} \quad g_{\text{BA}}/g_{\text{MSE}} \rightarrow 1.$$

c) *Moreover, if $|m_{\nu+2}|(L) < \infty$ and f has bounded continuous derivatives up to order $r + \nu + 2$, then there exist constants C , D and E such that*

$$g_0/g_{\text{MSE}} - 1 = C n^{-1/(r+\nu+1)}(1 + o(1)), \quad g_{\text{BA}}/g_{\text{MSE}} - 1 = D n^{-1/(r+\nu+1)}(1 + o(1)),$$

$$g_0/g_{\text{BA}} - 1 = E n^{-\min\{r+1,2\}/(r+\nu+1)}(1 + o(1)).$$

From the previous result we see that asymptotically g_0 and g_{BA} approximate g_{MSE} at the same rate. Following closely the proof of Theorem 6 we can establish that $\text{MSE}(g_0)$ and $\text{MSE}(g_{\text{BA}})$ also approximate $\text{MSE}(g_{\text{MSE}})$ at the same rate. In fact, for $g = g_0$ and $g = g_{\text{BA}}$ we have $\text{MSE}(g) - \text{MSE}(g_{\text{MSE}}) = O(n^{-(2\nu+2)/(r+\nu+1)})$. In the next result we restrict our attention to the order of convergence to zero of $\text{MSE}(g_0) - \text{MSE}(g_{\text{BA}})$ which enables us to compare the bandwidths g_0 and g_{BA} in the sense of the mean squared error.

Theorem 6. *Under assumptions (L1), (L2) and (D1), if f has bounded and continuous derivatives up to order $r + \nu + 2$, $|m_{\nu+2}|(L) < \infty$ and $\int |u|^3 \{L^{(r)}(u)\}^2 du < \infty$, then there exists a constant Λ such that*

$$\text{MSE}(g_0) - \text{MSE}(g_{\text{BA}}) = \Lambda E n^{-\min\{r+2\nu+2, 2\nu+3\}/(r+\nu+1)} (1 + o(1)),$$

where E is the constant appearing in Theorem 5.

Explicit formulas for the constants C , D , E and Λ appearing in Theorems 5 and 6 are given in Section 6, equations (22), (24), (25) and (26), respectively. From them we see that $C = D < 0$ and $\Lambda < 0$ for all densities f whenever $r \geq 2$, and also $E < 0$ if the kernel L is such that $(-1)^{\nu/2} m_{\nu+2}(L) < 0$ (which is in particular true for the Gaussian-based kernel L to be used in the next section). Consequently, from an asymptotic point of view we conclude that g_{BA} is not only a better approximation to g_{MSE} than g_0 but is also a better bandwidth than g_0 in the MSE sense because in this case the constant ΛE appearing in Theorem 6 is strictly positive. As we will see in the next section, even for small and moderate sample sizes $\text{MSE}(g_{\text{BA}})$ seems to be much closer to $\text{MSE}(g_{\text{MSE}})$ than $\text{MSE}(g_0)$.

A different situation may occur when $r = 0$. When the kernel L is of order ν , for all densities f satisfying $\int f^{(\nu)} f^2 / \int f^{(\nu)} f - \int f^2 > 0$ (which seems to be true for all sufficiently regular densities although we were not able to prove it) the constants C and D remain negative but in this case C is always bigger than D which implies that $E > 0$. Hence, the asymptotically optimal bandwidth g_0 is a better asymptotic approximation for g_{MSE} than g_{BA} . Also, we can prove that $\Lambda < 0$, so that in the MSE sense it follows that asymptotically g_0 is better than g_{BA} too. Although this is valid asymptotically, we will see in next section that for small and moderate sample sizes g_{BA} may still be preferable to g_0 in some cases.

5 Case study: normal mixture densities

Our goal in this section is to compare the performance of the three bandwidths, g_{MSE} , g_{BA} and g_0 , in a non-asymptotic way. To this end we work with the exact MSE formula within the class of normal mixture densities, that is, the class of densities f that can be written as $f(x) = \sum_{\ell=1}^k w_{\ell} \phi_{\sigma_{\ell}}(x - \mu_{\ell})$, where $\phi_{\sigma}(x - \mu)$ denotes the density of the normal distribution with mean μ and standard deviation σ . This class is very rich, containing densities with a wide variety of features, such as kurtosis, skewness, multimodality, etc, and has been

previously used for computing exact errors in the context of kernel density estimation (see Marron and Wand, 1992).

Below we find an explicit formula for the MSE given in (5) in the case where f is the aforementioned normal mixture density and L is the Gaussian-based kernel of even order ν considered in Wand and Schucany (1990), given by $L(x) = \sum_{s=0}^{\nu/2-1} a_s \phi^{(2s)}(x)$ with $a_s = (-1)^s (2^s s!)^{-1}$, which has $m_\nu(L)/\nu! = -a_{\nu/2}$. Note that we only need to obtain explicit formulas for $L^{(r)}(0)$, $R_{L,r,g}(f)$, ψ_r , $S_{L,r,g}(f)$ and $R_{\{L^{(r)}\}^2,0,g}(f)$.

For any even $r_1, r_2 \in \mathbb{N}$, $\mu_1, \mu_2, \mu_3 \in \mathbb{R}$ and $\sigma_1 > 0, \sigma_2 > 0, \sigma_3 > 0$, write

$$\begin{aligned} \tilde{\mu} &= \left\{ \sigma_1^{-2} \sigma_2^{-2} (\mu_1 - \mu_2)^2 + \sigma_1^{-2} \sigma_3^{-2} (\mu_1 - \mu_3)^2 + \sigma_2^{-2} \sigma_3^{-2} (\mu_2 - \mu_3)^2 \right\}^{1/2}, \\ \tilde{\sigma} &= \left\{ \sigma_1^{-2} + \sigma_2^{-2} + \sigma_3^{-2} \right\}^{1/2}, \quad \tilde{\mu} = \tilde{\sigma}^{-2} \left\{ \sigma_1^{-2} \mu_1 + \sigma_2^{-2} \mu_2 + \sigma_3^{-2} \mu_3 \right\}, \end{aligned}$$

and $\mu_k^\dagger = \mu_k - \tilde{\mu}$. Then, for $\boldsymbol{\mu} = (\mu_1, \mu_2, \mu_3)$ and $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \sigma_3)$ let us denote

$$\begin{aligned} I_{r_1, r_2}(\boldsymbol{\mu}; \boldsymbol{\sigma}) &= (2\pi)^{-1/2} \phi_{\tilde{\sigma}}(\tilde{\mu}) (\sigma_1 \sigma_2 \sigma_3)^{-1} \\ &\times \sum_{j_1=0}^{r_1} \sum_{j_2=0}^{r_2} \text{OF}(j_1 + j_2) \binom{r_1}{j_1} \binom{r_2}{j_2} H_{r_1-j_1}(\sigma_1^{-1} \mu_1^\dagger) H_{r_2-j_2}(\sigma_2^{-1} \mu_2^\dagger) \sigma_1^{-r_1-j_1} \sigma_2^{-r_2-j_2} \tilde{\sigma}^{-j_1-j_2}, \end{aligned}$$

where for any $p \in \mathbb{N}$ we write $\text{OF}(2p) = (2p-1)(2p-3) \cdots 3 \cdot 1 = (2p)!(2^p p!)^{-1}$, $\text{OF}(2p+1) = 0$ and $H_p(x)$ the p th Hermite polynomial, defined by $H_p(x) = (-1)^p \phi^{(p)}(x)/\phi(x)$.

Theorem 7. For $L(x) = \sum_{s=0}^{\nu/2-1} a_s \phi^{(2s)}(x)$ and $f(x) = \sum_{\ell=1}^k w_\ell \phi_{\sigma_\ell}(x - \mu_\ell)$ we have

$$\begin{aligned} B(g) &= (-1)^{r/2} n^{-1} g^{-r-1} (2\pi)^{-1/2} \sum_{s=0}^{\nu/2-1} (-1)^s a_s \text{OF}(2s+r) \\ &+ \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \left\{ (1-n^{-1}) \sum_{s=0}^{\nu/2-1} a_s g^{2s} \phi_{\sigma_{\ell\ell'}(g)}^{(2s+r)}(\mu_{\ell\ell'}) - \phi_{\sigma_{\ell\ell'}}^{(r)}(\mu_{\ell\ell'}) \right\} \end{aligned}$$

and

$$\begin{aligned} V(g) &= 4(n-2)(n-1)n^{-3} \sum_{\ell_1, \ell_2, \ell_3=1}^k w_{\ell_1} w_{\ell_2} w_{\ell_3} \\ &\times \sum_{s, s'=0}^{\nu/2-1} a_s a_{s'} g^{2s+2s'} I_{2s+r, 2s'+r}(\mu_{\ell_1}, \mu_{\ell_2}, \mu_{\ell_3}; \sigma_{\ell_1}(g), \sigma_{\ell_2}(g), \sigma_{\ell_3}) \\ &+ 2(n-1)n^{-3} \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \sum_{s, s'=0}^{\nu/2-1} a_s a_{s'} g^{2s+2s'} I_{2s+r, 2s'+r}(0, 0, \mu_{\ell\ell'}; g, g, \sigma_{\ell\ell'}) \\ &- (4n-6)(n-1)n^{-3} \left\{ \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \sum_{s=0}^{\nu/2-1} a_s g^{2s} \phi_{\sigma_{\ell\ell'}(g)}^{(2s+r)}(\mu_{\ell\ell'}) \right\}^2, \end{aligned}$$

where $\mu_{\ell\ell'} = \mu_\ell - \mu_{\ell'}$ and $\sigma_{\ell\ell'}^2 = \sigma_\ell^2 + \sigma_{\ell'}^2$ for $\ell, \ell' = 1, 2, \dots, k$ and for any $\sigma > 0$ we write $\sigma(g) = (\sigma^2 + g^2)^{1/2}$.

For L and f as given in the previous theorem we can also write the asymptotic bias as

$$B_0(g) = (-1)^{r/2} n^{-1} g^{-r-1} (2\pi)^{-1/2} \sum_{s=0}^{\nu/2-1} (-1)^s a_s \text{OF}(2s+r) - a_{\nu/2} \psi_{r+\nu} g^\nu$$

and its minimizer as

$$g_0 = \left| (2\pi)^{-1/2} \sum_{s=0}^{\nu/2-1} (-1)^s a_s \text{OF}(2s+r) a_{\nu/2}^{-1} \psi_{r+\nu}^{-1} n^{-1} \right|^{1/(r+\nu+1)},$$

with

$$\psi_{r+\nu} = \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \phi_{\sigma_{\ell\ell'}}^{(r+\nu)}(\mu_{\ell\ell'}),$$

as shown in the proof of Theorem 7 in Section 6.

The previous results allow us to compare the exact MSE function and its minimizer g_{MSE} with their asymptotic approximations. The first form of asymptotic MSE is defined as the dominant expression in equation (7), that is, $\text{AMSE}_0(g) = 4n^{-1} \text{Var} f^{(r)}(X_1) + B_0^2(g)$. A second asymptotic approximation of the MSE is obtained by combining the asymptotic variance and the exact bias, that is, $\text{AMSE}_1(g) = 4n^{-1} \text{Var} f^{(r)}(X_1) + B^2(g)$. Notice that g_0 and g_{BA} minimize $\text{AMSE}_0(g)$ and $\text{AMSE}_1(g)$, respectively.

To analyze the finite sample behavior of the asymptotic approximations we use two of the 15 normal mixture densities introduced in Marron and Wand (1992). Precisely, we focus on density #1 (standard normal density) and density #12 (asymmetric claw density), corresponding to the cases where the difficulty in estimating the density itself is low and high. The comments that follow regarding density #1 are equally applicable to the first ten Marron-Wand densities, whereas the situation for density #12 is shared by the last five of them. Therefore, we prefer to stay with these two representatives and not to reproduce the full graphics concerning all the densities to save space.

Figure 2 shows the graphs of $\text{MSE}(g)$ and its two approximations with g on a \log_{10} scale, for the normal density #1 and the asymmetric claw density #12. The kernel used in the estimator was the standard normal density, which has order 2. Note that in the case of density #1 the two asymptotic versions provide quite reliable approximations to the MSE, and the same holds for their minimizers. Even so, the quality of g_0 as surrogate for g_{MSE} deteriorates more rapidly than that of g_{BA} as r increases. For density #12, however, the situation is more clearly favorable to g_{BA} , which stays close to g_{MSE} whereas g_0 gives a poor approximation, even for $r = 0$.

To get further insight into the accuracy of the approximations to g_{MSE} it is quite informative to look at Figure 3, which shows, in a \log_{10} - \log_{10} scale the values of g_{MSE} , g_{BA} and g_0 as a function of the sample size. Eventually these graphs take the form of a straight line with slope $-1/(r+3)$ (see Theorem 5). As anticipated in Figure 2, whereas for the normal density both approximations are quite close to the exact g_{MSE} value, for density

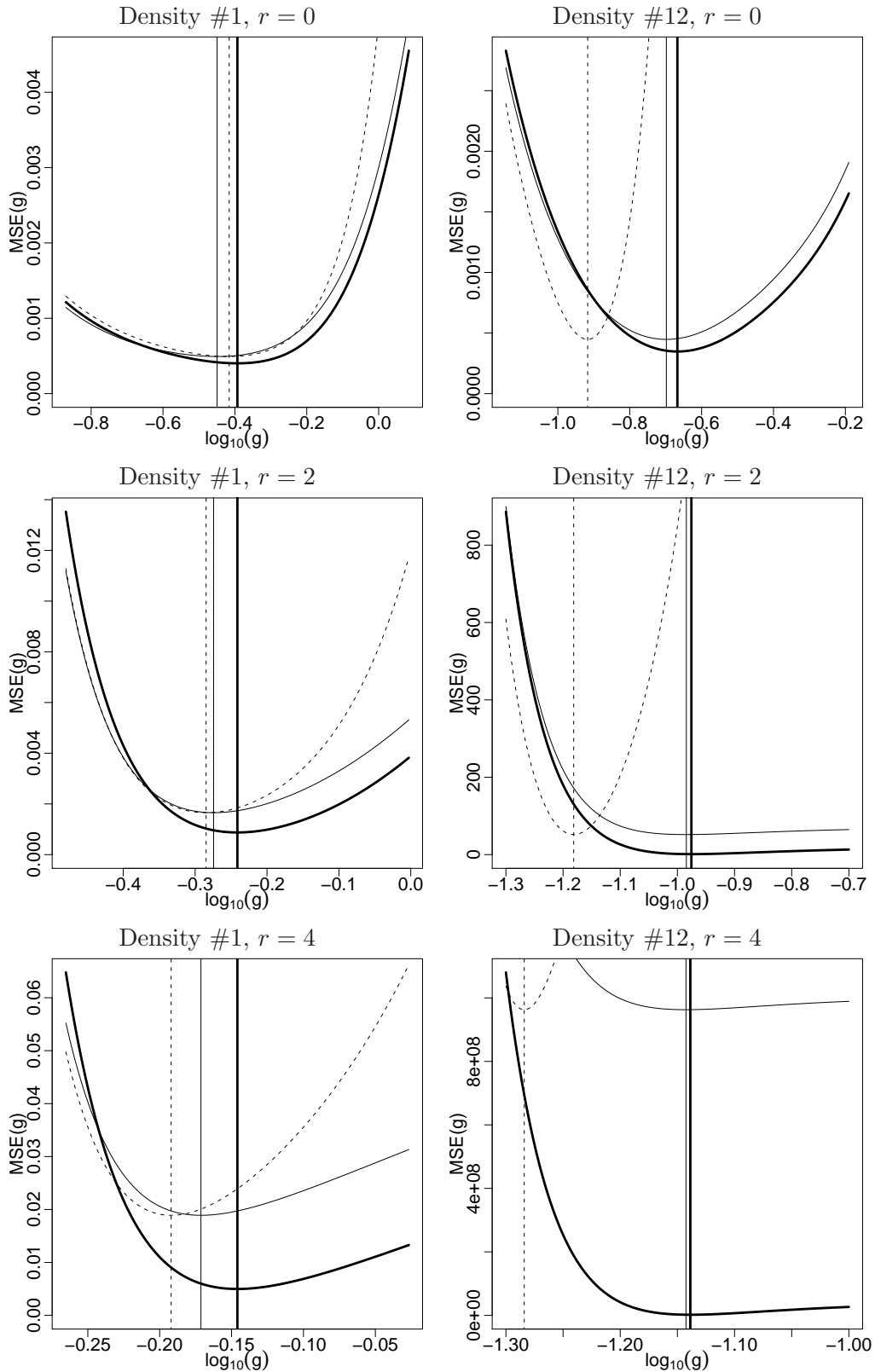


Figure 2: Comparison of $\text{MSE}(g)$ (thick solid line) and its asymptotic approximations $\text{AMSE}_0(g)$ (dashed line) and $\text{AMSE}_1(g)$ (thin solid line) for the normal density #1 (left column) and the asymmetric claw density #12 (right column) and for $r = 0, 2, 4$. Their respective minimizers, g_{MSE} , g_0 and g_{BA} , are indicated by the vertical lines. The kernel is the standard normal density and the sample size is $n = 100$.

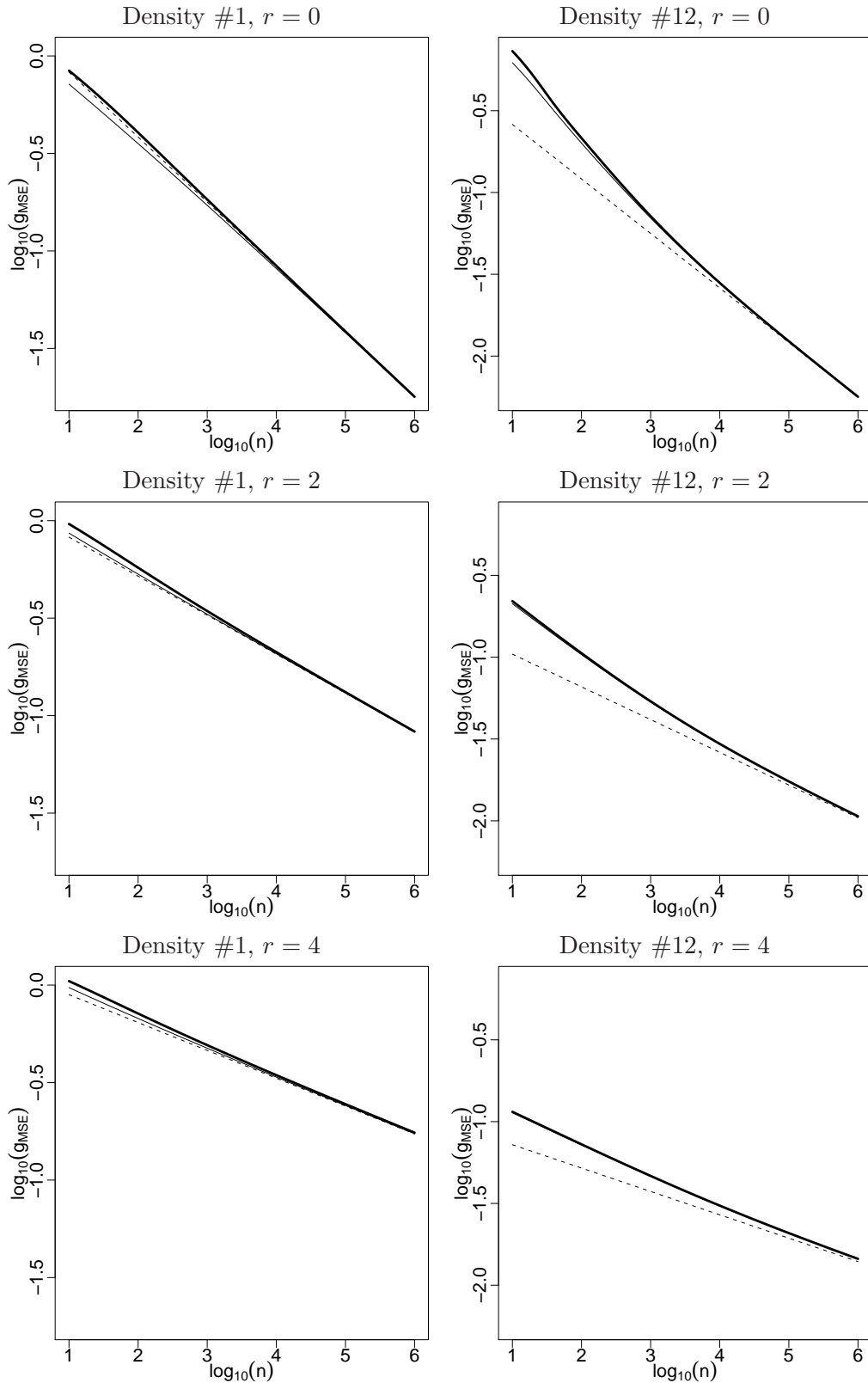


Figure 3: Comparison of g_{MSE} (thick solid line) and its asymptotic approximations g_0 (dashed line) and g_{BA} (thin solid line) for the normal density #1 (left column) and the asymmetric claw density #12 (right column) and for $r = 0, 2, 4$, as a function of the sample size. The kernel is the standard normal density.

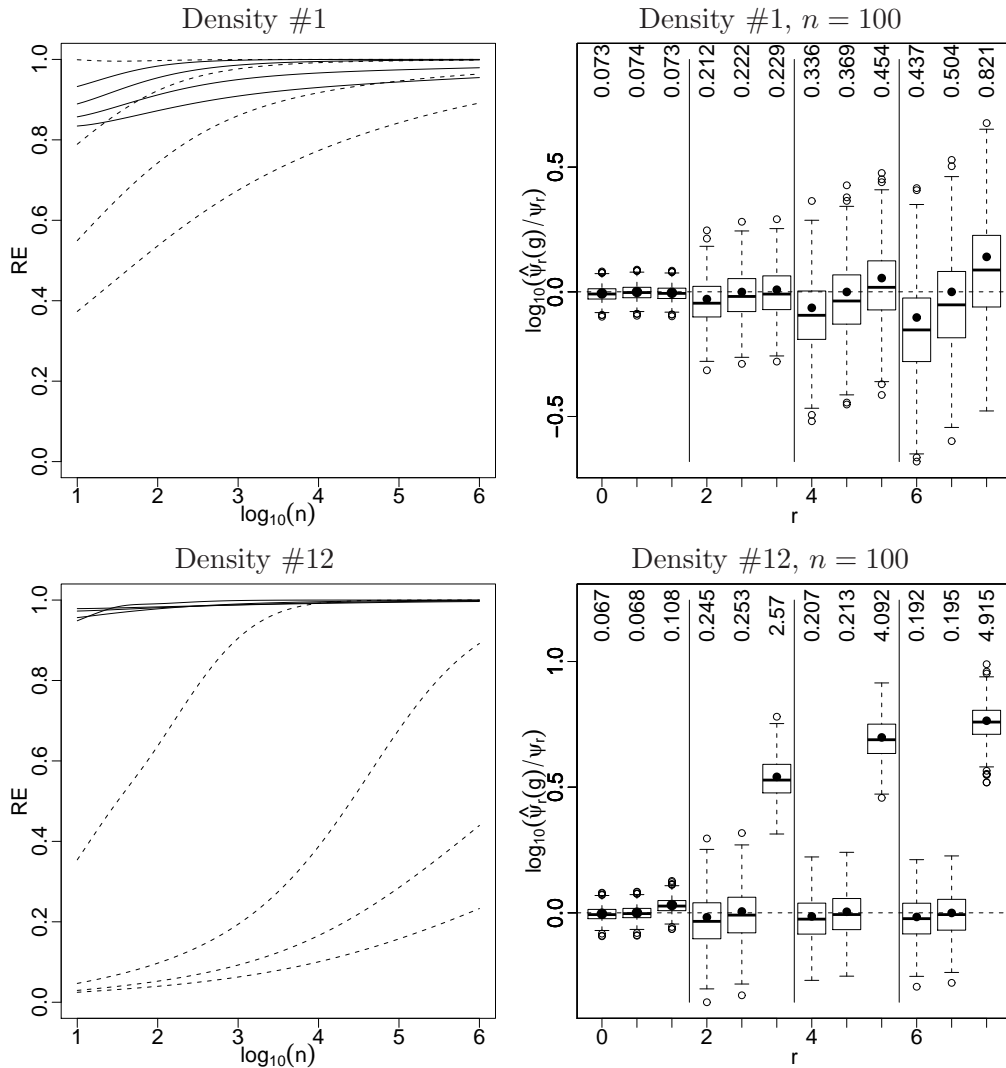


Figure 4: Relative efficiencies of g_0 and g_{BA} with respect g_{MSE} in terms of MSE (left) and distribution boxplots for the estimator $\hat{\psi}_r(g)$ with $g = g_{MSE}$, g_{BA} and g_0 , together with their respective root sample MSEs (right).

#12 a sample size between $n = 1000$ and $n = 10000$ is needed before g_0 can be considered a reasonable surrogate for g_{MSE} , as opposite to g_{BA} , which stays close to g_{MSE} even for small values of n .

Figure 3 shows how well g_0 and g_{BA} approximate g_{MSE} , but perhaps it is even more interesting to compare their performances in MSE terms. In Figure 4 (left column) we show the relative efficiencies $[\text{MSE}(g_{MSE})/\text{MSE}(g)]^{1/2}$ for $g = g_{BA}$ (solid lines) and $g = g_0$ (dashed lines) against $\log_{10}(n)$ for $r = 0, 2, 4, 6$. As expected, for each of $g = g_{BA}$ and $g = g_0$ the efficiency graphs are naturally placed in descending order as r increases, that is, for $g = g_{BA}$ the top solid curve in each plot corresponds to $r = 0$ and the bottom

solid curve corresponds to $r = 6$, and similarly for $g = g_0$. This reflects the fact that the approximations to g_{MSE} given by g_{BA} and g_0 get worse (in the MSE sense) as the degree of derivative r increases, as predicted by the asymptotic theory (see Theorems 5 and 6 above).

However, even though both $\text{MSE}(g_{\text{BA}})$ and $\text{MSE}(g_0)$ exhibit the same relative order of convergence to $\text{MSE}(g_{\text{MSE}})$, we can see in the left column of Figure 4 that for small and moderate sample sizes there are also marked differences between g_{BA} and g_0 . Whereas for $n \geq 10$ and the cases $r = 0, 2, 4, 6$ represented in Figure 4 the efficiency of g_{BA} is always greater than 90%, showing that the loss in changing the goal from g_{MSE} to g_{BA} is nearly negligible, in some cases Figure 4 shows that the use of the bandwidth g_0 may lead to a very disappointing performance of the estimator.

Specifically, for the normal density #1 g_0 is even more efficient than g_{BA} for density #1 when $r = 0$, and it is also quite acceptable for $r = 2$, but for $r \geq 4$ the efficiency of g_0 decays rapidly, and it is already lower than 70% (for $r = 4$) or 50% (for $r = 6$) for sample size $n = 100$. This effect is even more dramatic for the case of the asymmetric claw density #12: for sample size $n = 100$ the efficiency of g_0 is about 60% for $r = 0$ and it is lower than 10% for $r \geq 2$.

Our conclusion is that g_0 can be a bad surrogate for g_{MSE} , especially for $r \geq 4$. This is quite a striking conclusion, since g_0 is the usual target bandwidth for plug-in bandwidth selection methods for the estimation of ψ_r .

In the right column of Figure 4 we show the boxplots for the distribution of $\log_{10}(\hat{\psi}_r(g)/\psi_r)$ based on 500 generated samples of size $n = 100$. In each graph we have vertical lines dividing the cases according to $r = 0, 2, 4, 6$ and, for each of these cases, we have three boxplots corresponding to the use of the theoretical $g = g_{\text{MSE}}$, $g = g_{\text{BA}}$ and $g = g_0$ in the estimator, from left to right. We have also added a solid circle to each boxplot indicating the sample mean of the distribution and a number on top with the square root of the sample MSE of $\hat{\psi}_r(g)/\psi_r$.

The boxplots show the reasons for the bad efficiency results of g_0 . Although this bandwidth is meant to annihilate the bias term asymptotically, it looks like g_0 does not get close to this goal for moderate sample sizes, since $\hat{\psi}_r(g_0)$ clearly overestimates ψ_r in mean, especially for $r \geq 4$. Moreover, this occasionally large bias does not come with a reduction in variance, since in fact $\hat{\psi}_r(g_0)$ is more variable than the other two estimators. Both effects (in bias and variance) are highly stressed for the case of density #12. In contrast, it is possible to observe how the estimator using the bandwidth g_{BA} is unbiased, as it should be by definition, at the expense of only a slight increase of variance over g_{MSE} . Nevertheless, the distributions of the estimator with g_{MSE} or g_{BA} are very similar.

Finally, we explore the consequences of using a higher order kernel for estimating ψ_r . From expansion (7) it can be shown that $\text{MSE}(g_{\text{MSE}})$ is of order $n^{-(\min\{r,\nu\}+\nu+1)/(r+\nu+1)}$, thus achieving the optimal n^{-1} rate whenever $\nu \geq r$; moreover, the estimator has also the optimal asymptotic variance if $\nu \geq r + 2$ (see Wand and Jones, 1995, p. 67–70).

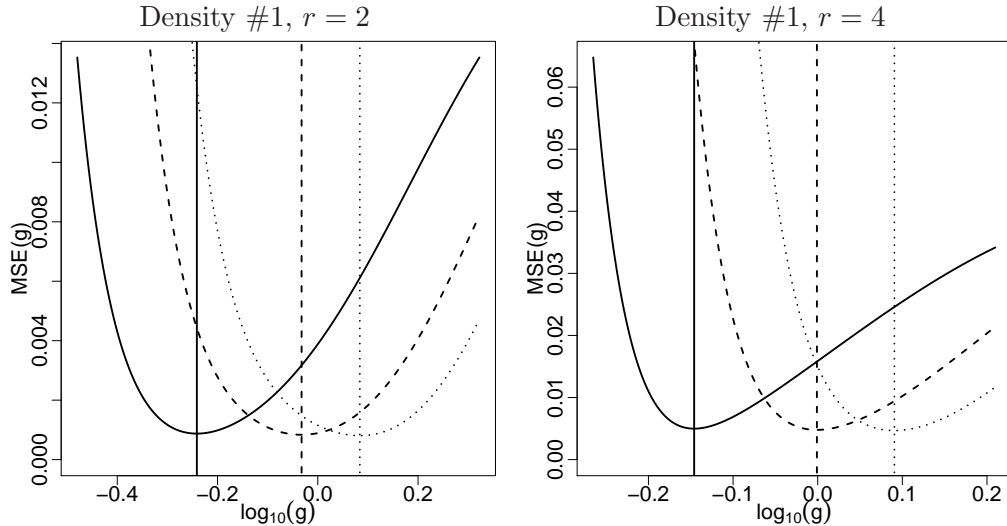


Figure 5: Comparison of $MSE(g)$ and g_{MSE} for normal kernels of order $\nu = 2$ (thick solid line), $\nu = 4$ (dashed line) and $\nu = 6$ (thin solid line) for the normal density #1 and for $r = 2$ (left) and $r = 4$ (right). The sample size is $n = 100$.

However, for the problem of density estimation Marron and Wand (1992) showed that these asymptotics may need a very large sample size before their effects begin to take place, and that the gains for small or moderate sample sizes are usually minor. To investigate the effect of using a higher order kernel we show in Figure 5 how the MSE function and its minimizer varies for kernel orders $\nu = 2, 4, 6$ for the normal density #1 and sample size $n = 100$.

The effect of using a higher order kernel is that the whole $MSE(g)$ curve and g_{MSE} are moved to the right. This is a consequence of the reduction in bias together with an increase in variance, as noted in Marron and Wand (1992). We investigated how the use of higher order kernels affected the asymptotic approximations to $MSE(g)$ and g_{MSE} as well, but we have omitted their plots in Figure 5 for the sake of clarity, since they led basically to the same conclusions. We note, however, that the approximation of g_{MSE} by g_0 seems to get worse in the case of higher order kernels.

Also in accordance with the case of density estimation, it is noticeable that the corresponding heights of the minima are not very different for the three kernel orders represented in Figure 5, meaning that for $n = 100$ there is little to be gained by using a higher order kernel. To explore this fact in more detail, Figure 6 shows $\log_{10} MSE(g_{MSE})$ as a function of $\log_{10}(n)$ for $r = 0, 2, 4, 6$ when kernels of orders $\nu = 2$ (solid line), $\nu = 4$ (dashed line) and $\nu = 6$ (dotted line) are used to estimate ψ_r . Notice that for $r = 0$ and $r = 2$, the three kernel orders lead to an optimal MSE rate n^{-1} , which is consistent with the corresponding graphs having eventually the shape of a straight line with slope -1 , and results in nearly identical MSE performance. For $r = 4$ and $r = 6$ the use of higher order kernels

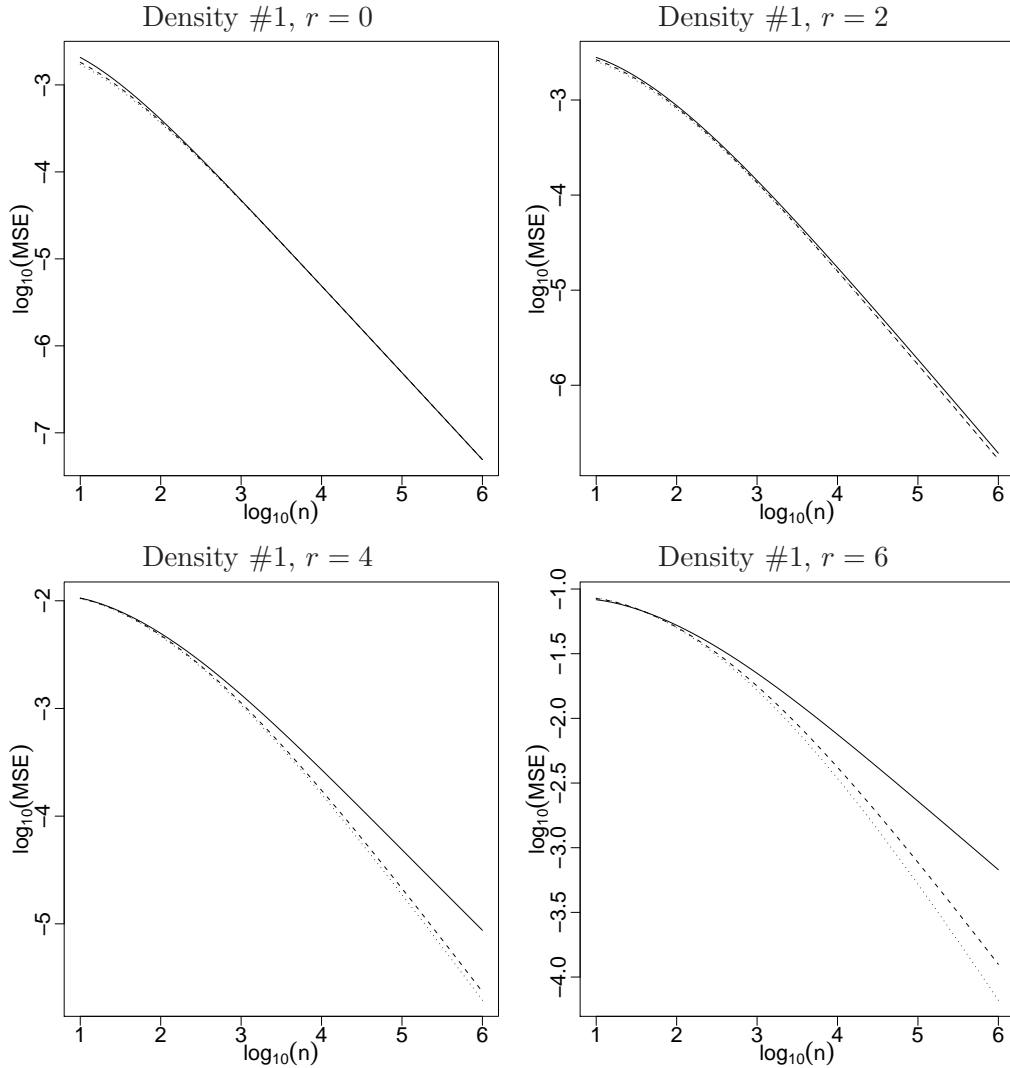


Figure 6: Comparison of $\log_{10} \text{MSE}(\text{gmSE})$ when kernels of orders $\nu = 2$ (solid line), $\nu = 4$ (dashed line) and $\nu = 6$ (dotted line) are used to estimate ψ_r for the normal density #1.

improves the performance asymptotically, as predicted by theory, but a large sample size is usually needed before these asymptotics become to produce a noticeable effect. The sample size at which higher order kernel estimators become clearly preferable varies with the true underlying density #1, it is about $n = 1000$ in the case of the normal density, but for the asymmetric claw density #12 (not shown) higher order kernels still have not become dominant for $n = 10^6$.

Remark 1. For the case $\nu = 2$ the exact MSE formula for normal mixture densities can be found in Aldershof's thesis (1991). Even so, its consequences (as extracted from the previous discussions) had not been fully explored yet.

6 Proofs

We start by presenting some properties of $R_{L,r,g}(f)$ and $S_{L,r,g}(f)$ as functions of g . Let us denote $\psi_{r,s} = \int f^{(r)} f^{(s)} f$.

Lemma 1. *Under assumptions (L1) and (D1), we have:*

- a) *The function $g \mapsto R_{L,r,g}(f)$ is continuous and such that $\lim_{g \rightarrow 0} R_{L,r,g}(f) = \psi_r \int L$ and $\lim_{g \rightarrow \infty} g^{r+1} R_{L,r,g}(f) = L^{(r)}(0)$.*
- b) *The function $g \mapsto S_{L,r,g}(f)$ is continuous and such that $\lim_{g \rightarrow 0} S_{L,r,g}(f) = \psi_{r,r} (\int L)^2$ and $\lim_{g \rightarrow \infty} g^{2r+2} S_{L,r,g}(f) = \{L^{(r)}(0)\}^2$.*

Proof. Using the fact that $L^{(j)}$ and $f^{(j)}$ are bounded and square integrable, for $j = 0, 1, \dots, r$, and the same tools as in Hall and Marron (1987), it is straightforward to check that we can write

$$R_{L,r,g}(f) = \int L(u) (f^{(r)} * \bar{f})(gu) du, \quad (9)$$

with $\bar{f}(x) = f(-x)$. Therefore, as $L \in L_1$ the continuity and the first limit in part a) follow from the Dominated Convergence Theorem (DCT) and the boundedness and continuity of the convolution product of square integrable functions, together with the fact that $(f^{(r)} * \bar{f})(0) = \psi_r$. For the second limit, using again the DCT, together with the boundedness and continuity of $L^{(r)}$ at zero, we obtain

$$\lim_{g \rightarrow \infty} g^{r+1} R_{L,r,g}(f) = \lim_{g \rightarrow \infty} \iint L^{(r)}\left(\frac{x-y}{g}\right) f(x) f(y) dx dy = L^{(r)}(0)$$

as stated.

The proof of part b) can be obtained in a similar way. For the first limit we start by writing

$$\begin{aligned} S_{L,r,g}(f) &= \iiint L(u) L(v) f^{(r)}(gu-x) f^{(r)}(gv-x) \bar{f}(x) dx dudv \\ &= \iint L(u) L(v) (f^{(r)} \odot f^{(r)} \odot \bar{f})(gu, gv) dudv, \end{aligned} \quad (10)$$

where we are denoting

$$(\alpha \odot \beta \odot \gamma)(y, z) = \int \alpha(y-x) \beta(z-x) \gamma(x) dx.$$

Reasoning as in the proof of Theorem 21.33 in Hewitt and Stromberg (1965), for $\alpha, \beta, \gamma \in L_3$ it can be shown that $\alpha \odot \beta \odot \gamma$ is a bounded continuous function. Consequently, as $f, f^{(r)} \in L_3$ (since they are bounded and square integrable), we get the stated limit by using again the DCT, together with the fact that $(f^{(r)} \odot f^{(r)} \odot \bar{f})(0, 0) = \psi_{r,r}$. The second limit again follows from a direct application of the DCT, since

$$\lim_{g \rightarrow \infty} g^{2r+2} S_{L,r,g}(f) = \lim_{g \rightarrow \infty} \iiint L^{(r)}\left(\frac{x-y}{g}\right) L^{(r)}\left(\frac{x-z}{g}\right) f(x) f(y) f(z) dx dy dz = \{L^{(r)}(0)\}^2.$$

□

Proof of Theorem 1. From the previous lemma, together with (3) and (4), we conclude that $B^2(g)$ and $V(g)$ are continuous functions such that

$$\lim_{g \rightarrow 0} B^2(g) = \infty, \quad \lim_{g \rightarrow \infty} B^2(g) = \psi_r^2, \quad \lim_{g \rightarrow 0} V(g) = \infty, \quad \lim_{g \rightarrow \infty} g^{2r+2}V(g) = 0.$$

So the MSE function, which equals $\text{MSE}(g) = B^2(g) + V(g)$, is a continuous function such that

$$\lim_{g \rightarrow 0} \text{MSE}(g) = \infty, \quad \lim_{g \rightarrow \infty} \text{MSE}(g) = \psi_r^2.$$

Therefore, to show that there exists a value $g_{\text{MSE}} = g_{\text{MSE},r,n}(f)$ minimizing the MSE function, it suffices to show that, for big enough g_* , we have $\text{MSE}(g_*) < \psi_r^2$. So if we define

$$D(g) = \text{MSE}(g) - \psi_r^2 = [B^2(g) - \psi_r^2] + V(g),$$

all that we need to show is that, for some $\rho > 0$, we have $\lim_{g \rightarrow \infty} g^\rho D(g) < 0$.

But using the previous lemma we have

$$\lim_{g \rightarrow \infty} g^{2r+2}[B^2(g) - \psi_r^2] = -\text{sgn}(\psi_r L^{(r)}(0)) \cdot \infty,$$

where $\text{sgn}(x) = x/|x|$ for $x \neq 0$. As our assumptions imply that $\text{sgn} \psi_r = \text{sgn} L^{(r)}(0) = (-1)^{r/2}$, it immediately follows that $\lim_{g \rightarrow \infty} g^{2r+2}[B^2(g) - \psi_r^2] = -\infty$. This, together with the limit properties of the variance allows us to conclude that $\lim_{g \rightarrow \infty} g^{2r+2}D(g) = -\infty$ and so the proof is complete. \square

Proof of Theorem 2. From the previous lemma, together with (3), we know that $B(g)$ is a continuous function such that

$$\lim_{g \rightarrow 0} B(g) = (\text{sgn} L^{(r)}(0)) \cdot \infty, \quad \lim_{g \rightarrow \infty} B(g) = -\psi_r.$$

Again, our assumptions imply that $\text{sgn}(-\psi_r \cdot L^{(r)}(0)) = -1$, which yields the proof using Bolzano's theorem. \square

Proof of Theorem 3. For $g = g_{\text{MSE}}$ or $g = g_{\text{BA}}$, suppose that ng^{r+1} does not converge to infinity. Then ng^{r+1} has a subsequence which is upper bounded by some positive constant C . Therefore, along that subsequence we have $g \rightarrow 0$.

For $g = g_{\text{MSE}}$ this implies that

$$\limsup_{n \rightarrow \infty} \text{MSE}(g_{\text{MSE}}) \geq \limsup_{n \rightarrow \infty} B^2(g_{\text{MSE}}) = \limsup_{n \rightarrow \infty} \{n^{-1} g_{\text{MSE}}^{-r-1} L^{(r)}(0)\}^2 \geq \{L^{(r)}(0)/C\}^2 > 0,$$

which contradicts the fact that $0 \leq \text{MSE}(g_{\text{MSE}}) \leq \text{MSE}(n^{-1/(r+2)}) \rightarrow 0$ that follows from (5) together with the previous lemma.

Similarly, for $g = g_{\text{BA}}$ we would obtain that

$$0 = B^2(g_{\text{BA}}) = \limsup_{n \rightarrow \infty} B^2(g_{\text{BA}}) = \limsup_{n \rightarrow \infty} \{n^{-1} g_{\text{BA}}^{-r-1} L^{(r)}(0)\}^2 \geq \{L^{(r)}(0)/C\}^2 > 0,$$

so that the result also follows by contradiction. \square

Proof of Theorem 4. Let us prove the result for g_{MSE} . Denote by $\Lambda_{f,L}$ the set of accumulation points of the sequence (g_{MSE}) . Take $0 < \lambda \in \Lambda_{f,L}$ and (g_{n_k}) a subsequence of (g_{MSE}) such that $\lambda = \lim_{k \rightarrow \infty} g_{n_k}$. Writing $B(g; n)$ and $\text{MSE}(g; n)$ for $B(g)$ and $\text{MSE}(g)$, respectively, from equalities (3) and (4) we get that, for fixed $g > 0$,

$$\lim_{n \rightarrow \infty} \text{MSE}(g; n) = \lim_{n \rightarrow \infty} B^2(g; n) = \{R_{L,r,g}(f) - \psi_r\}^2,$$

so that using Lemma 1 and Theorem 3, we obtain

$$\begin{aligned} 0 &= \lim_{g \rightarrow 0} \{R_{L,r,g}(f) - \psi_r\}^2 = \lim_{g \rightarrow 0} \lim_{k \rightarrow \infty} B^2(g; n_k) = \lim_{g \rightarrow 0} \lim_{k \rightarrow \infty} \text{MSE}(g; n_k) \\ &\geq \lim_{k \rightarrow \infty} \text{MSE}(g_{n_k}; n_k) \geq \lim_{k \rightarrow \infty} B^2(g_{n_k}; n_k) = \{R_{L,r,\lambda}(f) - \psi_r\}^2. \end{aligned}$$

Therefore

$$\Lambda_{f,L} \subset \{\lambda \geq 0 : R_{L,r,\lambda}(f) = \psi_r\} = \{\lambda \geq 0 : \int_0^\infty t^r |\varphi_f(t)|^2 [1 - \varphi_L(t\lambda)] dt = 0\}, \quad (11)$$

since from Parseval's formula, together with $\varphi_{f^{(r)}}(t) = (it)^r \varphi_f(t)$ (see Butzer and Nessel, 1971, Proposition 5.2.19) we easily get that

$$\psi_r = (-1)^{r/2} \pi^{-1} \int_0^\infty t^r |\varphi_f(t)|^2 dt, \quad R_{L,r,\lambda}(f) = (-1)^{r/2} \pi^{-1} \int_0^\infty t^r |\varphi_f(t)|^2 \varphi_L(t\lambda) dt. \quad (12)$$

Additionally, we also have

$$\Lambda_{f,L} \subset \left[0, \min \left(\frac{S_L}{C_f}, \frac{T_L}{D_f} \right) \right]. \quad (13)$$

This is because in fact, if $\lambda > 0$ is such that $\lambda \in \Lambda_{f,L}$, from (11) we have

$$\int_0^{C_f} t^r |\varphi_f(t)|^2 [1 - \varphi_L(t\lambda)] dt = 0 \quad \text{and} \quad \int_{T_L/\lambda}^\infty t^r |\varphi_f(t)|^2 [1 - \varphi_L(t\lambda)] dt = 0.$$

Taking into account that φ_L is a real function (for L being symmetric) and such that $1 - \varphi_L(t\lambda) \geq 0$, from the first equality we conclude that $\varphi_L(s) = 1$ for all $0 \leq s \leq \lambda C_f$, and then $S_L \geq \lambda C_f$, that is, $\lambda \leq S_L/C_f$. From the second equality we have $\varphi_f(t) = 0$ for all $t \geq T_L/\lambda$, and then $D_f \leq T_L/\lambda$, that is, $\lambda \leq T_L/D_f$.

From (13) we finally get

$$0 \leq \limsup_{n \rightarrow \infty} g_{\text{MSE}} \leq \min \left(\frac{S_L}{C_f}, \frac{T_L}{D_f} \right), \quad (14)$$

which concludes the proof for g_{MSE} .

Similarly, notice that any λ being an accumulation point of g_{BA} the equality $R_{L,r,\lambda}(f) - \psi_r = 0$ should also hold, due to Theorem 3, the continuity properties in Lemma 1 and the fact that $B(g_{\text{BA}}) = 0$. Consequently, (14) is also true for g_{BA} and so the desired result. \square

As a tool for the proof of Theorem 5 we will need the following lemma, which follows directly from expressions (9) and (10) for $R_{L,r,g}(f)$ and $S_{L,r,g}(f)$, respectively, the differentiation theorem under the integral sign and standard Taylor expansions.

Lemma 2. *Under assumptions (L1), (L2), (D1) and (D2) we have:*

a) *The function $g \mapsto R_{L,r,g}(f)$ is differentiable with*

$$\begin{aligned} R_{L,r,g}(f) &= \psi_r + g^\nu \psi_{r+\nu} m_\nu(L) / \nu! + o(g^\nu), \\ dR_{L,r,g}(f)/dg &= g^{\nu-1} \psi_{r+\nu} m_\nu(L) / (\nu-1)! + o(g^{\nu-1}). \end{aligned}$$

Additionally, if $|m_{\nu+2}|(L) < \infty$ and f has bounded continuous derivatives up to order $r+\nu+2$, the previous residual term $o(g^\nu)$ may be replaced by $g^{\nu+2} \psi_{r+\nu+2} m_{\nu+2}(L) / (\nu+2)! + o(g^{\nu+2})$.

b) *If $\int |u| \{L^{(r)}(u)\}^2 du < \infty$, the function $g \mapsto R_{\{L^{(r)}\}^2, 0, g}(f)$ is differentiable and such that $dR_{\{L^{(r)}\}^2, 0, g}(f)/dg = o(1)$.*

c) *The function $g \mapsto S_{L,r,g}(f)$ is differentiable and such that*

$$dS_{L,r,g}(f)/dg = 2g^{\nu-1} \psi_{r+\nu,r} m_\nu(L) / (\nu-1)! + o(g^{\nu-1}).$$

Proof of Theorem 5. a) From expansion (7) and taking for g the asymptotically optimal bandwidth (8), we easily get

$$n^{2\nu/(r+\nu+1)} \left(\text{MSE}(g_0) - 4n^{-1} \text{Var} f^{(r)}(X_1) \right) = o(1)$$

and then, as $\text{MSE}(g_{\text{MSE}}) \leq \text{MSE}(g_0)$,

$$\limsup n^{2\nu/(r+\nu+1)} \left(\text{MSE}(g_{\text{MSE}}) - 4n^{-1} \text{Var} f^{(r)}(X_1) \right) < \infty. \quad (15)$$

Moreover, using the fact that $g_{\text{MSE}} \rightarrow 0$ (that follows from Theorem 4, due to condition (L2)), from expansion (7) we also get

$$\begin{aligned} & n^{2\nu/(r+\nu+1)} \left(\text{MSE}(g_{\text{MSE}}) - 4n^{-1} \text{Var} f^{(r)}(X_1) \right) \\ &= \left((n^{1/(r+\nu+1)} g_{\text{MSE}})^{-r-1} L^{(r)}(0) + (n^{1/(r+\nu+1)} g_{\text{MSE}})^\nu m_\nu(L) \psi_{r+\nu} / \nu! \right)^2 \\ & \quad + o \left((n^{1/(r+\nu+1)} g_{\text{MSE}})^{-2r-2} + (n^{1/(r+\nu+1)} g_{\text{MSE}})^{\nu-r-1} + (n^{1/(r+\nu+1)} g_{\text{MSE}})^{2\nu} \right), \end{aligned}$$

which contradicts (15) if $\liminf n^{1/(r+\nu+1)} g_{\text{MSE}} = 0$ or $\limsup n^{1/(r+\nu+1)} g_{\text{MSE}} = \infty$. Therefore the proof for g_{MSE} is complete. The proof for g_{BA} can be obtained in a similar way by noting that, using (6),

$$\begin{aligned} 0 &= n^{\nu/(r+\nu+1)} B(g_{\text{BA}}) \\ &= \left\{ (n^{1/(r+\nu+1)} g_{\text{BA}})^{-r-1} L^{(r)}(0) + (n^{1/(r+\nu+1)} g_{\text{BA}})^\nu m_\nu(L) \psi_{r+\nu} / \nu! \right\} (1 + o(1)). \end{aligned}$$

b) From Lemma 2 and equalities (3) and (4) the functions $B(g)$ and $V(g)$, and therefore $\text{MSE}(g)$, are differentiable with

$$B'(g) = -(r+1)n^{-1}g^{-r-2}L^{(r)}(0) + \nu g^{\nu-1}\psi_{r+\nu}m_\nu(L)/\nu! + o(g^{\nu-1}) \quad (16)$$

and

$$V'(g) = 2c_{1,r}n^{-1}g^{\nu-1}m_\nu(L)/\nu! - 2c_{2,r}n^{-2}g^{-2r-2} + o(n^{-1}g^{\nu-1} + n^{-2}g^{-2r-2}), \quad (17)$$

with $c_{1,r} = 4\nu(\psi_{r+\nu,r} - \psi_{r+\nu}\psi_r)$ and $c_{2,r} = (2r+1)\psi_0 \int \{L^{(r)}\}^2$.

From these expansions together with (6), part a) of this result and equation $\text{MSE}'(g_{\text{MSE}}) = 2B(g_{\text{MSE}})B'(g_{\text{MSE}}) + V'(g_{\text{MSE}}) = 0$ we obtain

$$\begin{aligned} & ng_{\text{MSE}}^{r+1}B(g_{\text{MSE}})ng_{\text{MSE}}^{r+2}B'(g_{\text{MSE}}) \\ &= -n^2g_{\text{MSE}}^{2r+3}V'(g_{\text{MSE}})/2 \\ &= -c_{1,r}ng_{\text{MSE}}^{2r+\nu+2}m_\nu(L)/\nu! + c_{2,r}g_{\text{MSE}} + o(ng_{\text{MSE}}^{2r+\nu+2} + g_{\text{MSE}}) \end{aligned} \quad (18)$$

where

$$ng_{\text{MSE}}^{r+1}B(g_{\text{MSE}}) = L^{(r)}(0) + ng_{\text{MSE}}^{r+\nu+1}\psi_{r+\nu}m_\nu(L)/\nu! + o(1) \quad (19)$$

and

$$\begin{aligned} ng_{\text{MSE}}^{r+2}B'(g_{\text{MSE}}) &= -(r+1)L^{(r)}(0) + \nu ng_{\text{MSE}}^{r+\nu+1}\psi_{r+\nu}m_\nu(L)/\nu! + o(1) \\ &= (-1)^{r/2+1}\{(r+1)|L^{(r)}(0)| + \nu ng_{\text{MSE}}^{r+\nu+1}|\psi_{r+\nu}||m_\nu(L)|/\nu!\} + o(1). \end{aligned}$$

is such that $\liminf ng_{\text{MSE}}^{r+2}|B'(g_{\text{MSE}})| > 0$. Therefore, from (18) we finally get

$$L^{(r)}(0) + ng_{\text{MSE}}^{r+\nu+1}\psi_{r+\nu}m_\nu(L)/\nu! = o(1), \quad (20)$$

that concludes the proof for g_{MSE} . Also, notice that from $0 = ng_{\text{BA}}^{r+1}B(g_{\text{BA}})$ and (6) we obtain the same formula as in (20) with g_{BA} instead of g_{MSE} and thus the limit $g_0/g_{\text{BA}} \rightarrow 1$ and, consequently, $g_{\text{BA}}/g_{\text{MSE}} \rightarrow 1$.

c) Using the fact that f has a bounded derivative of order $r + \nu + 2$, from Lemma 2 we know that the residual term $o(g^\nu)$ appearing in the expansion of $B(g)$ can be replaced by $O(g^{\nu+2})$. This enables us to improve the order of convergence of the residual term in equation (19) which can be replaced by $O(g_{\text{MSE}}^2) = o(g_{\text{MSE}})$. Using again equation (18) and the fact that $ng_{\text{MSE}}^{r+2}B'(g_{\text{MSE}}) = -c_{3,r}(1 + o(1))$, where $c_{3,r} = (r + \nu + 1)L^{(r)}(0)$, we get

$$L^{(r)}(0) + ng_{\text{MSE}}^{r+\nu+1}\psi_{r+\nu}m_\nu(L)/\nu! = c_{1,r}c_{3,r}^{-1}ng_{\text{MSE}}^{2r+\nu+2}m_\nu(L)/\nu! - c_{2,r}c_{3,r}^{-1}g_{\text{MSE}} + o(g_{\text{MSE}}).$$

Taking into account that g_0 satisfies the equality

$$L^{(r)}(0) + ng_0^{r+\nu+1}\psi_{r+\nu}m_\nu(L)/\nu! = 0, \quad (21)$$

for some \bar{g} between g_0 and g_{MSE} we have

$$n(r + \nu + 1)\bar{g}^{r+\nu}(g_0/g_{\text{MSE}} - 1)\psi_{r+\nu}m_\nu(L)/\nu! = -c_{1,r}c_{3,r}^{-1}ng_{\text{MSE}}^{2r+\nu+1}m_\nu(L)/\nu! + c_{2,r}c_{3,r}^{-1} + o(1).$$

In order to conclude it suffices to remark that $n^{1/(r+\nu+1)}\bar{g} = c_{0,r}(1 + o(1))$ where $c_{0,r}^{r+\nu+1} = -\nu!L^{(r)}(0)/(m_\nu(L)\psi_{r+\nu})$. Therefore, the announced convergence for $g_0/g_{\text{MSE}} - 1$ takes place with

$$C = C_{L,r,\nu}(f) = -c_{0,r}c_{3,r}^{-2}\{c_{2,r} + 4\nu(\psi_{\nu,0}\psi_\nu^{-1} - \psi_0)L(0)\delta_{r0}\}, \quad (22)$$

where δ_{r0} is the Kronecker symbol, that is, $\delta_{r0} = 1$ for $r = 0$ and $\delta_{r0} = 0$ for $r \neq 0$.

On the other hand, starting from $0 = ng_{\text{BA}}^{r+1}B(g_{\text{BA}})$ and using (6) with the residual term $o(g^\nu)$ replaced by $O(g^{\nu+2})$ we come to

$$L^{(r)}(0) + ng_{\text{BA}}^{r+\nu+1}\psi_{r+\nu}m_\nu(L)/\nu! = g_{\text{BA}}^{r+1}\psi_r + O(ng_{\text{BA}}^{r+\nu+3}). \quad (23)$$

Reasoning as before we conclude that the announced convergence for $g_{\text{BA}}/g_{\text{MSE}} - 1$ takes place with

$$D = D_{L,r,\nu}(f) = -c_{0,r}c_{3,r}^{-2}\{c_{2,r} + (4\nu(\psi_{\nu,0}\psi_\nu^{-1} - \psi_0) + (\nu + 1)\psi_0)L(0)\delta_{r0}\}. \quad (24)$$

Finally, using the fact that f has a bounded continuous derivative of order $r + \nu + 2$, from Lemma 2 we know that the residual term $o(g^\nu)$ appearing in the expansion of $B(g)$ can be replaced by $g^{\nu+2}\psi_{r+\nu+2}m_{\nu+2}(L)/(\nu + 2)! + o(g^{\nu+2})$ which enables us write the residual term in equation (23) more precisely. Together with equation (21) we conclude that the announced convergence for $g_0/g_{\text{BA}} - 1$ takes place with

$$E = E_{L,r,\nu}(f) = -c_{0,r}c_{3,r}^{-1}\{c_{0,r}^{r+\nu+2}\psi_{r+\nu+2}m_{\nu+2}(L)/(\nu + 2)!(1 - \delta_{r0}) - \psi_0\delta_{r0}\}. \quad (25)$$

□

The orders of convergence for the higher order derivatives of $R_{L,r,g}(f)$, $R_{\{L^{(r)}\}^2,0,g}(f)$ and $S_{L,r,g}(f)$ given in the next lemma will be used in the proof of Theorem 6. They follow directly from expressions (9) and (10), the differentiation theorem under the integral sign and standard Taylor expansions.

Lemma 3. *Under assumptions (L1), (L2) and (D1), if f has bounded and continuous derivatives up to order $r + \nu + 2$, $|m_{\nu+2}|(L) < \infty$ and $\int |u|^3\{L^{(r)}(u)\}^2 du < \infty$, then the functions $g \mapsto R_{L,r,g}(f)$, $g \mapsto R_{\{L^{(r)}\}^2,0,g}(f)$ and $g \mapsto S_{L,r,g}(f)$ are three-times differentiable with*

$$\begin{aligned} d^2 R_{L,r,g}(f)/dg^2 &= O(g^{\nu-2}), & d^3 R_{L,r,g}(f)/dg^3 &= O(g^{\nu-3}), \\ d^2 R_{\{L^{(r)}\}^2,0,g}(f)/dg^2 &= O(1), & d^3 R_{\{L^{(r)}\}^2,0,g}(f)/dg^3 &= O(1), \\ d^2 S_{L,r,g}(f)/dg^2 &= O(g^{\nu-2}), & d^3 S_{L,r,g}(f)/dg^3 &= O(g^{\nu-3}). \end{aligned}$$

Proof of Theorem 6. From Lemmas 2 and 3 and equalities (3) and (4) the functions $B(g)$ and $V(g)$, and therefore $\text{MSE}(g)$, are three-times differentiable with

$$B''(g) = O(n^{-1}g^{-r-3} + g^{\nu-2}), \quad B'''(g) = O(n^{-1}g^{-r-4} + g^{\nu-3})$$

and

$$V''(g) = O(n^{-2}g^{-2r-3} + n^{-1}g^{\nu-2}), \quad V'''(g) = O(n^{-2}g^{-2r-4} + n^{-1}g^{\nu-3}).$$

Moreover, a Taylor expansion for $g \mapsto \text{MSE}(g)$ around $g = g_{\text{BA}}$ leads to

$$\begin{aligned} \text{MSE}(g_0) - \text{MSE}(g_{\text{BA}}) &= \text{MSE}'(g_{\text{BA}})g_{\text{BA}}(g_0/g_{\text{BA}} - 1) + \text{MSE}''(g_{\text{BA}})g_{\text{BA}}^2(g_0/g_{\text{BA}} - 1)^2/2 \\ &\quad + \text{MSE}'''(\tilde{g})g_{\text{BA}}^3(g_0/g_{\text{BA}} - 1)^3/3!, \end{aligned}$$

for some \tilde{g} between g_0 and g_{BA} . Taking into account that $B(g_{\text{BA}}) = 0$ and $n^{1/(r+\nu+1)}g_{\text{BA}} = c_{0,r}(1 + o(1))$, from the previous orders of convergence for $B''(g)$, $B'''(g)$, $V''(g)$ and $V'''(g)$, the expansions (16) and (17) for $B'(g)$ and $V'(g)$, respectively, and Theorem 5.c), we get

$$\text{MSE}'(g_{\text{BA}})g_{\text{BA}} = c_{0,r}^{-2r-1}d_r n^{-(2\nu+1)/(r+\nu+1)}(1 + o(1)),$$

$$\text{MSE}''(g_{\text{BA}})g_{\text{BA}}^2(g_0/g_{\text{BA}} - 1) = 2c_{0,r}^{-2r-2}c_{3,r}^2 E n^{-\min(r+2\nu+1, 2\nu+2)/(r+\nu+1)}(1 + o(1))$$

and

$$\text{MSE}'''(\tilde{g})g_{\text{BA}}^3(g_0/g_{\text{BA}} - 1)^2 = O(n^{-(2\nu+2)/(r+\nu+1)}),$$

where $d_r = -2\{c_{2,r} + 4\nu(\psi_{\nu,0}\psi_\nu^{-1} - \psi_0)L(0)\delta_{r0}\}$ and the constants $c_{0,r}$, $c_{2,r}$ and $c_{3,r}$ are defined in the proof of Theorem 5. Therefore, from Theorem 5.c), the announced convergence for $\text{MSE}(g_0) - \text{MSE}(g_{\text{BA}})$ will take place with

$$\Lambda = \Lambda_{L,r,\nu}(f) = c_{0,r}^{-2r-2}\{c_{0,r}d_r + c_{3,r}^2 E \delta_{r0}\}. \quad (26)$$

□

Proof of Theorem 7. As noted previously, to obtain an explicit formula for the MSE function we just need to provide explicit formulas for $L^{(r)}(0)$, $R_{L,r,g}(f)$, ψ_r , $S_{L,r,g}(f)$ and $R_{\{L^{(r)}\}^2,0,g}(f)$. From Fact C.1.6 in Appendix C in Wand and Jones (1995) we already have

$$L^{(r)}(0) = (-1)^{r/2}(2\pi)^{-1/2} \sum_{s=0}^{\nu/2-1} (-1)^s a_s \text{OF}(2s+r).$$

Also, using Fact C.1.12 there, taking into account that r is even,

$$\begin{aligned} \psi_r &= \int f^{(r)}(x)f(x)dx = \sum_{\ell,\ell'=1}^k w_\ell w_{\ell'} \int \phi_{\sigma_\ell}^{(r)}(x - \mu_\ell) \phi_{\sigma_{\ell'}}(x - \mu_{\ell'}) dx \\ &= \sum_{\ell,\ell'=1}^k w_\ell w_{\ell'} \phi_{\sigma_{\ell\ell'}}^{(r)}(\mu_{\ell\ell'}). \end{aligned}$$

And from the same result and Fact C.1.11 we have

$$\begin{aligned}
R_{L,r,g}(f) &= \int L_g^{(r)} * f(x) f(x) dx \\
&= \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \sum_{s=0}^{\nu/2-1} a_s g^{2s} \int (\phi_g^{(2s+r)} * \phi_{\sigma_\ell})(x - \mu_\ell) \phi_{\sigma_{\ell'}}(x - \mu_{\ell'}) dx \\
&= \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \sum_{s=0}^{\nu/2-1} a_s g^{2s} \int \phi_{\sigma_\ell(g)}^{(2s+r)}(x - \mu_\ell) \phi_{\sigma_{\ell'}}(x - \mu_{\ell'}) dx \\
&= \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \sum_{s=0}^{\nu/2-1} a_s g^{2s} \phi_{\sigma_{\ell\ell'}(g)}^{(2s+r)}(\mu_{\ell\ell'})
\end{aligned}$$

and so, the formula for the bias is complete.

On the other hand, Theorem 6.1 in Aldershof *et al.* (1995) with $m = 3$ and $r_3 = 0$ states that

$$\int \phi_{\sigma_1}^{(r_1)}(x - \mu_1) \phi_{\sigma_2}^{(r_2)}(x - \mu_2) \phi_{\sigma_3}(x - \mu_3) dx = I_{r_1, r_2}(\boldsymbol{\mu}; \boldsymbol{\sigma}). \quad (27)$$

But we have $L_g^{(r)} * f(x) = \sum_{\ell=1}^k w_\ell \sum_{s=0}^{\nu/2-1} a_s g^{2s} \phi_{\sigma_\ell(g)}^{(2s+r)}(x - \mu_\ell)$ so that from (27) we obtain

$$\begin{aligned}
S_{L,r,g}(f) &= \int \{L_g^{(r)} * f(x)\}^2 f(x) dx \\
&= \sum_{\ell_1, \ell_2, \ell_3=1}^k w_{\ell_1} w_{\ell_2} w_{\ell_3} \sum_{s, s'=0}^{\nu/2-1} a_s a_{s'} g^{2s+2s'} \\
&\quad \times \int \phi_{\sigma_{\ell_1}(g)}^{(2s+r)}(x - \mu_{\ell_1}) \phi_{\sigma_{\ell_2}(g)}^{(2s'+r)}(x - \mu_{\ell_2}) \phi_{\sigma_{\ell_3}}(x - \mu_{\ell_3}) dx \\
&= \sum_{\ell_1, \ell_2, \ell_3=1}^k w_{\ell_1} w_{\ell_2} w_{\ell_3} \sum_{s, s'=0}^{\nu/2-1} a_s a_{s'} g^{2s+2s'} \\
&\quad \times I_{2s+r, 2s'+r}(\mu_{\ell_1}, \mu_{\ell_2}, \mu_{\ell_3}; \sigma_{\ell_1}(g), \sigma_{\ell_2}(g), \sigma_{\ell_3}).
\end{aligned}$$

For the remaining term, it is easy to check that

$$f * \bar{f}(z) = \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \phi_{\sigma_{\ell\ell'}}(z - \mu_{\ell\ell'})$$

and we also know that

$$g^{-2r-1} R_{\{L^{(r)}\}^2, 0, g}(f) = \iint \{L_g^{(r)}(x - y)\}^2 f(x) f(y) dx dy = \int \{L_g^{(r)}(z)\}^2 (f * \bar{f})(z) dz$$

so that in the normal mixture case we have

$$g^{-2r-1} R_{\{L^{(r)}\}^2, 0, g}(f) = \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \sum_{s, s'=0}^{\nu/2-1} a_s a_{s'} g^{2s+2s'} \int \phi_g^{(2s+r)}(z) \phi_g^{(2s'+r)}(z) \phi_{\sigma_{\ell\ell'}}(z - \mu_{\ell\ell'}) dz$$

and we conclude by using formula (27) with $\mu_1 = \mu_2 = 0$, $\mu_3 = \mu_{\ell\ell'}$ and $\sigma_1 = \sigma_2 = g$, $\sigma_3 = \sigma_{\ell\ell'}$. \square

Acknowledgement. This research has been partially supported by the Spanish Ministerio de Ciencia y Tecnología project MTM2010-16660 and by the CMUC (Centre for Mathematics, University of Coimbra)/FCT.

References

- Aldershof, B. (1991) *Estimation of Integrated Squared Density Derivatives*. Ph.D. thesis, University of North Carolina, Chapel Hill.
- Aldershof, B., Marron, J. S., Park, B. U. and Wand, M. P. (1995). Facts about the Gaussian probability density function. *Appl. Anal.*, **59**, 289–306.
- Bhattacharya, G. K. and Roussas, G. G. (1969) Estimation of a certain functional of a probability density function. *Skand. Aktuarietidskr*, **1969**, 201–206.
- Bickel, P. J. and Ritov, Y. (1988) Estimating integrated squared density derivatives: sharp best order of convergence estimates. *Sankhya Ser. A*, **50**, 381–393.
- Butzer, P.L. and Nessel, R.J. (1971) *Fourier Analysis and Approximation, Vol. 1*. Birkhäuser Verlag, Basel.
- Chacón, J.E., Montanero, J. and Nogales, A.G. (2007a) A note on kernel density estimation at a parametric rate. *J. Nonparametr. Stat.*, **19**, 13–21.
- Chacón, J.E., Montanero, J., Nogales, A.G. and Pérez, P. (2007b) On the existence and limit behavior of the optimal bandwidth for kernel density estimation. *Statist. Sinica*, **17**, 289–300.
- Chacón, J.E. and Tenreiro, C. (2011) Data-based choice of the number of pilot stages for plug-in bandwidth selection. To appear in *Comm. Statist. Theory Methods*. doi:10.1080/03610926.2011.606486
- Dmitriev, Y.G. and Tarasenko, F.P. (1973) On the estimation of functionals of the probability density and its derivatives. *Theory Probab. Appl.*, **18**, 628–633.
- Dmitriev, Y.G. and Tarasenko, F.P. (1975) On a class of non-parametric estimates of non-linear functionals of density. *Theory Probab. Appl.*, **19**, 390–394.
- Giné, E. and Mason D.M. (2008) Uniform in bandwidth estimation of integral functionals of the density function. *Scand. J. Statist.*, **35**, 739–761.

- Giné, E. and Nickl, R. (2008) A simple adaptive estimator of the integrated square of a density. *Bernoulli*, **14**, 47–61.
- Hall, P. and Marron, J.S. (1987) Estimation of integrated squared density derivatives. *Statist. Probab. Lett.*, **6**, 109–115.
- Hettmansperger, T.P. (1984) *Statistical Inference Based on Ranks*. Wiley, New York.
- Hewitt, E. and Stromberg, K. (1965) *Real and Abstract Analysis*. Springer-Verlag, Berlin.
- Jones, M.C. and Sheather, S.J. (1991) Using non-stochastic terms to advantage in kernel-based estimation of integrated squared density derivatives. *Statist. Probab. Lett.*, **11**, 511–514.
- Laurent, B. (1997) Estimation of integral functionals of a density and its derivatives. *Bernoulli*, **3**, 181–211.
- Levit, B. Y. (1978) Asymptotically efficient estimation of nonlinear functionals. *Problems Inform. Transmission*, **14**, 65–72.
- Marron, J. S. and Wand, M. P. (1992) Exact mean integrated squared error. *Ann. Statist.*, **20**, 712–736.
- Politis, D. N. and Romano, J. P. (1999) Multivariate density estimation with general flat-top kernels of infinite order. *J. Multivariate Anal.*, **68**, 1–25.
- Prakasa Rao, B.L.S. (1983) *Nonparametric Functional Estimation*. Academic Press, New York.
- Prakasa Rao, B.L.S. (1999) Estimation of the integrated squared density derivatives by wavelets. *Bull. Inform. Cybernet.*, **31**, 47–65.
- Scott, D.W. (1992) *Multivariate Density Estimation: Theory, Practice and Visualization*. Wiley, New York.
- Schuster, E.F. (1974) On the rate of convergence of an estimate of a functional of a probability density. *Scand. Actuar. J.*, **1974**, 103–107.
- Sheather, S.J., Hettmansperger, T.P. and Donald, M.R. (1994) Data-based bandwidth selection for kernel estimators of the integral of $f^2(x)$. *Scand. J. Statist.*, **21**, 265–275.
- Tenreiro, C. (2003) On the asymptotic normality of multistage integrated density derivatives kernel estimators. *Statist. Probab. Lett.*, **64**, 311–322.
- van Es, B. (1992) Estimating functionals related to a density by a class of statistics based on spacings. *Scand. J. Statist.*, **19**, 61–72.

- Wand, M.P. and Jones, M.C. (1995) *Kernel Smoothing*. Chapman and Hall, London.
- Wand, M.P. and Schucany, W.R. (1990) Gaussian-based kernels. *Canad. J. Statist.*, **18**, 197–204.
- Wu, T.-J. (1995) Adaptive root n estimates of integrated squared density derivatives. *Ann. Statist.*, **23**, 1474–1495.