

ASYMPTOTICS FOR GENERAL MULTIVARIATE KERNEL DENSITY DERIVATIVE ESTIMATORS

José E. Chacón¹, Tarn Duong^{2,3} and M. P. Wand⁴

¹*Universidad de Extremadura*, ²*Institut Pasteur*, ³*Institut Curie*
and ⁴*University of Wollongong*

Abstract: We investigate kernel estimators of multivariate density derivative functions using general (or unconstrained) bandwidth matrix selectors. These density derivative estimators have been relatively less well researched than their density estimator analogues. A major obstacle for progress has been the intractability of the matrix analysis when treating higher order multivariate derivatives. With an alternative vectorization of these higher order derivatives, mathematical intractabilities are surmounted in an elegant and unified framework. The finite sample and asymptotic analysis of squared errors for density estimators are generalized to density derivative estimators. Moreover, we are able to exhibit a closed form expression for a normal scale bandwidth matrix for density derivative estimators. These normal scale bandwidths are employed in a numerical study to demonstrate the gain in performance of unconstrained selectors over their constrained counterparts.

Key words and phrases: Asymptotic mean integrated squared error, normal scale rule, optimal, unconstrained bandwidth matrices.

1. Introduction

Estimating probability density functions with kernel functions has had notable success due to their ease of interpretation and visualization. On the other hand, estimating derivatives of density functions has received less attention. This is partly because it is a more challenging theoretical problem (especially for multivariate data). Nonetheless there remains much information about the structure of a density function that is not easily ascertained from the density function itself, for example, local maxima and minima. One of the original papers on kernel density estimation (Parzen (1962)) was also concerned with estimating the global mode of the density function, though not from a density derivative point of view. It can be recast as finding the local maxima via derivative estimates, the global mode follows as the largest of these local maxima. The focus on a mode as a single point can be extended to the region immediately surrounding the mode, known as a bump or modal region. Modal regions can be used to determine the existence of multi-modality and/or clusters. The Godtlielsen, Marron, and Chaudhuri (2002) feature significance technique for bump-hunting relies on estimating and characterizing the first and second derivatives for bivariate data. In

an econometrics setting, the Engel curve describes the demand for a good/service as a function of income; it classifies goods/services based on the slope of the Engel curve and the first derivative is an essential component for interpreting them, see Hildenbrand and Hildenbrand (1986). In a more general setting, Singh (1977) suggests, as an application of multivariate density derivatives, to estimate the Fisher information matrix in its translation parameter form.

The first paper to be concerned with univariate kernel density derivative estimation appears to have been Bhattacharya (1967), followed by Schuster (1969) and Singh (1979, 1987). Singh (1976) studied a multivariate estimator with a diagonal bandwidth matrix, while Härdle, Marron, and Wand (1990) used a bandwidth parametrized as a constant times the identity matrix. This previous research has mostly focused on constrained parametrizations of the bandwidth matrix since this simplifies the matrix analysis compared to the general, unconstrained parametrization. Analyzing squared error measures for general kernel density derivative estimators has reached an impasse using the traditional vectorization of higher order derivatives of multivariate functions to transform the higher order derivative tensor into a more tractable vector form. Here we introduce an alternative vectorization of higher order derivatives, a subtle rearrangement of the traditional vectorization that allows us to easily write down the usual error expressions from the density estimation case. Thus we generalize the squared error analysis for kernel density derivative estimators. Furthermore, we are able to write down a closed form expression for a normal scale bandwidth matrix, i.e., the optimal bandwidth for the r th order derivative of a normal density with normal kernels. Normal scale selectors were the step that eventually led to the development of the now widely used bandwidth selectors for density estimation; we set up a similar starting point for future bandwidth selection in density derivative estimation.

In Section 2 we define a kernel estimator of a multivariate density derivative, and provide the results for mean integrated square convergence both asymptotically and for finite samples. The influence of the bandwidth matrix on convergence is established here. In Section 3 we focus on the class of normal mixture densities. Estimation of these densities with normal kernels produces further simplified special cases of the results in the Section 2, where we develop a normal scale bandwidth selector. We use these normal scale selectors to quantify the improvement when using unconstrained matrices in asymptotic performance in Section 4, and in finite sample performance in Section 5. We illustrate normal scale selectors on data arising from high throughput biotechnology in Section 6. The usual normal scale selectors based on the density function may lead to insufficient smoothing when estimating the density curvature (or second derivative). We conclude with a discussion in Section 7.

2. Kernel Density Derivative Estimation

Multivariate kernel density estimation has reached maturity, and recent advances there can be carried over to the density derivative case. To proceed, we use the linearity of the kernel density estimator to define a kernel density derivative estimator. The usual performance measure for kernel density estimation is the mean integrated squared error (MISE), which is easily extended to cover density derivatives.

We need some notation. For a matrix \mathbf{A} , let

$$\mathbf{A}^{\otimes r} = \bigotimes_{i=1}^r \mathbf{A} = \overbrace{\mathbf{A} \otimes \cdots \otimes \mathbf{A}}^{r \text{ matrices}}$$

denote the r th Kronecker power of \mathbf{A} . If $\mathbf{A} \in \mathcal{M}_{m \times n}$ (i.e., \mathbf{A} is a matrix of order $m \times n$) then $\mathbf{A}^{\otimes r} \in \mathcal{M}_{m^r \times n^r}$; we adopt the convention $\mathbf{A}^{\otimes 1} = \mathbf{A}$ and $\mathbf{A}^{\otimes 0} = \mathbf{1} \in \mathbb{R}$.

If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a real function of a vector variable, let $\mathbf{D}^{\otimes r} f(\mathbf{x}) \in \mathbb{R}^{d^r}$ be the vector containing all the partial derivatives of order r of f at \mathbf{x} , arranged so that we can formally write

$$\mathbf{D}^{\otimes r} f = \frac{\partial f}{(\partial \mathbf{x})^{\otimes r}}.$$

Thus we write the r th derivative of f as a vector of length d^r , and not as an r -fold tensor array or as a matrix. Moreover, if $f: \mathbb{R}^d \rightarrow \mathbb{R}^p$ is a vector function of a vector variable with components $f = (f_1, \dots, f_p)$, then we set

$$\mathbf{D}^{\otimes r} f(\mathbf{x}) = \begin{pmatrix} \mathbf{D}^{\otimes r} f_1(\mathbf{x}) \\ \vdots \\ \mathbf{D}^{\otimes r} f_p(\mathbf{x}) \end{pmatrix}.$$

Notice that, using this notation, we have $\mathbf{D}(\mathbf{D}^{\otimes r} f) = \mathbf{D}^{\otimes(r+1)} f$. Also, the gradient of f is just $\mathbf{D}^{\otimes 1} f$ and the Hessian $\mathbf{H}f = \partial^2 f / (\partial \mathbf{x} \partial \mathbf{x}^T)$ is such that $\text{vec } \mathbf{H}f = \mathbf{D}^{\otimes 2} f$, where vec denotes the vector operator (see Henderson and Searle (1979)). This vectorization carries some redundancy, for example, $\mathbf{D}^{\otimes 2} f$ contains repeated mixed partial derivatives whereas the usual vectorization $\text{vech } \mathbf{H}f$ contains only the unique second order partial derivatives, with vech denoting the vector half operator (see Henderson and Searle (1979)). The latter is usually preferred since it is minimal and its matrix analysis is not more complicated than the former. However for $r > 2$, it appears that the matrix analysis using a generalization of the vector half operator becomes intractable. Other authors have used the same vectorization we propose to develop results for higher order derivatives: Holmquist (1996a) computes derivatives of the multivariate normal

density function, and Chacón and Duong (2010) compute kernel estimators of multivariate integrated density derivative functionals.

Suppose now that $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a density and we wish to estimate $D^{\otimes r} f(\mathbf{x})$ for some $r \in \mathbb{N}$. To this end, we use the kernel estimator $\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H}) = D^{\otimes r} \widehat{f}(\mathbf{x}; \mathbf{H})$ where, given a random sample $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_d$ drawn from f ,

$$\widehat{f}(\mathbf{x}; \mathbf{H}) = n^{-1} \sum_{i=1}^n K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i) \quad (2.1)$$

denotes the multivariate kernel density estimator with kernel K and bandwidth matrix \mathbf{H} , with $K_{\mathbf{H}}(\mathbf{x}) = |\mathbf{H}|^{-1/2} K(\mathbf{H}^{-1/2} \mathbf{x})$. Conditions on K and \mathbf{H} are given later. Notice that alternative expressions for $\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H})$ are

$$\begin{aligned} \widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H}) &= n^{-1} \sum_{i=1}^n D^{\otimes r} K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i) \\ &= n^{-1} (\mathbf{H}^{-1/2})^{\otimes r} \sum_{i=1}^n (D^{\otimes r} K)_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i) \\ &= n^{-1} |\mathbf{H}|^{-1/2} (\mathbf{H}^{-1/2})^{\otimes r} \sum_{i=1}^n D^{\otimes r} K(\mathbf{H}^{-1/2}(\mathbf{x} - \mathbf{X}_i)). \end{aligned} \quad (2.2)$$

The last expression here is quite helpful for implementing the estimator since it separates the roles of K and \mathbf{H} . This is the multivariate generalization of the kernel estimator that appears, for instance, in Härdle, Marron, and Wand (1990) and Jones (1994).

Here the most general (unconstrained) form of the bandwidth matrix is used. In contrast, earlier papers considered this estimator, but with constrained parametrizations of the bandwidth matrix: Härdle, Marron, and Wand (1990) used a parametrization where \mathbf{H} is h^2 multiplied by the identity matrix; Singh (1976) used $\mathbf{H} = \text{diag}(h_1^2, \dots, h_d^2)$. However, in the case $r = 0$ (estimation of the density itself), Wand and Jones (1993) gave examples that show that a significant gain may be achieved by using unconstrained parametrizations over constrained ones; see also Chacón (2009). We generalize these results to arbitrary derivatives in Section 4.

We measure the error of the kernel density derivative estimator at the point \mathbf{x} by using the mean squared error

$$\text{MSE}(\mathbf{x}; \mathbf{H}) \equiv \text{MSE}\{\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H})\} = \mathbb{E}\|\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H}) - D^{\otimes r} f(\mathbf{x})\|^2,$$

with $\|\cdot\|$ the Euclidean norm. We have the two forms $\text{MSE}(\mathbf{x}; \mathbf{H})$ and $\text{MSE}\{\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H})\}$, depending on whether we wish to suppress the explicit dependence on $\widehat{D^{\otimes r} f}$ or not. It is easy to check that we can write $\text{MSE}(\mathbf{x}; \mathbf{H}) =$

$B^2(\mathbf{x}; \mathbf{H}) + V(\mathbf{x}; \mathbf{H})$, where

$$B^2(\mathbf{x}; \mathbf{H}) = \|\mathbb{E}\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H}) - D^{\otimes r} f(\mathbf{x})\|^2,$$

$$V(\mathbf{x}; \mathbf{H}) = \mathbb{E}\|\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H}) - \mathbb{E}\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H})\|^2.$$

Analogously, as a global measure of the performance of the estimator, we use the mean integrated squared error

$$\text{MISE}(\mathbf{H}) \equiv \text{MISE}\{\widehat{D^{\otimes r} f}(\cdot; \mathbf{H})\} = \int \text{MSE}\{\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H})\} d\mathbf{x}.$$

Our results rely on the following assumptions on the bandwidth matrix, the density function and the kernel function.

- (A1) \mathbf{H} , the bandwidth matrix, is symmetric and positive-definite, and such that every element of $\mathbf{H} \rightarrow 0$ and $n^{-1}|\mathbf{H}|^{-1/2}(\mathbf{H}^{-1})^{\otimes r} \rightarrow 0$ as $n \rightarrow \infty$.
- (A2) f is a density function for which all partial derivatives to order $r + 2$ exist, all its partial derivatives of order r are square integrable, and all its partial derivatives of order $r + 2$ are bounded, continuous, and square integrable.
- (A3) K is a positive, symmetric, square integrable density function such that $\int \mathbf{x}\mathbf{x}^T K(\mathbf{x})d\mathbf{x} = m_2(K)\mathbf{I}_d$ for some real number $m_2(K)$ with \mathbf{I}_d the identity matrix of order d , and all its partial derivatives of order r are square integrable.

This is not the minimal set of assumptions required, but it provides a useful starting point for quantifying the squared error results. We leave it to future research to reduce these assumptions.

For any vector function $g: \mathbb{R}^d \rightarrow \mathbb{R}^p$ let

$$\mathbf{R}(g) = \int g(\mathbf{x})g(\mathbf{x})^T d\mathbf{x} \in \mathcal{M}_{p \times p}.$$

For an arbitrary kernel L we write

$$\mathbf{R}_{L, \mathbf{H}, r}(f) = \int L_{\mathbf{H}} * D^{\otimes r} f(\mathbf{x})D^{\otimes r} f(\mathbf{x})^T d\mathbf{x} \in \mathcal{M}_{d^r \times d^r},$$

where the convolution operator is applied to each component of $D^{\otimes r} f(\mathbf{x})$ separately. Expanding all the terms in the bias-variance decomposition we obtain the following exact representation of the MISE function.

Theorem 1. *Suppose (A1)–(A3) hold. Then*

$$\begin{aligned} \text{MISE}\{\widehat{D^{\otimes r} f}(\cdot; \mathbf{H})\} &= \text{tr } \mathbf{R}(D^{\otimes r} f) + n^{-1}|\mathbf{H}|^{-1/2} \text{tr } ((\mathbf{H}^{-1})^{\otimes r} \mathbf{R}(D^{\otimes r} K)) \\ &\quad + (1 - n^{-1}) \text{tr } \mathbf{R}_{K * K, \mathbf{H}, r}(f) - 2 \text{tr } \mathbf{R}_{K, \mathbf{H}, r}(f). \end{aligned}$$

The proof of this, along with the proofs for other theorems in this paper, are deferred to the Appendix.

The form of the MISE given in the theorem involves a complicated dependence on the bandwidth matrix \mathbf{H} via the \mathbf{R} functionals. To show this dependence more clearly, we search for a more mathematically tractable form. The next result gives an asymptotic representation of the MISE function that can be viewed as an extension of (4.10) in Wand and Jones (1995, p. 98), which corresponds to the case $r = 0$.

Theorem 2. *Suppose (A1)–(A3) hold. Then*

$$\text{MISE}\{\widehat{\mathbf{D}^{\otimes r} f}(\cdot; \mathbf{H})\} = \text{AMISE}\{\widehat{\mathbf{D}^{\otimes r} f}(\cdot; \mathbf{H})\} + o(n^{-1}|\mathbf{H}|^{-1/2}\text{tr}^r(\mathbf{H}^{-1}) + \text{tr}^2\mathbf{H}),$$

where

$$\begin{aligned} \text{AMISE}\{\widehat{\mathbf{D}^{\otimes r} f}(\cdot; \mathbf{H})\} &= n^{-1}|\mathbf{H}|^{-1/2} \text{tr}((\mathbf{H}^{-1})^{\otimes r} \mathbf{R}(\mathbf{D}^{\otimes r} K)) \\ &\quad + \frac{m_2(K)^2}{4} \text{tr}[(\mathbf{I}_{dr} \otimes \text{vec}^T \mathbf{H}) \mathbf{R}(\mathbf{D}^{\otimes(r+2)} f)(\mathbf{I}_{dr} \otimes \text{vec} \mathbf{H})]. \end{aligned}$$

The AMISE-optimal bandwidth matrix $\mathbf{H}_{\text{AMISE}}$ is defined to be the matrix, amongst all symmetric positive definite matrices, that minimizes $\text{AMISE}\{\widehat{\mathbf{D}^{\otimes r} f}(\cdot; \mathbf{H})\}$. Next we give the order of such a matrix, and the order of the resulting minimal AMISE.

Theorem 3. *Suppose (A1)–(A3) hold. Every entry of the optimal bandwidth matrix $\mathbf{H}_{\text{AMISE}}$ is of order $O(n^{-2/(d+2r+4)})$, and $\min_{\mathbf{H}} \text{AMISE}\{\widehat{\mathbf{D}^{\otimes r} f}(\cdot; \mathbf{H})\}$ is of order $O(n^{-4/(d+2r+4)})$.*

Theorems 1, 2 and 3 generalize the existing MISE, AMISE, and $\mathbf{H}_{\text{AMISE}}$ results for constrained to unconstrained bandwidth matrices. These rates of convergence reveal the dependence on dimension d and derivative order r . As either of these increase, the minimum achievable AMISE increases. At least asymptotically, the marginal increase in difficulty in estimating a derivative an order higher is the same as estimating a density two dimensions higher.

To appreciate the ramifications of our theorems, we briefly review the literature for kernel density derivative estimation. Bhattacharya (1967) showed that the rate of convergence in probability of a univariate kernel density derivative estimator with a second order kernel is bounded by $n^{-1/(2r+4)}$, without developing intermediate squared error convergence results. Schuster (1969) established the same rate for a wider class of kernels. Singh (1979, 1987) showed, for his specially constructed univariate kernel estimator, that $n^{-2(p-r)/(2p+1)}$ is the MSE and MISE rate of convergence. This estimator employs kernels whose r th moment is 1 and all other j th moments are zero, $j = 0, 1, \dots, p-1$, $j \neq r$ for $p > r$;

the order of these kernels is greater than or equal to r . Since we assume second order kernels (assumption (A3)), Theorem 3 does not generalize this result for $r > 2$. For second order kernels, Wand and Jones (1995, p. 49) showed that $p = r + 2$ which gives the MISE to be $O(n^{-4/(2r+5)})$, which is indeed Theorem 3 with $d = 1$.

For d -variate density derivative estimation, Stone (1980) had that for any linear non-parametric estimator, the minimum achievable MSE of an estimator of g , a scalar functional of $\mathbf{D}^{\otimes r} f$, is $O(n^{-2(p-r)/(2p+d)})$, where $p > r$ and the p th order derivative of g is bounded. From (A2), we have $p = r + 2$ so the MISE rate is $O(n^{-4/(d+2r+4)})$, the rate in Theorem 3. However this general result cannot elucidate such questions as the relationship between the convergence rate and the bandwidth matrix, which Theorems 1, 2, and 3 demonstrate clearly. Singh (1976) showed that a multivariate kernel density derivative estimator with a diagonal bandwidth matrix $\mathbf{H} = \text{diag}(h_1^2, h_2^2, \dots, h_d^2)$ is mean square convergent, though the rate of convergence was not specified. More recently Duong et al. (2008) established that the MISE convergence rate for kernel estimators with unconstrained bandwidth matrices, is $n^{-4/(d+6)}$ and $n^{-4/(d+8)}$ for $r = 1, 2$. Theorem 3 extends these results to arbitrary r .

So far we have only considered scalar functionals of the expected value and the variance of $\widehat{\mathbf{D}^{\otimes r} f}$. For completeness, the following theorem gives these quantities in their vector and matrix form. This theorem is a generalization of the results obtained by Duong et al. (2008) for $r = 1, 2$ to arbitrary r .

Theorem 4. *Suppose (A1)–(A3) hold. Then*

$$\begin{aligned} \mathbb{E}\widehat{\mathbf{D}^{\otimes r} f}(\mathbf{x}; \mathbf{H}) &= \mathbf{D}^{\otimes r} f(\mathbf{x}) + \frac{1}{2} m_2(K)(\mathbf{I}_d \otimes \text{vec}^T \mathbf{H}) \mathbf{D}^{\otimes(r+2)} f(\mathbf{x}) + O(\text{tr}^2 \mathbf{H}) \mathbf{1}_d, \\ \text{Var} \widehat{\mathbf{D}^{\otimes r} f}(\mathbf{x}; \mathbf{H}) &= n^{-1} |\mathbf{H}|^{-1/2} (\mathbf{H}^{-1/2})^{\otimes r} \mathbf{R}(\mathbf{D}^{\otimes r} K) (\mathbf{H}^{-1/2})^{\otimes r} f(\mathbf{x}) \\ &\quad + o(n^{-1} |\mathbf{H}|^{-1/2} \text{tr}^r \mathbf{H}^{-1}) \mathbf{J}_d, \end{aligned}$$

where the elements of $\mathbf{1}_p \in \mathbb{R}^p$ and $\mathbf{J}_p \in \mathcal{M}_{p \times p}$ are all ones.

3. Normal Mixture Densities

In this section we study the normal mixture case in detail. We start with $K = \phi$ as the density of the standard d -variate normal, $\phi(\mathbf{x}) = (2\pi)^{-d/2} \exp(-\mathbf{x}^T \mathbf{x} / 2)$, and with $f = \phi_{\Sigma}(\cdot - \boldsymbol{\mu})$ a normal density with mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\Sigma}$.

The MISE and AMISE expressions in the normal case are closely related to the moments of quadratic forms in normal variables. Given two symmetric matrices \mathbf{A} and \mathbf{B} in $\mathcal{M}_{d \times d}$, we write $\mu_{r,s}(\mathbf{A}, \mathbf{B}) \equiv \mathbb{E}[(\mathbf{z}^T \mathbf{A}^{-1} \mathbf{z})^r (\mathbf{z}^T \mathbf{B}^{-1} \mathbf{z})^s]$ and $\mu_r(\mathbf{A}) \equiv \mu_{r,0}(\mathbf{A}, \mathbf{I})$, where \mathbf{z} is a d -variate random vector with standard normal distribution.

Theorem 5. *Suppose (A1) holds. Let f be a normal density with mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\Sigma}$, and K be the normal kernel. Then*

$$\begin{aligned} \text{MISE}\{\widehat{\text{D}}^{\otimes r} f(\cdot; \mathbf{H})\} &= 2^{-(d+r)} \pi^{-d/2} \{ |\boldsymbol{\Sigma}|^{-1/2} \mu_r(\boldsymbol{\Sigma}) + n^{-1} |\mathbf{H}|^{-1/2} \mu_r(\mathbf{H}) \\ &\quad + (1 - n^{-1}) |\mathbf{H} + \boldsymbol{\Sigma}|^{-1/2} \mu_r(\mathbf{H} + \boldsymbol{\Sigma}) \\ &\quad - 2^{(d+2r+2)/2} |\mathbf{H} + 2\boldsymbol{\Sigma}|^{-1/2} \mu_r(\mathbf{H} + 2\boldsymbol{\Sigma}) \}. \end{aligned}$$

An explicit form of the minimizer of the MISE is not available.

To rewrite the MISE without the μ_r functionals, we use Theorem 1 in Holmquist (1996b) to write

$$\mu_r(\mathbf{A}) = \text{OF}(2r) (\text{vec}^T \mathbf{A}^{-1})^{\otimes r} \mathcal{S}_{d,2r} (\text{vec} \mathbf{I}_d)^{\otimes r},$$

where $\text{OF}(2r) = (2r - 1)(2r - 3) \cdots 5 \cdot 3 \cdot 1$ denotes the odd factorial and $\mathcal{S}_{m,n} \in \mathcal{M}_{m^n \times m^n}$ is the symmetrizer matrix of order m, n ; see Holmquist (1996a,b). These references contain technical definitions of the symmetrizer matrix that we do not reproduce here. Instead, we focus on the action of the symmetrizer matrix on Kronecker products of vectors. Let $\mathbf{x}_i \in \mathbb{R}^p$, $i = 1, \dots, n$, and $\mathbf{X}^* = \{\mathbf{x}_1^*, \dots, \mathbf{x}_n^*\}$ be a permutation of $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. The symmetrizer matrix maps the product $\bigotimes_{i=1}^n \mathbf{x}_i$ to a linear combination of products of all possible permutations of $\mathbf{x}_1, \dots, \mathbf{x}_n$,

$$\mathcal{S}_{m,n} \left(\bigotimes_{i=1}^n \mathbf{x}_i \right) = \frac{1}{n!} \sum_{\text{all } \mathbf{X}^*} \bigotimes_{i=1}^n \mathbf{x}_i^*.$$

More explicitly for a 3-fold product, $\mathcal{S}_{m,3}(\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3) = (1/6)[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3 + \mathbf{x}_1 \otimes \mathbf{x}_3 \otimes \mathbf{x}_2 + \mathbf{x}_2 \otimes \mathbf{x}_1 \otimes \mathbf{x}_3 + \mathbf{x}_2 \otimes \mathbf{x}_3 \otimes \mathbf{x}_1 + \mathbf{x}_3 \otimes \mathbf{x}_1 \otimes \mathbf{x}_2 + \mathbf{x}_3 \otimes \mathbf{x}_2 \otimes \mathbf{x}_1]$.

Corollary 1. *Under the conditions of Theorem 5,*

$$\begin{aligned} \text{MISE}\{\widehat{\text{D}}^{\otimes r} f(\cdot; \mathbf{H})\} &= 2^{-(d+r)} \pi^{-d/2} \text{OF}(2r) \{ |\boldsymbol{\Sigma}|^{-1/2} (\text{vec}^T \boldsymbol{\Sigma}^{-1})^{\otimes r} + |\mathbf{H}|^{-1/2} (\text{vec}^T \mathbf{H}^{-1})^{\otimes r} \\ &\quad + (1 - n^{-1}) |\mathbf{H} + \boldsymbol{\Sigma}|^{-1/2} (\text{vec}^T (\mathbf{H} + \boldsymbol{\Sigma})^{-1})^{\otimes r} \\ &\quad - 2^{(d+2r+2)/2} |\mathbf{H} + 2\boldsymbol{\Sigma}|^{-1/2} (\text{vec}^T (\mathbf{H} + 2\boldsymbol{\Sigma})^{-1})^{\otimes r} \} \mathcal{S}_{d,2r} (\text{vec} \mathbf{I}_d)^{\otimes r}. \end{aligned}$$

This corollary shows the explicit dependence of the MISE on the bandwidth matrix \mathbf{H} . However, the direct computation of $\mathcal{S}_{d,2r} \in \mathcal{M}_{d^{2r} \times d^{2r}}$ may be an onerous task; for example, for $d = r = 4$, $\mathcal{S}_{4,8}$ is a 65536×65536 matrix. In contrast, although Theorem 5 does not express the explicit dependence of the MISE on \mathbf{H} due to the use of the μ_r functionals, formula (6) in Holmquist (1996b) gives the computationally efficient recursive expression

$$\mu_r(\mathbf{A}) = (r - 1)! 2^{r-1} \sum_{j=0}^{r-1} \frac{\text{tr}(\mathbf{A}^{-(r-j)})}{j! 2^j} \mu_j(\mathbf{A}). \tag{3.1}$$

Here, we understand $\mathbf{A}^{-p} = (\mathbf{A}^{-1})^p$ for $p > 0$ and, consequently, $\mathbf{A}^0 = \mathbf{I}_d$.

An analogous expression of the AMISE is given below. In this case, we can also give an explicit expression for its minimizer.

Theorem 6. *Under the conditions of Theorem 5,*

$$\begin{aligned} & \text{AMISE}\{\widehat{\mathbf{D}}^{\otimes r} f(\cdot; \mathbf{H})\} \\ &= 2^{-(d+r)} \pi^{-d/2} \{n^{-1} |\mathbf{H}|^{-1/2} \mu_r(\mathbf{H}) + \frac{1}{16} |\boldsymbol{\Sigma}|^{-1/2} \mu_{r,2}(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}^{1/2} \mathbf{H}^{-1} \boldsymbol{\Sigma}^{1/2})\}. \end{aligned}$$

The value of \mathbf{H} that minimizes this function is

$$\mathbf{H}_{\text{AMISE}} = \left(\frac{4}{d + 2r + 2} \right)^{2/(d+2r+4)} \boldsymbol{\Sigma} n^{-2/(d+2r+4)}.$$

The AMISE expression in Theorem 6 can be rewritten without μ_r functions. From Theorem 5 in Holmquist (1996b),

$$\begin{aligned} & \mu_{r,s}(\mathbf{A}, \mathbf{B}) \\ &= \text{OF}(2r + 2s) [(\text{vec}^T \mathbf{A}^{-1})^{\otimes r} \otimes (\text{vec}^T \mathbf{B}^{-1})^{\otimes s}] \mathcal{S}_{d,2r+2s}(\text{vec} \mathbf{I}_d)^{\otimes(r+s)}. \end{aligned}$$

Corollary 2. *Under the conditions of Theorem 5,*

$$\begin{aligned} & \text{AMISE}\{\widehat{\mathbf{D}}^{\otimes r} f(\cdot; \mathbf{H})\} \\ &= 2^{-(d+r)} \pi^{-d/2} \{ \text{OF}(2r) n^{-1} |\mathbf{H}|^{-1/2} (\text{vec}^T \mathbf{H}^{-1})^{\otimes r} \mathcal{S}_{d,2r}(\text{vec} \mathbf{I}_d)^{\otimes r} \\ & \quad + \frac{1}{16} |\boldsymbol{\Sigma}|^{-1/2} \text{OF}(2r + 4) [(\text{vec}^T \boldsymbol{\Sigma}^{-1})^{\otimes r} \otimes \{\text{vec}^T (\boldsymbol{\Sigma}^{-1/2} \mathbf{H} \boldsymbol{\Sigma}^{-1/2})\}^{\otimes 2}] \\ & \quad \times \mathcal{S}_{d,2r+4}(\text{vec} \mathbf{I}_d)^{\otimes(r+2)} \}. \end{aligned}$$

To facilitate the comparison of the extra amount of smoothing induced for higher dimensions and higher derivatives for standard normal densities, we examine the ratio

$$\frac{h_{\text{AMISE}}(d, r, n)}{h_{\text{AMISE}}(1, 0, n)} = \left(\frac{3}{4} \right)^{1/5} \left(\frac{4}{d + 2r + 2} \right)^{1/(d+2r+4)} n^{(d+2r-1)/(5d+10r+20)},$$

where $\mathbf{H}_{\text{AMISE}}(d, r, n) = h_{\text{AMISE}}^2(d, r, n) \mathbf{I}_d$ is the AMISE-optimal bandwidth for the standard normal density given d, r and n . These values are tabulated in Table 1 for $d = 1, \dots, 6$ and $r = 0, 1, 2, 3$. These ratios are the same for different d and r if $d + 2r$ is fixed.

The normal scale bandwidth selector is obtained by replacing the variance $\boldsymbol{\Sigma}$ in $\mathbf{H}_{\text{AMISE}}$ from Theorem 6 by an estimate $\hat{\boldsymbol{\Sigma}}$,

$$\hat{\mathbf{H}}_{\text{NS}} = \left(\frac{4}{d + 2r + 2} \right)^{2/(d+2r+4)} \hat{\boldsymbol{\Sigma}} n^{-2/(d+2r+4)}. \tag{3.2}$$

Table 1. Comparison of extra smoothing induced for higher dimensions and higher order derivatives for a standard normal density. Each table entry contains the value of $h_{\text{AMISE}}(d, r, n)/h_{\text{AMISE}}(1, 0, n)$.

		$d = 1$	$d = 2$	$d = 3$	$d = 4$	$d = 5$	$d = 6$
$n = 1,000$	$r = 0$	1.00	1.19	1.36	1.51	1.64	1.76
	$r = 1$	1.36	1.51	1.64	1.76	1.86	1.96
	$r = 2$	1.64	1.76	1.86	1.96	2.04	2.12
	$r = 3$	1.86	1.96	2.04	2.12	2.19	2.26
$n = 10,000$	$r = 0$	1.00	1.28	1.55	1.79	2.01	2.21
	$r = 1$	1.55	1.79	2.01	2.21	2.40	2.56
	$r = 2$	2.01	2.21	2.40	2.56	2.71	2.85
	$r = 3$	2.40	2.56	2.71	2.85	2.98	3.10
$n = 100,000$	$r = 0$	1.00	1.39	1.77	2.13	2.47	2.79
	$r = 1$	1.77	2.13	2.47	2.79	3.08	3.35
	$r = 2$	2.47	2.79	3.08	3.35	3.60	3.84
	$r = 3$	3.08	3.35	3.60	3.84	4.05	4.25

We can use (3.2) as a starting point to develop consistent data-driven bandwidth matrices. Consistent univariate selectors for density derivatives include the unbiased cross validation selector of Härdle, Marron, and Wand (1990), and the selector of Wu (1997) and Wu and Lin (2000). The performance of the multivariate versions of these selectors is yet to be established and we do not pursue this further here.

We now consider general normal mixture densities $f(\mathbf{x}) = \sum_{\ell=1}^k w_{\ell} \phi_{\Sigma_{\ell}}(\mathbf{x} - \boldsymbol{\mu}_{\ell})$ where $w_{\ell} > 0$ and $\sum_{\ell=1}^k w_{\ell} = 1$. Normal mixture densities are widely employed in simulation studies since they provide a rich class of densities with tractable exact squared error expressions. Normal mixture densities were used in early attempts for data-based bandwidth selection for multivariate kernel density estimation, see Ćwik and Koronacki (1997). They proposed a 2-step procedure: (i) a preliminary normal mixture is fitted to the data and (ii) the MISE- and AMISE-optimal bandwidths are computed from the closed form expressions for the MISE and AMISE for this normal mixture fit. We provide MISE and AMISE expressions for density derivatives to provide the basis for an analogous selector.

Theorem 7. *Suppose (A1) holds, that f is the normal mixture density $\sum_{\ell=1}^k w_{\ell} \phi_{\Sigma_{\ell}}(\cdot - \boldsymbol{\mu}_{\ell})$, and that K is the normal kernel. Then*

$$\begin{aligned} & \text{MISE}\{\widehat{\text{D}}^{\otimes r} f(\cdot; \mathbf{H})\} \\ &= 2^{-r} n^{-1} (4\pi)^{-d/2} |\mathbf{H}|^{-1/2} \boldsymbol{\mu}_r(\mathbf{H}) + \mathbf{w}^T \{(1 - n^{-1}) \boldsymbol{\Omega}_2 - 2\boldsymbol{\Omega}_1 + \boldsymbol{\Omega}_0\} \mathbf{w}, \end{aligned}$$

where $\mathbf{w} = (w_1, w_2, \dots, w_k)^T$ and $\boldsymbol{\Omega}_a \in \mathcal{M}_{k \times k}$ with (ℓ, ℓ') entry given by $(\boldsymbol{\Omega}_a)_{\ell, \ell'} = (-1)^r (\text{vec}^T \mathbf{I}_{d^r}) \text{D}^{\otimes 2r} \phi_{\mathbf{a}\mathbf{H} + \Sigma_{\ell\ell'}}(\boldsymbol{\mu}_{\ell\ell'})$, with $\boldsymbol{\mu}_{\ell\ell'} = \boldsymbol{\mu}_{\ell} - \boldsymbol{\mu}_{\ell'}$, $\Sigma_{\ell\ell'} = \Sigma_{\ell} + \Sigma_{\ell'}$. Then

$$\begin{aligned} & \text{MISE}\{\widehat{\mathbf{D}}^{\otimes r} f(\cdot; \mathbf{H})\} \\ &= 2^{-r} \text{OF}(2r) n^{-1} (4\pi)^{-d/2} |\mathbf{H}|^{-1/2} (\text{vec}^T \mathbf{H}^{-1})^{\otimes r} \mathcal{S}_{d,2r}(\text{vec} \mathbf{I}_d)^{\otimes r} \\ & \quad + \mathbf{w}^T \{(1 - n^{-1}) \boldsymbol{\Omega}_2 - 2\boldsymbol{\Omega}_1 + \boldsymbol{\Omega}_0\} \mathbf{w}, \end{aligned}$$

where

$$\begin{aligned} (\boldsymbol{\Omega}_a)_{\ell, \ell'} &= (-1)^r \phi_{a\mathbf{H} + \boldsymbol{\Sigma}_{\ell\ell'}}(\boldsymbol{\mu}_{\ell\ell'}) (\text{vec}^T (a\mathbf{H} + \boldsymbol{\Sigma}_{\ell\ell'})^{-2})^{\otimes r} \mathcal{S}_{d,2r} \\ & \quad \times \sum_{j=0}^r (-1)^j \text{OF}(2j) \binom{2r}{2j} [\boldsymbol{\mu}_{\ell\ell'}^{\otimes(2r-2j)} \otimes (\text{vec}(a\mathbf{H} + \boldsymbol{\Sigma}_{\ell\ell'}))^{\otimes j}]. \end{aligned}$$

Theorem 7 is the analogue of Theorem 1 in Wand and Jones (1993). The following AMISE formula is the analogue of Theorem 1 in Wand (1992).

Theorem 8. *Under the conditions of Theorem 7,*

$$\text{AMISE}\{\widehat{\mathbf{D}}^{\otimes r} f(\cdot; \mathbf{H})\} = 2^{-r} n^{-1} (4\pi)^{-d/2} |\mathbf{H}|^{-1/2} \mu_r(\mathbf{H}) + \frac{1}{4} \mathbf{w}^T \tilde{\boldsymbol{\Omega}} \mathbf{w},$$

where $\tilde{\boldsymbol{\Omega}} \in \mathcal{M}_{k \times k}$ with (ℓ, ℓ') entry given by $\tilde{\boldsymbol{\Omega}}_{\ell, \ell'} = (-1)^r \text{vec}^T(\mathbf{I}_{dr} \otimes (\text{vec} \mathbf{H} \text{vec}^T \mathbf{H})) \mathbf{D}^{\otimes 2r+4} \phi_{\boldsymbol{\Sigma}_{\ell\ell'}}(\boldsymbol{\mu}_{\ell\ell'})$, with $\boldsymbol{\mu}_{\ell\ell'} = \boldsymbol{\mu}_\ell - \boldsymbol{\mu}_{\ell'}$, $\boldsymbol{\Sigma}_{\ell\ell'} = \boldsymbol{\Sigma}_\ell + \boldsymbol{\Sigma}_{\ell'}$. Then

$$\begin{aligned} & \text{AMISE}\{\widehat{\mathbf{D}}^{\otimes r} f(\cdot; \mathbf{H})\} \\ &= 2^{-r} \text{OF}(2r) n^{-1} (4\pi)^{-d/2} |\mathbf{H}|^{-1/2} (\text{vec}^T \mathbf{H}^{-1})^{\otimes r} \mathcal{S}_{d,2r}(\text{vec} \mathbf{I}_d)^{\otimes r} \\ & \quad + \frac{1}{4} \mathbf{w}^T \tilde{\boldsymbol{\Omega}} \mathbf{w}, \end{aligned}$$

where

$$\begin{aligned} \tilde{\boldsymbol{\Omega}}_{\ell, \ell'} &= (-1)^r \phi_{\boldsymbol{\Sigma}_{\ell\ell'}}(\boldsymbol{\mu}_{\ell\ell'}) [(\text{vec}^T \boldsymbol{\Sigma}_{\ell\ell'}^{-2})^{\otimes r} \otimes (\text{vec}^T (\boldsymbol{\Sigma}_{\ell\ell'}^{-1} \mathbf{H} \boldsymbol{\Sigma}_{\ell\ell'}^{-1}))^{\otimes 2}] \\ & \quad \times \mathcal{S}_{d,2r+4} \sum_{j=0}^{r+2} (-1)^j \text{OF}(2j) \binom{2r+4}{2j} [\boldsymbol{\mu}_{\ell\ell'}^{\otimes(2r-2j+4)} \otimes (\text{vec} \boldsymbol{\Sigma}_{\ell\ell'})^{\otimes j}]. \end{aligned}$$

For $r = 0$, Theorem 8 gives $\text{AMISE}\{\hat{f}(\cdot; \mathbf{H})\} = n^{-1} (4\pi)^{-d/2} |\mathbf{H}|^{-1/2} + \mathbf{w}^T \tilde{\boldsymbol{\Omega}} \mathbf{w} / 4$ with $\tilde{\boldsymbol{\Omega}}_{\ell, \ell'} = \phi_{\boldsymbol{\Sigma}_{\ell\ell'}}(\boldsymbol{\mu}_{\ell\ell'}) (\text{vec}^T (\boldsymbol{\Sigma}_{\ell\ell'}^{-1} \mathbf{H} \boldsymbol{\Sigma}_{\ell\ell'}^{-1}))^{\otimes 2} \mathcal{S}_{d,4} [\boldsymbol{\mu}_{\ell\ell'}^{\otimes 4} - 6\boldsymbol{\mu}_{\ell\ell'}^{\otimes 2} \otimes \text{vec} \boldsymbol{\Sigma}_{\ell\ell'} + 3(\text{vec} \boldsymbol{\Sigma}_{\ell\ell'})^{\otimes 2}]$. This appears to be completely different, though equivalent to, Theorem 1 in Wand (1992).

4. Asymptotic Relative Efficiency

We examine the gain in density estimation performance when using the added flexibility of unconstrained selectors. The usual measure of asymptotic performance is the minimal achievable AMISE. We compare this minimal AMISE

for these parametrization classes: \mathcal{F} is the class of all positive-definite matrices, \mathcal{D} is the class of all positive-definite diagonal matrices, and \mathcal{I} is the class of positive scalar multiples of the identity matrix. We consider f to be a single normal density to simplify the mathematical analysis.

Corollary 3. *Suppose the conditions of Theorem 5 hold. For the class of unconstrained bandwidth matrices \mathcal{F} ,*

$$\min_{\mathbf{H} \in \mathcal{F}} \text{AMISE}(\mathbf{H}) = 2^{-(d+r+4)} 2^{8/(d+2r+4)} \pi^{-d/2} (d+2r+4)(d+2r+2)^{(d+2r)/(d+2r+4)} \\ \times |\boldsymbol{\Sigma}|^{-1/2} \mu_r(\boldsymbol{\Sigma}) n^{-4/(d+2r+4)}.$$

The asymptotic rate of convergence of the minimal achievable AMISE was previously stated in Theorem 3 for general f . This corollary provides its constants when f is a single normal density, generalizing the result from Wand and Jones (1995, p. 112) to general r .

Corollary 4. *Suppose the conditions of Theorem 5 hold. For the class \mathcal{I} , $\mathbf{H} = h^2 \mathbf{I}_d$,*

$$\text{AMISE}(\mathbf{H}) = 2^{-(d+r)} \pi^{-d/2} \{n^{-1} h^{-d-2r} \mu_r(\mathbf{I}_d) + \frac{1}{16} |\boldsymbol{\Sigma}|^{-1/2} \mu_{r+2}(\boldsymbol{\Sigma}) h^4\}.$$

The bandwidth that minimizes the AMISE is $\mathbf{H}_{\text{AMISE}} = h_{\text{AMISE}}^2 \mathbf{I}_d$, where

$$h_{\text{AMISE}} = \left(\frac{4(d+2r) |\boldsymbol{\Sigma}|^{1/2} \mu_r(\mathbf{I}_d)}{\mu_{r+2}(\boldsymbol{\Sigma}) n} \right)^{1/(d+2r+4)}.$$

The minimal achievable AMISE is

$$\min_{\mathbf{H} \in \mathcal{I}} \text{AMISE}(\mathbf{H}) \\ = 2^{-(d+r+4)} 2^{8/(d+2r+4)} \pi^{-d/2} \{ |\boldsymbol{\Sigma}|^{-(d+2r)/2} \mu_{r+2}(\boldsymbol{\Sigma})^{d+2r} \mu_{r+1}(\mathbf{I}_d)^4 \}^{1/(d+2r+4)} \\ \times (d+2r+4)(d+2r)^{-1} n^{-4/(d+2r+4)}.$$

Comparing this corollary to the previous one, the rate of AMISE convergence does not depend on the parametrization class. The difference in finite sample performance is due to the different coefficients of the minimal AMISE. The gain in density estimation efficiency using an unconstrained bandwidth matrix over constrained bandwidths was established in the bivariate case by Wand and Jones (1993). Their main measure of this gain is the Asymptotic Relative Efficiency

$$\text{ARE}(\mathcal{F} : \mathcal{D}) = \left[\frac{\min_{\mathbf{H} \in \mathcal{F}} \text{AMISE}(\mathbf{H})}{\min_{\mathbf{H} \in \mathcal{D}} \text{AMISE}(\mathbf{H})} \right]^{(d+2r+4)/4}.$$

The interpretation of $\text{ARE}(\mathcal{F} : \mathcal{D})$ is that, for large n , the minimum error using n observations with a diagonal bandwidth can be achieved using only $\text{ARE}(\mathcal{F} : \mathcal{D}) \times n$ observations with an unconstrained \mathbf{H} . Analogous definitions and interpretations apply to $\text{ARE}(\mathcal{F} : \mathcal{I})$ and $\text{ARE}(\mathcal{D} : \mathcal{I})$.

Computing these AREs analytically for general densities is mathematically intractable, so we focus on the case where f is a normal density, making use of Corollaries 3 and 4.

Corollary 5. *Suppose the conditions of Theorem 5 hold. Then*

$$\begin{aligned} \text{ARE}(\mathcal{F} : \mathcal{I}) &= [(d+2r+2)(d+2r)]^{(d+2r)/4} |\boldsymbol{\Sigma}|^{-1/2} \mu_r(\boldsymbol{\Sigma})^{(d+2r+4)/4} \mu_{r+2}(\boldsymbol{\Sigma})^{-(d+2r)/4} \mu_r(\mathbf{I}_d)^{-1}. \end{aligned}$$

Corollary 6. *Suppose the conditions of Theorem 5 hold, and that f is a bivariate normal density function having variances equal and with correlation coefficient ρ . Then*

$$\text{ARE}(\mathcal{F} : \mathcal{D}) = \text{ARE}(\mathcal{F} : \mathcal{I}) = [4(r+2)(r+1)]^{(r+1)/2} \frac{(1-\rho^2)^{1/2} Q(r, \rho)^{(r+3)/2}}{Q(r, 0) Q(r+2, \rho)^{(r+1)/2}},$$

where

$$Q(r, \rho) = \sum_{j=0}^r \sum_{j'=0}^j \binom{r}{j} \binom{j}{j'} (-2\rho)^{j-j'} m_{j+j'} m_{2r-j-j'}$$

and $m_k = (1/2)\{(-1)^k + 1\}\text{OF}(k)$ for $k = 0, 1, 2, \dots$

In Figure 1, we compare the AREs for four families of bivariate normal densities with mean zero, marginal variances σ_1^2, σ_2^2 , and correlation coefficient ρ ranging from -1 to $+1$. Without loss of generality, $\sigma_1 = 1$, and we let $\sigma_2 = 1, 2, 5, 10$. For each value of ρ , we compute $\text{ARE}(\mathcal{F} : \mathcal{D})$ numerically and $\text{ARE}(\mathcal{F} : \mathcal{I})$ analytically. The former are plotted as the black curves, and the latter as grey curves. We consider derivatives of order $r = 0, 1, \dots, 4$ which are drawn in the solid, short dashed, dotted, dot-dashed and long dashed lines, respectively. This figure generalizes the plots in Wand and Jones (1993). Note that, apart from equal marginal variances $\sigma_1 = \sigma_2$ with weak correlation, $\text{ARE}(\mathcal{F} : \mathcal{I})$ is close to zero, indicating that the class \mathcal{I} is inadequate for moderate to high correlation. Another observed feature is that the rate that both AREs tend to zero, as $|\rho|$ tends to 1, increases as the derivative order increases. This indicates that the gain from unconstrained bandwidths for higher derivatives exceeds the known gain for $r = 0$ (Wand and Jones (1993)).

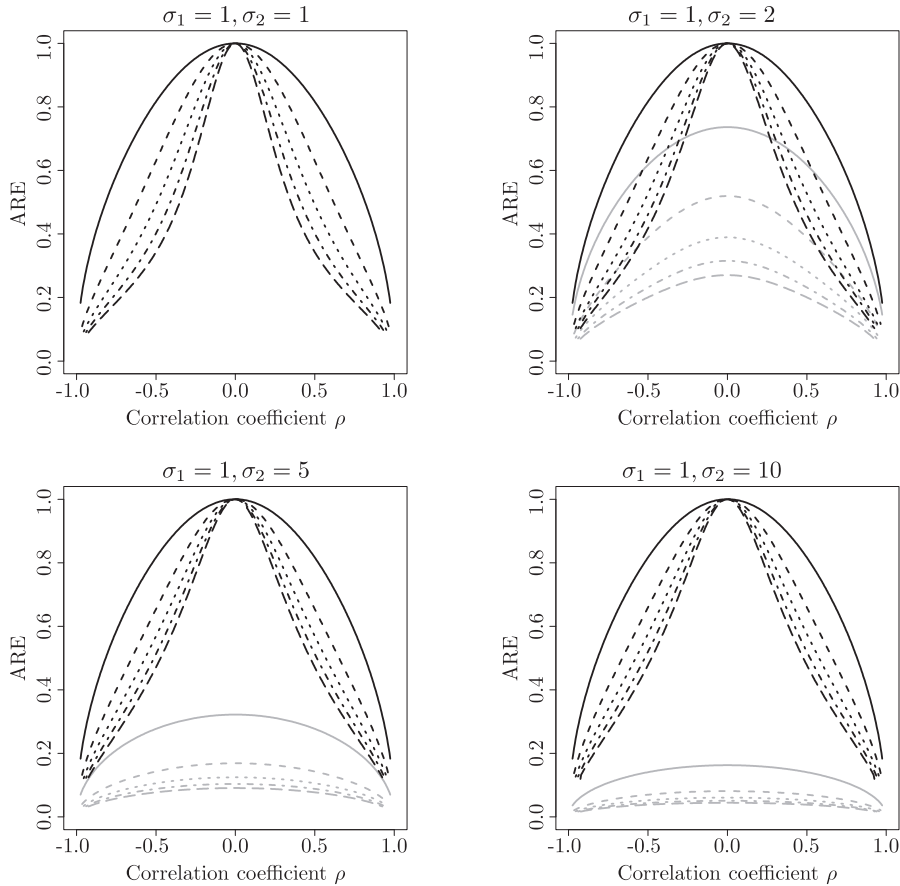


Figure 1. Each family of target densities is a normal density with marginal variances σ_1 and σ_2 , and correlation coefficient ρ . The black curves are $\text{ARE}(\mathcal{F} : \mathcal{D})$, the grey curves $\text{ARE}(\mathcal{F} : \mathcal{I})$. For $\sigma_1 = \sigma_2$ the two AREs coincide. The derivatives are: $r = 0$ (solid), $r = 1$ (short dashed), $r = 2$ (dotted), $r = 3$ (dot-dashed), $r = 4$ (long dashed).

5. Numerical Study

We conducted a numerical simulation study to compare the finite sample performance for unconstrained and constrained bandwidths. We used the normal scale selector $\hat{\mathbf{H}}_{\text{NS}}$, labelled NSF and NSD, where the ‘F’ and ‘D’ denote the \mathcal{F} and \mathcal{D} parametrization classes defined at the beginning of Section 4, respectively. In the `ks` library in R, the functions `Hmise.mixt()` and `Hamise.mixt()` compute unconstrained minimizers of the MISE and AMISE (from Theorems 7 and 8.). Their constrained versions are `Hmise.mixt.diag()` and `Hamise.mixt.diag()`. To compute the normal scale selectors, we took a zero mean and the sample variance as the inputs to these R functions. For each replicate, we com-

puted the Integrated Squared Error (ISE) between the resulting kernel density estimate and the target density. The ISE of the normal mixture density $f(\mathbf{x}) = \sum_{j=1}^k w_j \phi_{\Sigma_j}(\mathbf{x} - \boldsymbol{\mu}_j)$ has the explicit form

$$\begin{aligned} \text{ISE}(\mathbf{H}) &= \int_{\mathbb{R}^d} [\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H}) - D^{\otimes r} f(\mathbf{x})]^T [\widehat{D^{\otimes r} f}(\mathbf{x}; \mathbf{H}) - D^{\otimes r} f(\mathbf{x})] d\mathbf{x} \\ &= (-1)^r (\text{vec}^T \mathbf{I}_{d^r}) \left[n^{-2} \sum_{i=1}^n \sum_{j=1}^n D^{\otimes 2r} \phi_{2\mathbf{H}}(\mathbf{X}_i - \mathbf{X}_j) - 2n^{-1} \sum_{i=1}^n \sum_{j=1}^k w_j D_{\mathbf{H} + \Sigma_j}^{\otimes 2r}(\mathbf{X}_i - \boldsymbol{\mu}_j) \right. \\ &\quad \left. + \sum_{i=1}^k \sum_{j=1}^k w_i w_j D_{\Sigma_i + \Sigma_j}^{\otimes 2r}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) \right]. \end{aligned} \tag{5.1}$$

5.1. Bivariate target densities

We took sample sizes $n = 100$ and $n = 1,000$ for derivative orders $r = 0, 1, 2$, and each selector was computed on 100 replicates. The target densities we investigated are taken from Chacón (2009), with their original names and numbering, cover a wide range of density shapes, as shown in Figure 2. Target density #1 is a normal density and can be considered a base case. Density #2 is a correlated normal, densities #7 and #11 are multimodal, with varying degrees of intricate structure.

The box plots of the ISE of the 100 replicates for each target density, sample size n , and derivative order r , are displayed in Figure 3. As expected, unconstrained selectors make minimal gains for density #1 since most of its probability mass is oriented parallel to the co-ordinate axes. For density #7, NSF is detrimental compared to NSD. For $n = 1,000$, typical bandwidths were

$$\mathbf{H}_{\text{AMISE}} = \begin{bmatrix} 0.056 & 0.034 \\ 0.034 & 0.056 \end{bmatrix}, \quad \hat{\mathbf{H}}_{\text{NSD}} = \begin{bmatrix} 0.108 & 0 \\ 0 & 0.108 \end{bmatrix}, \quad \hat{\mathbf{H}}_{\text{NSF}} = \begin{bmatrix} 0.144 & -0.073 \\ -0.073 & 0.144 \end{bmatrix}.$$

According the AMISE selector, the kernel should be positively correlated, with $\rho = 0.60$. This is natural in view of the shape of density #7. The NSD selector gives uncorrelated kernels by construction, whereas the fact that the overall correlation coefficient of density #7 is negative makes the NSF select a kernel with negative correlation ($\rho = -0.51$), thus resulting in higher ISEs than with the diagonal selector. This shows that the normal scale is a poor surrogate in this case. For densities #2 and #11, NSF shows the most improvement. We conclude that these mixed results are not due to the inherent inadequacy of unconstrained as compared to diagonal selectors, but rather due to confounding with the crudeness of normal scale selectors. This serves as motivation to develop unconstrained versions of ‘hi-tech’ selectors in the future.

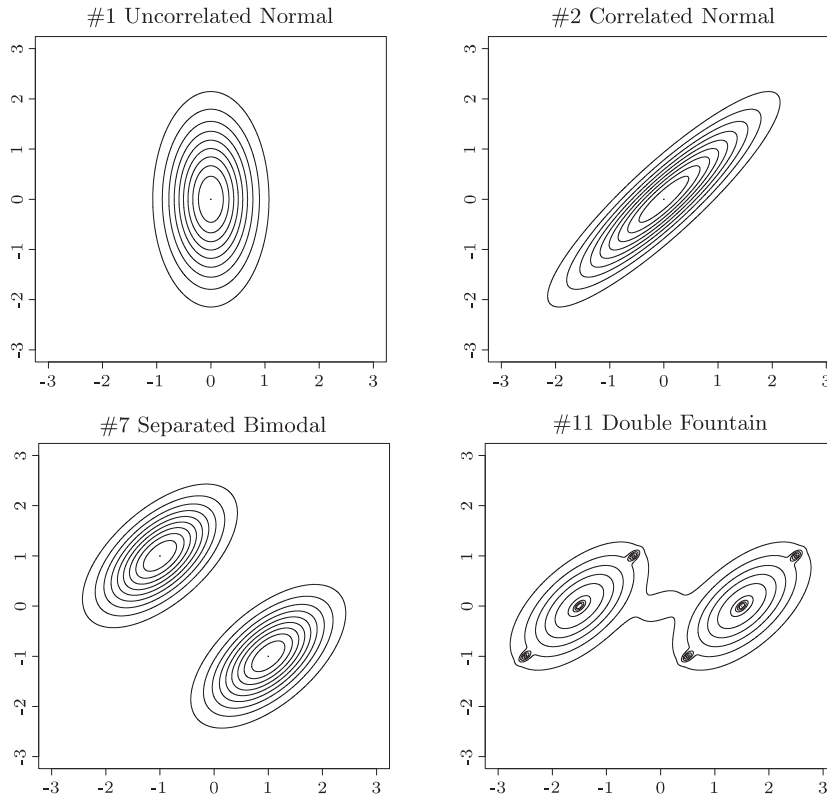


Figure 2. Contour plots for the four target densities.

5.2. Multivariate target densities

For a multivariate study, we generalized density #2: a correlated normal density for dimensions $d = 2, 3, 4$, with zero mean, variances one, and correlations 0.9. We took sample size $n = 100$, derivative orders $r = 0, 1, 2, 3$, with 500 replicates. From Figure 4, for this correlated normal density, there is a uniform decrease in the ISE for the unconstrained parametrization. For a fixed dimension d , the gain in performance decreases as the derivative order r increases. On the other hand, for a fixed order r , the gain in performance increases as d increases.

6. Application: High-Throughput Flow Cytometry

Flow cytometry is a method by which multiple characteristics of single cells or other particles are simultaneously measured as they pass through a laser beam in a fluid stream (Shapiro (2003)). The last few years have seen a major change in technology toward what has become known as *high-throughput* (or *high-content*) flow cytometry (e.g., Le Meur et al. (2007)). This technology combines robotic

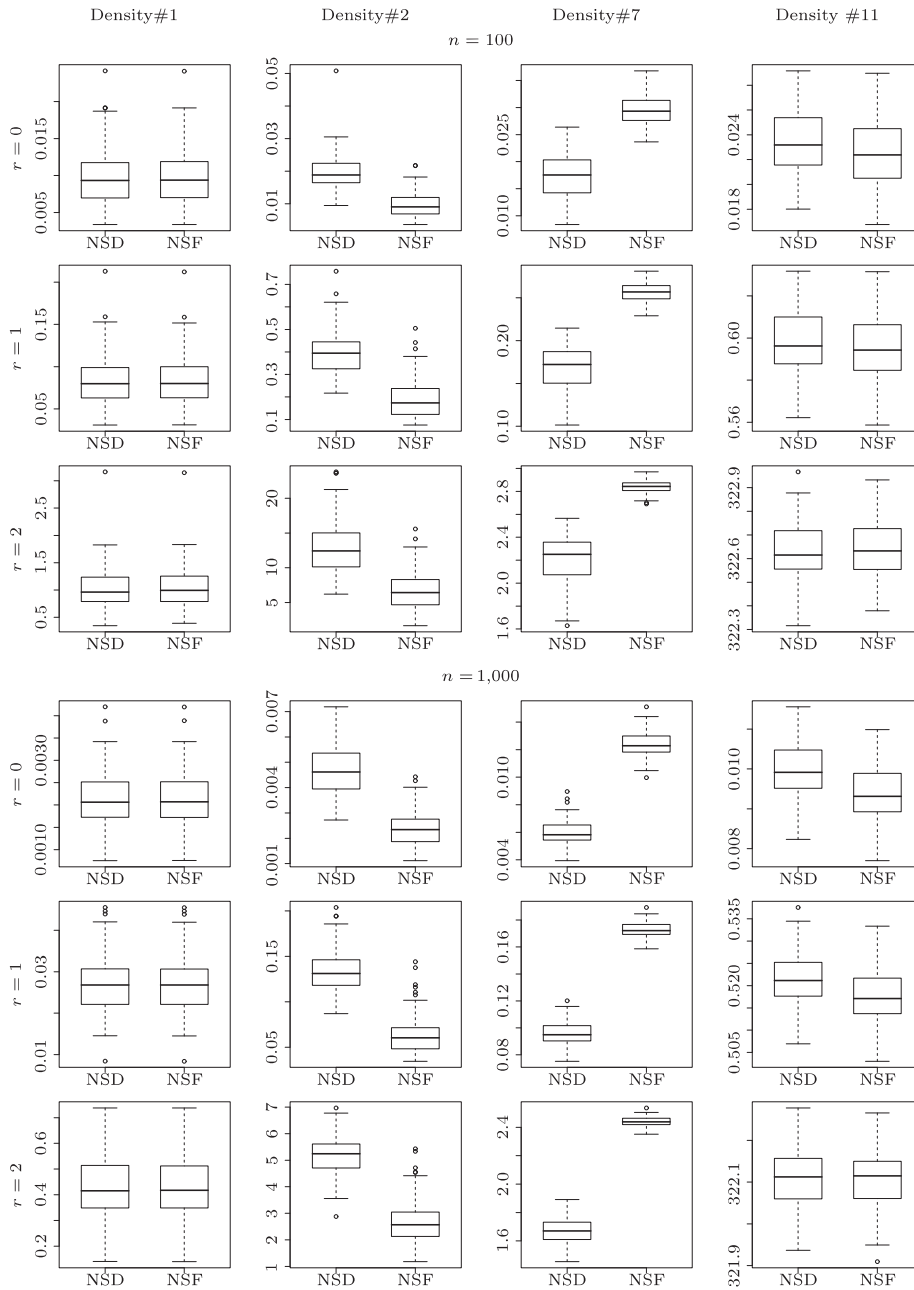


Figure 3. ISE box plots to compare the performance of diagonal and unconstrained parametrizations of normal scale selectors for bivariate densities. There are four target densities, sample sizes $n = 100, 1,000$ and derivative orders $r = 0, 1, 2$.

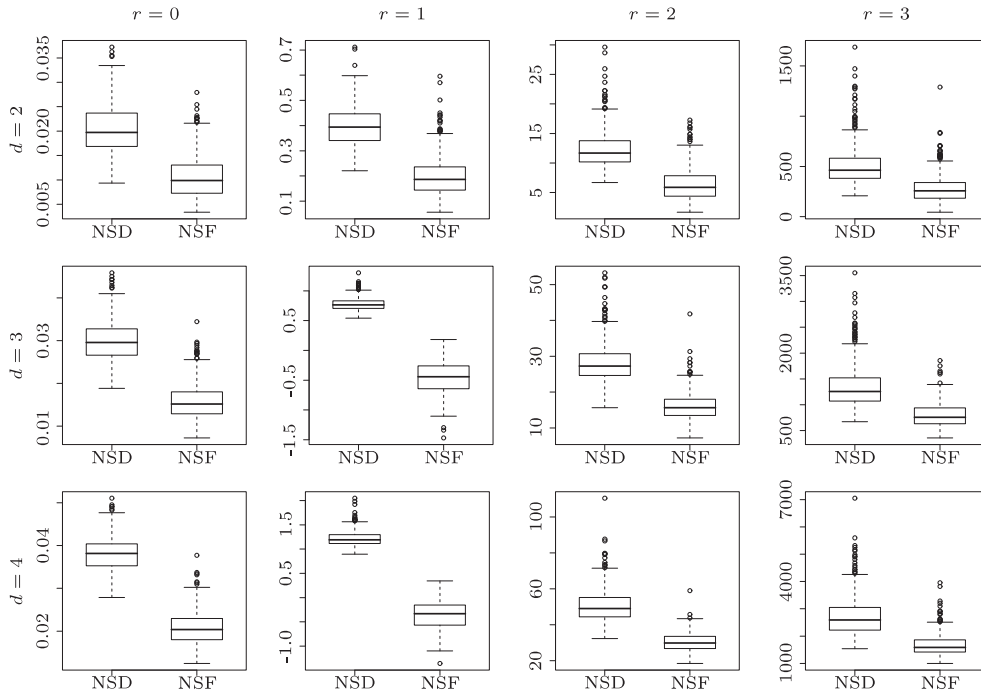


Figure 4. ISE box plots to compare the performance of diagonal and unconstrained parametrizations of normal scale selectors for multivariate correlated normal densities, $d = 2, 3, 4$. The sample size is $n = 100$, and the derivative orders are $r = 0, 1, 2, 3$.

fluid handling, flow cytometric instrumentation, and bioinformatics software so that relatively large numbers of cells can be processed and analyzed in a short period of time. With such massive amounts of data, there is a high premium on good automatic methods for pre-processing and extraction of clinically relevant information.

Figure 5 is a subset of data from the flow cytometry experiment described in Brinkman et al. (2007). The left panel is cellular fluorescence measurements – corresponding to antibodies CD4 and CD8 β , after gating on CD3-positive cells – on a patient who develops graft-versus-host disease. The right panel corresponds to a control. The data were collected 32 days after each patient had a blood and marrow transplant. The goal is to identify cell populations that differ between control and disease groups and, hence, constitute valid disease biomarkers, e.g., CD4-positive, CD8 β -positive, CD3-positive; where ‘positive’ indicates fluorescence of the relevant antibody above a threshold. The shapes in Figure 5 correspond to regions of high significant negative curvature based on the methodology of Godtlielsen, Marron, and Chaudhuri (2002), and refined by

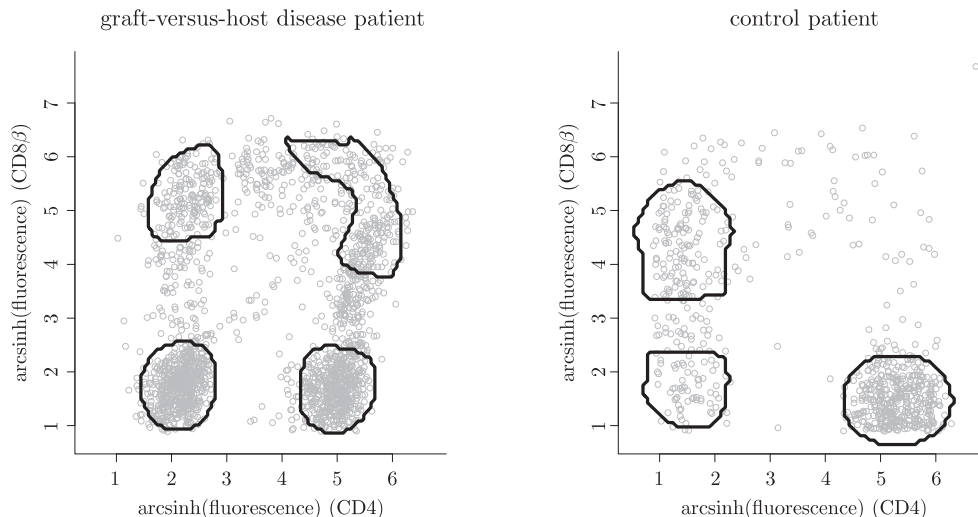


Figure 5. Cellular fluorescence measurements, after undergoing the arcsinh transformation, corresponding to antibodies CD4 and CD8 β after subsetting on CD3-positive cells. The left panel is data from a patient who develops graft-versus-host disease. The right panel is data from a control patient. Further details about the data are given in Brinkman et al. (2007). The shapes correspond to significant negative density curvature regions using the methodology of Duong et al. (2008) with the bandwidth chosen via the normal scale rule (3.2).

Duong et al. (2008). The bandwidth matrix is chosen according to the normal scale rule (3.2) with $d = r = 2$:

$$\hat{\mathbf{H}}_{\text{NS}} = (1/2)^{1/5} \hat{\Sigma} n^{-1/5}, \quad (6.1)$$

where $\hat{\Sigma}$ is the sample variance. Since the normal density is close to being the density which gives the largest optimal amount of smoothing given a fixed variance for density estimation (Terrell (1990)), (6.1) corresponds approximately to the largest bandwidth matrix which should be considered for curvature estimation. The absence of significant curvature for CD4-positive and CD8 β -positive cells in the control patient, despite use of this maximal bandwidth, represents an important clinical difference and gives rise to useful cellular signatures for graft-versus-host disease. Using the $r = 0$ normal scale rule, as illustrated in Table 1, could lead to insufficient smoothing for large sample sizes. In more comprehensive analyses of these data, described in Naumann and Wand (2009), more sophisticated filters for identifying cellular signatures are employed. The normal scale rule for second derivative estimation plays an important role in the initial phases of these filters, identifying candidate modal regions of possible interest.

The plots in Figure 5 were computed using the R library `feature` whose main function uses (6.1) as the upper limit on the default bandwidth matrix range.

7. Discussion

Kernel smoothing is a widely used non-parametric method for multivariate density estimation. It has the potential to be as equally successful for density derivative estimation. The relative lack of theoretical development for density derivatives compared to densities has hindered this progress. One obstacle is the specification of higher order derivatives. By writing the r th order array of r th order differentials as an r -fold Kronecker product of first order differentials, we maintain an intuitive, systematic vectorization of all derivatives. This allows the derivation of such quantities as the MISE and AMISE for kernel density estimators for general derivatives.

The single most important factor in the performance of kernel estimators is the choice of the bandwidth. For density estimation, there is now a solid body of work for reliable bandwidth matrix selection. Using the theoretical simplifications afforded by our vector form derivatives, we can write down an unconstrained data-driven selector based on normal scales. These normal scale selectors facilitate the quantification of the possible gain in performance in using the unconstrained bandwidth matrices compared to more constrained parametrizations. These selectors are a starting point from which more advanced unconstrained bandwidth selectors can be now developed, and for the second derivative, they are a starting point from which to estimate modal regions.

Acknowledgements

J. E. Chacón has been partially supported by Spanish Ministerio de Ciencia y Tecnología project MTM2006-06172, T. Duong was funded by Institut Pasteur through a ‘Programme Transversal de Recherches’ grant (PTR No. 218) and by Institut Curie through a Mayent-Rothschild Fellowship, and M.P. Wand received support from Australian Research Council Grant DP055651. The authors are grateful for the assistance received from Nolwenn Le Meur and Richard White.

Appendix: Proofs

A.1. Proof of the results in Section 2

A.1.1. Proof of Theorem 1

Proof of Theorem 1. First notice that we can write $\widehat{\mathbb{E}D^{\otimes r}f(\mathbf{x}; \mathbf{H})} = D^{\otimes r}K_{\mathbf{H}} * f(\mathbf{x}) = K_{\mathbf{H}} * D^{\otimes r}f(\mathbf{x})$. Therefore,

$$\int B^2(\mathbf{x}; \mathbf{H})d\mathbf{x} = \int \|K_{\mathbf{H}} * D^{\otimes r}f(\mathbf{x}) - D^{\otimes r}f(\mathbf{x})\|^2d\mathbf{x}$$

$$= \text{tr } \mathbf{R}(\mathbf{D}^{\otimes r} f) + \text{tr } \mathbf{R}_{K^*K, \mathbf{H}, r}(f) - 2 \text{tr } \mathbf{R}_{K, \mathbf{H}, r}(f),$$

as it is not difficult to check that $\int K_{\mathbf{H}} * \mathbf{D}^{\otimes r} f(\mathbf{x}) K_{\mathbf{H}} * \mathbf{D}^{\otimes r} f(\mathbf{x})^T d\mathbf{x} = \mathbf{R}_{K^*K, \mathbf{H}, r}(f)$. About the variance term, it is clear that

$$\int V(\mathbf{x}; \mathbf{H}) d\mathbf{x} = n^{-1} \int \mathbb{E} \|\mathbf{D}^{\otimes r} K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_1)\|^2 d\mathbf{x} - n^{-1} \int \|\mathbb{E} \mathbf{D}^{\otimes r} K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_1)\|^2 d\mathbf{x}. \tag{A.1}$$

The second integral in the right side is easily recognized as $\mathbf{R}_{K^*K, \mathbf{H}, r}(f)$ also, and for the first one we have

$$\begin{aligned} & \int \mathbb{E} \|\mathbf{D}^{\otimes r} K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_1)\|^2 d\mathbf{x} \\ &= \text{tr} \iint \mathbf{D}^{\otimes r} K_{\mathbf{H}}(\mathbf{x} - \mathbf{y}) \mathbf{D}^{\otimes r} K_{\mathbf{H}}(\mathbf{x} - \mathbf{y})^T f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= \text{tr} \int \mathbf{D}^{\otimes r} K_{\mathbf{H}}(\mathbf{x}) \mathbf{D}^{\otimes r} K_{\mathbf{H}}(\mathbf{x})^T d\mathbf{x} \\ &= \text{tr} \left[(\mathbf{H}^{-1/2})^{\otimes r} \int (\mathbf{D}^{\otimes r} K)_{\mathbf{H}}(\mathbf{x}) (\mathbf{D}^{\otimes r} K)_{\mathbf{H}}(\mathbf{x})^T d\mathbf{x} (\mathbf{H}^{-1/2})^{\otimes r} \right] \\ &= \text{tr} \left[(\mathbf{H}^{-1})^{\otimes r} |\mathbf{H}|^{-1/2} \int \mathbf{D}^{\otimes r} K(\mathbf{z}) \mathbf{D}^{\otimes r} K(\mathbf{z})^T d\mathbf{z} \right] \\ &= |\mathbf{H}|^{-1/2} \text{tr} \left((\mathbf{H}^{-1})^{\otimes r} \mathbf{R}(\mathbf{D}^{\otimes r} K) \right). \end{aligned}$$

A.1.2. Proof of Theorem 2

Notice that, for a function $f: \mathbb{R}^d \rightarrow \mathbb{R}^p$ such that every element in $\mathbf{D}^{\otimes q} f(\mathbf{x})$ is piecewise continuous, we can write its Taylor polynomial expansion as

$$f(\mathbf{x} + \mathbf{h}) = \sum_{r=0}^q \frac{1}{r!} [\mathbf{I}_p \otimes (\mathbf{h}^T)^{\otimes r}] \mathbf{D}^{\otimes r} f(\mathbf{x}) + o(\|\mathbf{h}\|^q) \mathbf{1}_p, \quad \mathbf{x}, \mathbf{h} \in \mathbb{R}^d.$$

See Baxandall and Liebeck (1986, p. 164). The proof of Theorem 2 then follows from Lemmas 1 and 2 below, together with the bias-variance decomposition of the MSE.

Write $\text{IB}^2(\mathbf{H}) = \int \text{B}^2(\mathbf{x}; \mathbf{H}) d\mathbf{x}$ and $\text{IV}(\mathbf{H}) = \int V(\mathbf{x}; \mathbf{H}) d\mathbf{x}$ as the integrated squared bias and integrated variance of the kernel estimator, respectively, so that $\text{MISE}(\mathbf{H}) = \text{IB}^2(\mathbf{H}) + \text{IV}(\mathbf{H})$.

Lemma 1. *Suppose (A1)–(A3) hold. Then*

$$\text{IB}^2(\mathbf{H}) = \frac{m_2(K)^2}{4} \text{tr} \left[(\mathbf{I}_{dr} \otimes \text{vec}^T \mathbf{H}) \mathbf{R}(\mathbf{D}^{\otimes(r+2)} f) (\mathbf{I}_{dr} \otimes \text{vec } \mathbf{H}) \right] + o(\text{tr}^2 \mathbf{H}).$$

Proof. We can write $\mathbb{E}\widehat{D^{\otimes r}f(\mathbf{x}; \mathbf{H})} = \int K(\mathbf{z})D^{\otimes r}f(\mathbf{x} - \mathbf{H}^{1/2}\mathbf{z})d\mathbf{z}$. Now, make use of a Taylor expansion to get

$$D^{\otimes r}f(\mathbf{x} - \mathbf{H}^{1/2}\mathbf{z}) = D^{\otimes r}f(\mathbf{x}) - [\mathbf{I}_{dr} \otimes (\mathbf{z}^T \mathbf{H}^{1/2})]D^{\otimes(r+1)}f(\mathbf{x}) + \frac{1}{2}[\mathbf{I}_{dr} \otimes (\mathbf{z}^T \mathbf{H}^{1/2})^{\otimes 2}]D^{\otimes(r+2)}f(\mathbf{x}) + o(\text{tr } \mathbf{H})\mathbf{1}_{dr}.$$

Substitute this in the previous formula and use (A3) to obtain

$$\begin{aligned} B(\mathbf{x}; \mathbf{H}) &= \frac{m_2(K)}{2} \|\mathbf{I}_{dr} \otimes \{(\text{vec}^T \mathbf{I}_d)(\mathbf{H}^{1/2})^{\otimes 2}\}\| D^{\otimes(r+2)}f(\mathbf{x}) + o(\text{tr } \mathbf{H})\| \\ &= \frac{m_2(K)}{2} \|\mathbf{I}_{dr} \otimes \text{vec}^T \mathbf{H}\| D^{\otimes(r+2)}f(\mathbf{x}) + o(\text{tr } \mathbf{H})\|. \end{aligned}$$

We finish the proof by squaring and integrating the previous expression, taking into account (A2).

Lemma 2. *Suppose (A1) holds. Then*

$$IV(\mathbf{H}) = n^{-1}|\mathbf{H}|^{-1/2} \text{tr} ((\mathbf{H}^{-1})^{\otimes r} \mathbf{R}(D^{\otimes r} K)) + o(n^{-1}|\mathbf{H}|^{-1/2} \text{tr}^r(\mathbf{H}^{-1})).$$

Proof. From the proof of Theorem 1 and the arguments in the previous lemma, we have

$$\begin{aligned} \int V(\mathbf{x}; \mathbf{H})d\mathbf{x} &= n^{-1} \int \mathbb{E}\|D^{\otimes r} K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_1)\|^2 d\mathbf{x} + O(n^{-1}) \\ &= n^{-1}|\mathbf{H}|^{-1/2} \text{tr} ((\mathbf{H}^{-1})^{\otimes r} \mathbf{R}(D^{\otimes r} K)) + o(n^{-1}|\mathbf{H}|^{-1/2} \text{tr}^r(\mathbf{H}^{-1})). \end{aligned}$$

A.1.3. Proof of Theorem 3

Proof. Similar to the decomposition $MISE(\mathbf{H}) = IB^2(\mathbf{H}) + IV(\mathbf{H})$, we can write $AMISE(\mathbf{H}) = AIB^2(\mathbf{H}) + AIV(\mathbf{H})$, where $AIB^2(\mathbf{H})$ and $AIV(\mathbf{H})$ are the leading terms of $IB^2(\mathbf{H})$ and $IV(\mathbf{H})$, respectively, from Lemmas 1 and 2.

Let $\mathbf{K}_{r,s} \in \mathcal{M}_{rs \times rs}$ be the commutation matrix of order r, s ; see Magnus and Neudecker (1979). The commutation matrix allows us to commute the order of the matrices in a Kronecker product, e.g., if $\mathbf{A} \in \mathcal{M}_{n \times r}$ and $\mathbf{B} \in \mathcal{M}_{m \times s}$, then $\mathbf{K}_{m,n}(\mathbf{A} \otimes \mathbf{B})\mathbf{K}_{r,s} = \mathbf{B} \otimes \mathbf{A}$.

To determine the derivative, we first find the differentials. Differentials of a function $f : \mathbb{R}^d \rightarrow \mathbb{R}^p$ have the advantage that they are always the same dimension as f itself, as opposed to derivatives whose dimension depends on the order of the derivative. So higher order differentials are easier to manipulate. The first identification theorem of Magnus and Neudecker (1999, p. 87) states if the differential of $f(\mathbf{y})$ can be expressed as $df(\mathbf{y}) = A(\mathbf{y})\mathbf{y}$ for some matrix $A(\mathbf{y}) \in \mathcal{M}_{p \times d}$ then the derivative is $Df(\mathbf{y}) = A(\mathbf{y})$. The differential of $AIB^2(\mathbf{H})$ is

$$dAIB^2(\mathbf{H}) = \frac{m_2(K)^2}{4} (\text{vec}^T \mathbf{R}(D^{\otimes(r+2)} f)) \mathbf{K}_{dr,d^2}^{\otimes 2} [(\mathbf{K}_{d^{r+2},d^2} \otimes \mathbf{I}_{dr})$$

$$+ \mathbf{I}_{d^{2r+4}}] \times (\mathbf{I}_{d^2} \otimes \text{vec}\mathbf{H} \otimes \text{vec}\mathbf{I}_{d^r}) d\text{vec}\mathbf{H},$$

since

$$\begin{aligned} & d \text{vec}\{\mathbf{I}_{d^r} \otimes (\text{vec}\mathbf{H} \text{vec}^T \mathbf{H})\} \\ &= d \text{vec}[\mathbf{K}_{d^r, d^2} \{(\text{vec}\mathbf{H} \text{vec}^T \mathbf{H}) \otimes \mathbf{I}_{d^r}\} \mathbf{K}_{d^2, d^r}] \\ &= (\mathbf{K}_{d^r, d^2} \otimes \mathbf{K}_{d^r, d^2}) d \text{vec}\{(\text{vec}\mathbf{H} \text{vec}^T \mathbf{H}) \otimes \mathbf{I}_{d^r}\} \\ &= \mathbf{K}_{d^r, d^2}^{\otimes 2} [(\mathbf{K}_{d^{r+2}, d^2} \otimes \mathbf{I}_{d^r}) + \mathbf{I}_{d^{2r+4}}](\mathbf{I}_{d^2} \otimes \text{vec}\mathbf{H} \otimes \text{vec}\mathbf{I}_{d^r}) d\text{vec}\mathbf{H}, \end{aligned}$$

where the last line follows by using a similar reasoning to determine Equation (11) in the proof of Theorem 2 in Chacón and Duong (2010).

The differential of AIV(\mathbf{H}) is

$$\begin{aligned} d\text{AIV}(\mathbf{H}) = & - \left\{ \frac{1}{2} \text{AIV}(\mathbf{H})(\text{vec}^T \mathbf{H}^{-1}) + n^{-1} |\mathbf{H}|^{-1/2} (\text{vec}^T \mathbf{R}(\mathbf{D}^{\otimes r} K)) (\mathbf{H}^{-1})^{\otimes 2r} \right. \\ & \left. \times \mathbf{\Lambda}_r [(\mathbf{I}_{d^{r-1}} \otimes \mathbf{K}_{d, d^{r-1}})(\text{vec}\mathbf{H}^{\otimes(r-1)} \otimes \mathbf{I}_d) \otimes \mathbf{I}_d] \right\} d\text{vec}\mathbf{H}, \end{aligned}$$

where $\mathbf{\Lambda}_r = \sum_{i=1}^r \mathbf{K}_{d^i, d^{r-i}}^{\otimes 2}$. The reasoning follows similar lines to computing Equations (9) and (10) in the proof of Theorem 2 in Chacón and Duong (2010).

Let every entry of \mathbf{H} be $O(n^{-\beta})$ for $\beta > 0$. Then $d\text{AIV}(\mathbf{H}) = O(n^{\beta(d/2+r+1)-1}) d \text{vec}\mathbf{H}$ and $d\text{AIB}^2(\mathbf{H}) = O(n^{-\beta}) d \text{vec}\mathbf{H}$. Equating powers gives $\beta = 2/(d+2r+4)$ and thus $d\text{AMISE}(\mathbf{H}) = O(n^{-2/(d+2r+4)}) d\text{vec}\mathbf{H}$. The optimal \mathbf{H} is a solution of the equation $\partial\text{AMISE}(\mathbf{H})/(\partial\text{vec}\mathbf{H}) = 0$, so all its entries are $O(n^{-2/(d+2r+4)})$, which implies that $\min_{\mathbf{H}} \text{AMISE}(\mathbf{H}) = O(n^{-4/(d+2r+4)})$.

A.1.4. Proof of Theorem 4

The proof follows directly from Lemmas 1 and 2.

Proof. From Lemma 1, $\widehat{\mathbb{E}\mathbf{D}^{\otimes r} f(\mathbf{x}; \mathbf{H})} = \mathbf{D}^{\otimes r} f(\mathbf{x}) + (m_2(K)/2)(\mathbf{I}_{d^r} \otimes \text{vec}^T \mathbf{H}) \mathbf{D}^{\otimes(r+2)} f(\mathbf{x}) [1 + O(\text{tr}\mathbf{H})]$. For the variance, we have

$$\text{Var} \widehat{\mathbf{D}^{\otimes r} f(\mathbf{x}; \mathbf{H})} = n^{-1} \mathbf{D}^{\otimes r} K_{\mathbf{H}} \mathbf{D}^{\otimes r} K_{\mathbf{H}}^T * f(\mathbf{x}) - n^{-1} [K_{\mathbf{H}} * \mathbf{D}^{\otimes r} f(\mathbf{x})] [K_{\mathbf{H}} * \mathbf{D}^{\otimes r} f(\mathbf{x})]^T.$$

From Lemma 2, the convolution in the first term dominates the convolution in the second term, since the value of the former is

$$\mathbf{D}^{\otimes r} K_{\mathbf{H}} \mathbf{D}^{\otimes r} K_{\mathbf{H}}^T * f(\mathbf{x}) = |\mathbf{H}|^{-1/2} (\mathbf{H}^{-1/2})^{\otimes r} \mathbf{R}(\mathbf{D}^{\otimes r} K) (\mathbf{H}^{-1/2})^{\otimes r} f(\mathbf{x}) [1 + o(1)]$$

and the proof is complete.

A.2. Proof of the results in Section 3

The proofs in Sections A.2.1 and A.2.2 assume, without loss of generality, that $f = \phi_{\Sigma}$. Results for the general normal density $f = \phi_{\Sigma}(\cdot - \boldsymbol{\mu})$ remain valid since they are invariant under this translation.

A.2.1. Proof of Theorem 5

The proof of Theorem 5 is based on the exact formula given in Theorem 1. Notice that, in the normal case, we have $\mathbf{R}_{\phi*\phi,\mathbf{H},r}(\phi_\Sigma) = \mathbf{R}_{\phi,2\mathbf{H},r}(\phi_\Sigma)$ and $\mathbf{R}(\mathbf{D}^{\otimes r}\phi_\Sigma) = \mathbf{R}_{\phi,\mathbf{0},r}(\phi_\Sigma)$, so that it follows that all we need to have an explicit expression for the MISE function in the normal case is an explicit formula for $\text{tr}\mathbf{R}_{\phi,\mathbf{H},r}(\phi_\Sigma)$ and $\text{tr}((\mathbf{H}^{-1})^{\otimes r}\mathbf{R}(\mathbf{D}^{\otimes r}\phi))$. These are provided in the following lemma.

Lemma 3. *For any symmetric positive definite matrix \mathbf{H} , we have*

- (i) $\mathbf{R}_{\phi,\mathbf{H},r}(\phi_\Sigma) = 2^{d/2+r}\mathbf{R}(\mathbf{D}^{\otimes r}\phi_{\mathbf{H}+2\Sigma})$,
- (ii) $\text{tr}((\mathbf{H}^{-1})^{\otimes r}\mathbf{R}(\mathbf{D}^{\otimes r}\phi_\Sigma)) = 2^{-(d+r)}\pi^{-d/2}|\Sigma|^{-1/2}\mu_r(\Sigma^{1/2}\mathbf{H}\Sigma^{1/2})$.

From (i) and (ii) we obtain

- (iii) $\text{tr}\mathbf{R}_{\phi,\mathbf{H},r}(\phi_\Sigma) = (2\pi)^{-d/2}|\mathbf{H} + 2\Sigma|^{-1/2}\mu_r(\mathbf{H} + 2\Sigma)$,
- (iv) $\text{tr}((\mathbf{H}^{-1})^{\otimes r}\mathbf{R}(\mathbf{D}^{\otimes r}\phi)) = 2^{-(d+r)}\pi^{-d/2}\mu_r(\mathbf{H})$.

Proof. (i) Reasoning as in Chacón and Duong (2010), it is easy to check that

$$\text{vec}\mathbf{R}(\mathbf{D}^{\otimes r}\phi_\Sigma) = (-1)^r\mathbf{D}^{\otimes 2r}\phi_{2\Sigma}(\mathbf{0}) = (-1)^r2^{-(d/2+r)}\mathbf{D}^{\otimes 2r}\phi_\Sigma(\mathbf{0}).$$

With this in mind, an element-wise application of some of the results in Appendix C of Wand and Jones (1995) leads to

$$\begin{aligned} \text{vec}\mathbf{R}_{\phi,\mathbf{H},r}(\phi_\Sigma) &= \text{vec}\int_{\mathbb{R}^d}(\phi_{\mathbf{H}}*\mathbf{D}^{\otimes r}\phi_\Sigma)(\mathbf{x})\mathbf{D}^{\otimes r}\phi_\Sigma(\mathbf{x})^T d\mathbf{x} \\ &= \text{vec}\int_{\mathbb{R}^d}\mathbf{D}^{\otimes r}\phi_{\mathbf{H}+\Sigma}(\mathbf{x})\mathbf{D}^{\otimes r}\phi_\Sigma(\mathbf{x})^T d\mathbf{x} \\ &= \int_{\mathbb{R}^d}\mathbf{D}^{\otimes r}\phi_\Sigma(\mathbf{x})\otimes\mathbf{D}^{\otimes r}\phi_{\mathbf{H}+\Sigma}(\mathbf{x})d\mathbf{x} \\ &= (-1)^r\int_{\mathbb{R}^d}\mathbf{D}^{\otimes 2r}\phi_\Sigma(\mathbf{x})\phi_{\mathbf{H}+\Sigma}(\mathbf{x})d\mathbf{x} \\ &= (-1)^r\mathbf{D}^{\otimes 2r}\phi_{\mathbf{H}+2\Sigma}(\mathbf{0}) \\ &= 2^{d/2+r}\text{vec}\mathbf{R}(\mathbf{D}^{\otimes r}\phi_{\mathbf{H}+2\Sigma}), \end{aligned}$$

which yields the result.

(ii) Chacón and Duong (2010) also show that

$$\text{vec}\mathbf{R}(\mathbf{D}^{\otimes r}\phi_\Sigma) = 2^{-(d+r)}\pi^{-d/2}\text{OF}(2r)|\Sigma|^{-1/2}\mathcal{S}_{d,2r}(\text{vec}\Sigma^{-1})^{\otimes r}. \tag{A.2}$$

Moreover, it is not hard to check that the symmetrizer matrix satisfies $\mathcal{S}_{d,2r}\text{vec}[(\mathbf{H}^{-1})^{\otimes r}] = \mathcal{S}_{d,2r}(\text{vec}\mathbf{H}^{-1})^{\otimes r}$. This is because $(\text{vec}\mathbf{H}^{-1})^{\otimes r}$ can be

obtained from $\text{vec}[(\mathbf{H}^{-1})^{\otimes r}]$ by multiplying it by Kronecker products of commutation and identity matrices, and multiplication of this kind of matrices by the symmetrizer matrix has no effect, as seen from part (iv) of Theorem 1 in Schott (2003). Therefore, if \mathbf{z} denotes a d -variate vector with standard normal distribution and $\mathbf{x} = \boldsymbol{\Sigma}^{-1/2}\mathbf{z}$, then

$$\begin{aligned} & \text{tr}((\mathbf{H}^{-1})^{\otimes r} \mathbf{R}(\mathbf{D}^{\otimes r} \phi_{\boldsymbol{\Sigma}})) \\ &= \text{vec}^T[(\mathbf{H}^{-1})^{\otimes r}] \text{vec} \mathbf{R}(\mathbf{D}^{\otimes r} \phi_{\boldsymbol{\Sigma}}) \\ &= 2^{-(d+r)} \pi^{-d/2} \text{OF}(2r) |\boldsymbol{\Sigma}|^{-1/2} \text{vec}^T[(\mathbf{H}^{-1})^{\otimes r}] \mathcal{S}_{d,2r}(\text{vec} \boldsymbol{\Sigma}^{-1})^{\otimes r} \\ &= 2^{-(d+r)} \pi^{-d/2} \text{OF}(2r) |\boldsymbol{\Sigma}|^{-1/2} (\text{vec}^T \mathbf{H}^{-1})^{\otimes r} \mathcal{S}_{d,2r}(\text{vec} \boldsymbol{\Sigma}^{-1})^{\otimes r} \\ &= 2^{-(d+r)} \pi^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \mathbb{E}[(\mathbf{x}^T \mathbf{H}^{-1} \mathbf{x})^r] \\ &= 2^{-(d+r)} \pi^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \mathbb{E}[(\mathbf{z}^T \boldsymbol{\Sigma}^{-1/2} \mathbf{H}^{-1} \boldsymbol{\Sigma}^{-1/2} \mathbf{z})^r]. \end{aligned}$$

Here, the fourth line follows from Theorem 1 in Holmquist (1996b).

A.2.2. Proof of Theorem 6

The proof of Theorem 6 starts from the AMISE expression given in Theorem 2. The term appearing in the asymptotic integrated variance was already computed in Lemma 3. For the asymptotic integrated squared bias, it is clear that for the normal kernel we have $m_2(K) = 1$. From (A.2) and the results in Holmquist (1996a) it follows that

$$\text{vec} \mathbf{R}(\mathbf{D}^{\otimes r} \phi_{\boldsymbol{\Sigma}}) = 2^{-(d+r)} \pi^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} (\boldsymbol{\Sigma}^{-1/2})^{\otimes 2r} \mathbb{E}[\mathbf{z}^{\otimes 2r}],$$

with \mathbf{z} a d -variate standard normal random vector. Or, in matrix form,

$$\mathbf{R}(\mathbf{D}^{\otimes r} \phi_{\boldsymbol{\Sigma}}) = 2^{-(d+r)} \pi^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} (\boldsymbol{\Sigma}^{-1/2})^{\otimes r} \mathbb{E}[(\mathbf{z}\mathbf{z}^T)^{\otimes r}] (\boldsymbol{\Sigma}^{-1/2})^{\otimes r}.$$

Therefore, using $(\boldsymbol{\Sigma}^{-1/2} \otimes \boldsymbol{\Sigma}^{-1/2}) \text{vec} \mathbf{H} = \text{vec}(\boldsymbol{\Sigma}^{-1/2} \mathbf{H} \boldsymbol{\Sigma}^{-1/2})$ and some other matrix results from Magnus and Neudecker (1999, p. 48), we come to

$$\begin{aligned} & \text{tr}[(\mathbf{I}_{dr} \otimes \text{vec}^T \mathbf{H}) \mathbf{R}(\mathbf{D}^{\otimes(r+2)} \phi_{\boldsymbol{\Sigma}}) (\mathbf{I}_{dr} \otimes \text{vec} \mathbf{H})] \\ &= 2^{-(d+r+2)} \pi^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \text{tr} \left[\{(\boldsymbol{\Sigma}^{-1})^{\otimes r} \otimes (\text{vec} \mathbf{B} \text{vec}^T \mathbf{B})\} \mathbb{E}[(\mathbf{z}\mathbf{z}^T)^{\otimes(r+2)}] \right], \end{aligned}$$

with $\mathbf{B} = \boldsymbol{\Sigma}^{-1/2} \mathbf{H} \boldsymbol{\Sigma}^{-1/2}$. Now, the trace in the right hand side can be written as

$$\begin{aligned} & \mathbb{E} \text{tr} \left[(\boldsymbol{\Sigma}^{-1} \mathbf{z}\mathbf{z}^T)^{\otimes r} \otimes \{\text{vec} \mathbf{B} \text{vec}^T \mathbf{B} (\mathbf{z}\mathbf{z}^T)^{\otimes 2}\} \right] \\ &= \mathbb{E} \left[\text{tr}^r(\boldsymbol{\Sigma}^{-1} \mathbf{z}\mathbf{z}^T) \text{tr} \{ \text{vec} \mathbf{B} \text{vec}^T (\mathbf{z}\mathbf{z}^T \mathbf{B} \mathbf{z}\mathbf{z}^T) \} \right] \end{aligned}$$

$$\begin{aligned} &= \mathbb{E} \left[(\mathbf{z}^T \boldsymbol{\Sigma}^{-1} \mathbf{z})^r \{ \text{vec}^T(\mathbf{z} \mathbf{z}^T \mathbf{B} \mathbf{z} \mathbf{z}^T) \text{vec} \mathbf{B} \} \right] \\ &= \mathbb{E} \left[(\mathbf{z}^T \boldsymbol{\Sigma}^{-1} \mathbf{z})^r (\mathbf{z}^T \mathbf{B} \mathbf{z})^2 \right]. \end{aligned}$$

This yields the proof for the AMISE formula.

If we evaluate the AMISE formula in Theorem 6 at $\mathbf{H} = c\boldsymbol{\Sigma}$ for some $c > 0$, we obtain

$$\text{AMISE}(c\boldsymbol{\Sigma}) = 2^{-(d+r)} \pi^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \left\{ n^{-1} c^{-(d+2r)/2} \mu_r(\boldsymbol{\Sigma}) + \frac{1}{16} c^2 \mu_{r,2}(\boldsymbol{\Sigma}, \mathbf{I}_d) \right\}.$$

But we show below that

$$\mu_{r,2}(\boldsymbol{\Sigma}, \mathbf{I}_d) = (d + 2r + 2)(d + 2r) \mu_r(\boldsymbol{\Sigma}), \tag{A.3}$$

leading to

$$\text{AMISE}(c\boldsymbol{\Sigma}) = 2^{-(d+r)} \pi^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \mu_r(\boldsymbol{\Sigma}) \left\{ n^{-1} c^{-(d+2r)/2} + \frac{1}{16} c^2 (d+2r+2)(d+2r) \right\};$$

this function is minimized by setting

$$c = \left(\frac{4}{(d + 2r + 2)n} \right)^{2/(d+2r+4)}.$$

Therefore, to finish the proof the only thing left is to show (A.3).

This task, however, is harder than it may seem at first sight. It is relatively easy if $\boldsymbol{\Sigma} = \mathbf{I}_d$ because, in this case, it suffices to show that $\mu_{r+1}(\mathbf{I}_d) = (d + 2r) \mu_r(\mathbf{I}_d)$, which is an immediate consequence of the recursive formula (3.1). Therefore, to show (A.3) we need a recursive formula similar to (3.1), but for the joint moments $\mu_{r,s}(\mathbf{A}, \mathbf{B})$. To that end, we first derive a technical lemma.

Lemma 4. Consider $g_\alpha(t) \equiv g_\alpha(t; \mathbf{A}, \mathbf{B}, \mathbf{C}) = \text{tr} [\{ \mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1} \}^\alpha]$ for suitable matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}$ and arbitrary $\alpha \in \mathbb{N}$. Then

$$g_\alpha^{(p)}(t; \mathbf{A}, \mathbf{B}, \mathbf{C}) = (-1)^p \frac{(\alpha + p - 1)!}{(\alpha - 1)!} \text{tr} [\{ \mathbf{A}(\mathbf{C} + t\mathbf{A})^{-1} \}^p \{ \mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1} \}^\alpha],$$

so that $g_\alpha^{(p)}(0; \mathbf{A}, \mathbf{B}, \mathbf{C}) = (-1)^p \frac{(\alpha + p - 1)!}{(\alpha - 1)!} \text{tr} [(\mathbf{A}\mathbf{C}^{-1})^p (\mathbf{B}\mathbf{C}^{-1})^\alpha]$.

Proof. The result is proved by induction on p . For $p = 1$, noting that the differential of $\mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1}$ is $d[\mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1}] = -\mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1} \mathbf{A}(\mathbf{C} + t\mathbf{A})^{-1} dt$, we have

$$\begin{aligned} &d \text{tr} [\{ \mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1} \}^\alpha] \\ &= \text{tr} d[\{ \mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1} \}^\alpha] \end{aligned}$$

$$\begin{aligned} &= \text{tr} \sum_{i=1}^{\alpha} \{\mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1}\}^{i-1} \cdot d[\mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1}] \cdot \{\mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1}\}^{\alpha-i} \\ &= -\text{tr} \sum_{i=1}^{\alpha} \{\mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1}\}^i \cdot \mathbf{A}(\mathbf{C} + t\mathbf{A})^{-1} \cdot \{\mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1}\}^{\alpha-i} dt \\ &= -\alpha \text{tr} [\mathbf{A}(\mathbf{C} + t\mathbf{A})^{-1} \{\mathbf{B}(\mathbf{C} + t\mathbf{A})^{-1}\}^{\alpha}] dt, \end{aligned}$$

and we are done. The case of arbitrary p follows from $g_{\alpha}^{(p)}(t) = \frac{d}{dt} g_{\alpha}^{(p-1)}(t)$.

Theorem 9. *We can write*

$$\begin{aligned} &\mu_{r,s}(\mathbf{A}, \mathbf{B}) \\ &= \sum_{i=0}^r \sum_{j=0}^{s-1} \binom{r}{i} \binom{s-1}{j} (r+s-i-j-1)! 2^{r+s-i-j-1} \text{tr}(\mathbf{A}^{-(r-i)} \mathbf{B}^{-(s-j)}) \mu_{i,j}(\mathbf{A}, \mathbf{B}). \end{aligned}$$

Proof. For ease of notation we prove the result for $\mu_{r,s}(\mathbf{A}^{-1}, \mathbf{B}^{-1})$; that is, we show that, for $q_{\mathbf{A}} = \mathbf{z}^T \mathbf{A} \mathbf{z}$ and $q_{\mathbf{B}} = \mathbf{z}^T \mathbf{B} \mathbf{z}$ with \mathbf{z} a d -variate standard normal random vector,

$$\mathbb{E}[q_{\mathbf{A}}^r q_{\mathbf{B}}^s] = \sum_{i=0}^r \sum_{j=0}^{s-1} \binom{r}{i} \binom{s-1}{j} (r+s-i-j-1)! 2^{r+s-i-j-1} \text{tr}(\mathbf{A}^{r-i} \mathbf{B}^{s-j}) \mathbb{E}[q_{\mathbf{A}}^i q_{\mathbf{B}}^j].$$

It is well known that the joint moment generating function of $q_{\mathbf{A}}$ and $q_{\mathbf{B}}$ is

$$M(t_1, t_2) = \mathbb{E}[e^{t_1 q_{\mathbf{A}} + t_2 q_{\mathbf{B}}}] = |\mathbf{I}_d - 2t_1 \mathbf{A} - 2t_2 \mathbf{B}|^{-1/2},$$

see Magnus (1986). From that, we can write

$$\mathbb{E}[q_{\mathbf{A}}^r q_{\mathbf{B}}^s] = \frac{\partial^{r+s} M}{\partial t_1^r \partial t_2^s}(0, 0),$$

so that all we need is to find a recursive formula for the partial derivatives of M . With the notations of the previous lemma, it is easy to show that

$$\frac{\partial M}{\partial t_2}(t_1, t_2) = M(t_1, t_2) \cdot g_1(t_2; -2\mathbf{B}, \mathbf{B}, \mathbf{I}_d - 2t_1 \mathbf{A}).$$

This way, using the formulas for the derivatives of g_1 and Leibniz formula for the derivatives of a product,

$$\begin{aligned} &\frac{\partial^s M}{\partial t_2^s}(t_1, t_2) \\ &= \frac{\partial^{s-1}}{\partial t_2^{s-1}} \left(M(t_1, t_2) \cdot g_1(t_2; -2\mathbf{B}, \mathbf{B}, \mathbf{I}_d - 2t_1 \mathbf{A}) \right) \end{aligned}$$

$$\begin{aligned}
 &= \sum_{j=1}^{s-1} \binom{s-1}{j} \frac{\partial^j M}{\partial t_2^j}(t_1, t_2) \cdot g_1^{(s-j-1)}(t_2; -2\mathbf{B}, \mathbf{B}, \mathbf{I}_d - 2t_1\mathbf{A}) \\
 &= \sum_{j=1}^{s-1} \binom{s-1}{j} (s-j-1)! 2^{s-j-1} \frac{\partial^j M}{\partial t_2^j}(t_1, t_2) \cdot \text{tr} [\{\mathbf{B}(\mathbf{I}_d - 2t_1\mathbf{A} - 2t_2\mathbf{B})^{-1}\}^{s-j}] \\
 &= \sum_{j=1}^{s-1} \frac{(s-1)!}{j!} 2^{s-j-1} \frac{\partial^j M}{\partial t_2^j}(t_1, t_2) \cdot g_{s-j}(t_1; -2\mathbf{A}, \mathbf{B}, \mathbf{I}_d - 2t_2\mathbf{B}).
 \end{aligned}$$

Now, if we compute the r th partial derivative with respect to t_1 we have

$$\begin{aligned}
 &\frac{\partial^{r+s} M}{\partial t_1^r \partial t_2^s}(t_1, t_2) \\
 &= \sum_{i=0}^r \sum_{j=1}^{s-1} \binom{r}{i} \frac{(s-1)!}{j!} 2^{s-j-1} \frac{\partial^{i+j} M}{\partial t_1^i \partial t_2^j}(t_1, t_2) \cdot g_{s-j}^{(r-i)}(t_1; -2\mathbf{A}, \mathbf{B}, \mathbf{I}_d - 2t_2\mathbf{B}).
 \end{aligned}$$

Substituting (t_1, t_2) in this expression for $(0, 0)$ and using again the previous lemma, we get the desired formula.

As a consequence of this result, we are able to prove formula (A.3).

Corollary 7. *For any symmetric matrix \mathbf{A} we have $\mu_{r,2}(\mathbf{A}, \mathbf{I}_d) = (d + 2r + 2)(d + 2r)\mu_r(\mathbf{A})$.*

Proof. First, notice that from the previous theorem, taking into account that $\mathbf{A}^0 = \mathbf{I}_d$,

$$\begin{aligned}
 \mu_{r,1}(\mathbf{A}, \mathbf{I}_d) &= \sum_{i=0}^r \binom{r}{i} (r-i)! 2^{r-i} \text{tr}(\mathbf{A}^{-(r-i)}) \mu_{i,0}(\mathbf{A}, \mathbf{I}_d) \\
 &= d\mu_r(\mathbf{A}) + \sum_{i=0}^{r-1} \frac{r!}{i!} 2^{r-i} \text{tr}(\mathbf{A}^{-(r-i)}) \mu_i(\mathbf{A}) \\
 &= (d + 2r)\mu_r(\mathbf{A}),
 \end{aligned}$$

where the last equality follows from (3.1). Using this and Theorem 9,

$$\begin{aligned}
 &\mu_{r,2}(\mathbf{A}, \mathbf{I}_d) \\
 &= \sum_{i=0}^r \sum_{j=0}^1 \binom{r}{i} \binom{1}{j} (r-i-j+1)! 2^{r-i-j+1} \text{tr}(\mathbf{A}^{-(r-i)}) \mu_{i,j}(\mathbf{A}, \mathbf{I}_d) \\
 &= \sum_{i=0}^r \left[\binom{r}{i} (r-i+1)! 2^{r-i+1} \text{tr}(\mathbf{A}^{-(r-i)}) \mu_{i,0}(\mathbf{A}, \mathbf{I}_d) \right.
 \end{aligned}$$

$$\begin{aligned}
 & + \binom{r}{i} (r-i)! 2^{r-i} \operatorname{tr}(\mathbf{A}^{-(r-i)}) \mu_{i,1}(\mathbf{A}, \mathbf{I}_d) \Big] \\
 & = \sum_{i=0}^r \left[\frac{r!}{i!} (r-i+1) 2^{r-i+1} \operatorname{tr}(\mathbf{A}^{-(r-i)}) \mu_i(\mathbf{A}) + \frac{r!}{i!} 2^{r-i} \operatorname{tr}(\mathbf{A}^{-(r-i)}) (d+2i) \mu_i(\mathbf{A}) \right] \\
 & = \sum_{i=0}^r \frac{r!}{i!} [2(r-i+1) + d+2i] 2^{r-i} \operatorname{tr}(\mathbf{A}^{-(r-i)}) \mu_i(\mathbf{A}) \\
 & = (d+2r+2) \mu_{r,1}(\mathbf{A}, \mathbf{I}_d) \\
 & = (d+2r+2)(d+2r) \mu_r(\mathbf{A}).
 \end{aligned}$$

A.2.3. Proof of Theorem 7

Proof. From the proof of Theorem 5, we have that for $K = \phi$, $\mathbf{R}_{\phi*\phi, \mathbf{H}, r}(f) = \mathbf{R}_{\phi, 2\mathbf{H}, r}(f)$ and $\mathbf{R}(\mathbf{D}^{\otimes r} f) = \mathbf{R}_{\phi, \mathbf{0}, r}(f)$. Combining this with Theorem 1 and part *iv*) of Lemma 3 we come to

$$\text{MISE}\{\widehat{(\mathbf{D}^{\otimes r} f)}(\cdot; \mathbf{H})\} = 2^{-r} n^{-1} (4\pi)^{-d/2} |\mathbf{H}|^{-1/2} \mu_r(\mathbf{H}) + (1-n^{-1}) \varpi_2 - 2\varpi_1 + \varpi_0,$$

where $\varpi_a = \operatorname{tr} \mathbf{R}_{\phi, a\mathbf{H}, r}(f) = (\operatorname{vec}^T \mathbf{I}_{dr}) \operatorname{vec} \mathbf{R}_{\phi, a\mathbf{H}, r}(f)$. Now,

$$\begin{aligned}
 \operatorname{vec} \mathbf{R}_{\phi, a\mathbf{H}, r}(f) & = \operatorname{vec} \int \phi_{a\mathbf{H}} * \mathbf{D}^{\otimes r} f(\mathbf{x}) \mathbf{D}^{\otimes r} f(\mathbf{x})^T d\mathbf{x} \\
 & = \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} \int \mathbf{D}^{\otimes r} \phi_{\Sigma_{\ell'}}(\mathbf{x} - \boldsymbol{\mu}_{\ell'}) (\mathbf{x}) \otimes \mathbf{D}^{\otimes r} \phi_{a\mathbf{H} + \Sigma_\ell}(\mathbf{x} - \boldsymbol{\mu}_\ell) d\mathbf{x} \\
 & = \sum_{\ell, \ell'=1}^k w_\ell w_{\ell'} (-1)^r \mathbf{D}^{\otimes 2r} \phi_{a\mathbf{H} + \Sigma_\ell + \Sigma_{\ell'}}(\boldsymbol{\mu}_\ell - \boldsymbol{\mu}_{\ell'}),
 \end{aligned}$$

so that we can write $\varpi_a = \mathbf{w}^T \boldsymbol{\Omega}_a \mathbf{w}$, where $\boldsymbol{\Omega}_a$ is the $k \times k$ matrix with (ℓ, ℓ') entry given by $(\boldsymbol{\Omega}_a)_{\ell, \ell'} = (-1)^r (\operatorname{vec}^T \mathbf{I}_{dr}) \mathbf{D}^{\otimes 2r} \phi_{a\mathbf{H} + \Sigma_\ell + \Sigma_{\ell'}}(\boldsymbol{\mu}_\ell - \boldsymbol{\mu}_{\ell'})$.

The second expression of $(\boldsymbol{\Omega}_a)_{\ell, \ell'}$ is derived from an identity in Holmquist (1996a),

$$\mathbf{D}^{\otimes 2r} \phi_\Sigma(\boldsymbol{\mu}) = \phi_\Sigma(\boldsymbol{\mu}) (\boldsymbol{\Sigma}^{-1})^{\otimes 2r} \mathcal{S}_{d, 2r} \sum_{j=0}^r (-1)^j \text{OF}(2j) \binom{2r}{2j} (\boldsymbol{\mu}^{\otimes (2r-2j)} \otimes (\operatorname{vec} \boldsymbol{\Sigma})^{\otimes j}), \tag{A.4}$$

with the use of $(\operatorname{vec}^T \mathbf{I}_{dr}) (\boldsymbol{\Sigma}^{-1})^{\otimes 2r} = \operatorname{vec}^T (\boldsymbol{\Sigma}^{-2})^{\otimes r}$.

A.2.4. Proof of Theorem 8

Proof. To determine the AMISE formula, it suffices to find an expression for the asymptotic integrated squared bias

$$\text{AIB}^2(\mathbf{H}) = \frac{1}{4} \operatorname{tr} [(\mathbf{I}_{dr} \otimes \operatorname{vec}^T \mathbf{H}) \mathbf{R}(\mathbf{D}^{\otimes (r+2)} f) (\mathbf{I}_{dr} \otimes \operatorname{vec} \mathbf{H})]$$

$$= \frac{1}{4} \text{vec}^T (\mathbf{I}_{dr} \otimes (\text{vec } \mathbf{H} \text{vec}^T \mathbf{H})) \text{vec } \mathbf{R}(\mathbf{D}^{\otimes(r+2)} f).$$

We can write

$$\text{vec } \mathbf{R}(\mathbf{D}^{\otimes(r+2)} f) = \text{vec } \mathbf{R}_{\phi, 0, \mathbf{H}, r+2}(f) = \sum_{\ell, \ell'=1}^k w_{\ell} w_{\ell'} (-1)^{r+2} \mathbf{D}^{\otimes 2r+4} \phi_{\Sigma_{\ell} + \Sigma_{\ell'}}(\boldsymbol{\mu}_{\ell} - \boldsymbol{\mu}_{\ell'})$$

so that $4\text{AIB}^2(\mathbf{H}) = \mathbf{w}^T \tilde{\boldsymbol{\Omega}} \mathbf{w}$, with $\boldsymbol{\mu}_{\ell\ell'} = \boldsymbol{\mu}_{\ell} - \boldsymbol{\mu}_{\ell'}$, $\Sigma_{\ell\ell'} = \Sigma_{\ell} + \Sigma_{\ell'}$,

$$\begin{aligned} \tilde{\boldsymbol{\Omega}}_{\ell, \ell'} &= (-1)^r \text{vec}^T (\mathbf{I}_{dr} \otimes (\text{vec } \mathbf{H} \text{vec}^T \mathbf{H})) \mathbf{D}^{\otimes 2r+4} \phi_{\Sigma_{\ell\ell'}}(\boldsymbol{\mu}_{\ell\ell'}) \\ &= (-1)^r \phi_{\Sigma_{\ell\ell'}}(\boldsymbol{\mu}_{\ell\ell'}) [(\text{vec}^T \mathbf{I}_d)^{\otimes r} \otimes (\text{vec}^T \mathbf{H})^{\otimes 2}] (\Sigma_{\ell\ell'}^{-1})^{\otimes (2r+4)} \mathcal{S}_{d, 2r+4} \\ &\quad \times \sum_{j=0}^{r+2} (-1)^j \text{OF}(2j) \binom{2r+4}{2j} [\boldsymbol{\mu}_{\ell\ell'}^{\otimes (2r-2j+4)} \otimes (\text{vec } \Sigma_{\ell\ell'})^{\otimes j}] \\ &= (-1)^r \phi_{\Sigma_{\ell\ell'}}(\boldsymbol{\mu}_{\ell\ell'}) [(\text{vec}^T \Sigma_{\ell\ell'}^{-2})^{\otimes r} \otimes (\text{vec}^T (\Sigma_{\ell\ell'}^{-1} \mathbf{H} \Sigma_{\ell\ell'}^{-1}))^{\otimes 2}] \mathcal{S}_{d, 2r+4} \\ &\quad \times \sum_{j=0}^{r+2} (-1)^j \text{OF}(2j) \binom{2r+4}{2j} [\boldsymbol{\mu}_{\ell\ell'}^{\otimes (2r-2j+4)} \otimes (\text{vec } \Sigma_{\ell\ell'})^{\otimes j}], \end{aligned}$$

using (A.4).

A.3. Proof of the results in Section 4

A.3.1. Proof of Corollary 3

Proof. From Theorem 6, $\mathbf{H}_{\text{AMISE}} = c_{\text{AMISE}} \boldsymbol{\Sigma}$, where $c_{\text{AMISE}} = \{4/[(d + 2r + 2)n]\}^{2/(d+2r+4)}$. Substituting this into the equation immediately following below (A.3),

$$\begin{aligned} &\min_{\mathbf{H} \in \mathcal{F}} \text{AMISE}(\mathbf{H}) \\ &= 2^{-(d+r)} \pi^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \mu_r(\boldsymbol{\Sigma}) c_{\text{AMISE}}^2 [n^{-1} c_{\text{AMISE}}^{-(d+2r+4)/2} + \frac{1}{16} (d+2r+2)(d+2r)] \\ &= 2^{-(d+r)} \pi^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \mu_r(\boldsymbol{\Sigma}) \left(\frac{4}{d+2r+2} \right)^{4/(d+2r+4)} \\ &\quad \times \left[\frac{1}{4} (d+2r+2) + \frac{1}{16} (d+2r+2)(d+2r) \right] n^{-4/(d+2r+4)} \\ &= 2^{-(d+r+4)} \pi^{-d/2} (d+2r+2)(d+2r+4) \left(\frac{4}{d+2r+2} \right)^{4/(d+2r+4)} \\ &\quad \times |\boldsymbol{\Sigma}|^{-1/2} \mu_r(\boldsymbol{\Sigma}) n^{-4/(d+2r+4)} \\ &= 2^{-(d+r+4)} 2^{8/(d+2r+4)} \pi^{-d/2} (d+2r+4)(d+2r+2)^{(d+2r)/(d+2r+4)} \\ &\quad \times |\boldsymbol{\Sigma}|^{-1/2} \mu_r(\boldsymbol{\Sigma}) n^{-4/(d+2r+4)}. \end{aligned}$$

A.3.2. Proof of Corollary 4

Proof. Substituting $\mathbf{H} = h^2\mathbf{I}_d$ into the AMISE formula in Theorem 6,

$$\begin{aligned} \text{AMISE}(h^2\mathbf{I}_d) &= 2^{-(d+r)}\pi^{-d/2}\{n^{-1}h^{-d-2r}\mu_r(\mathbf{I}_d) + \frac{1}{16}|\boldsymbol{\Sigma}|^{-1/2}\mu_{r,2}(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}^{1/2}(h^{-2}\mathbf{I}_d)\boldsymbol{\Sigma}^{1/2})\} \\ &= 2^{-(d+r)}\pi^{-d/2}\{n^{-1}h^{-d-2r}\mu_r(\mathbf{I}_d) + \frac{1}{16}|\boldsymbol{\Sigma}|^{-1/2}\mu_{r+2}(\boldsymbol{\Sigma})h^4\}, \end{aligned}$$

since $\mu_{r,2}(\boldsymbol{\Sigma}, \boldsymbol{\Sigma}) = \mu_{r+2}(\boldsymbol{\Sigma})$. Differentiating with respect to h and setting to zero, $-(d+2r)n^{-1}h^{-d-2r-1}\mu_r(\mathbf{I}_d) + \frac{1}{4}|\boldsymbol{\Sigma}|^{-1/2}\mu_{r+2}(\boldsymbol{\Sigma})h^3 = 0$ gives

$$h_{\text{AMISE}} = \left(\frac{4(d+2r)|\boldsymbol{\Sigma}|^{1/2}\mu_r(\mathbf{I}_d)}{\mu_{r+2}(\boldsymbol{\Sigma})n} \right)^{1/(d+2r+4)}.$$

The minimal AMISE is

$$\begin{aligned} \min_{\mathbf{H} \in \mathcal{I}} \text{AMISE}(\mathbf{H}) &= 2^{-(d+r)}\pi^{-d/2}h_{\text{AMISE}}^4\{n^{-1}h_{\text{AMISE}}^{-d-2r-4} + \frac{1}{16}|\boldsymbol{\Sigma}|^{-1/2}\mu_{r+2}(\boldsymbol{\Sigma})\} \\ &= 2^{-(d+r)}\pi^{-d/2} \left(\frac{4(d+2r)|\boldsymbol{\Sigma}|^{1/2}\mu_r(\mathbf{I}_d)}{\mu_{r+2}(\boldsymbol{\Sigma})} \right)^{4/(d+2r+4)} \left(\frac{\mu_{r+2}(\boldsymbol{\Sigma})}{4(d+2r)|\boldsymbol{\Sigma}|^{1/2}} + \frac{\mu_{r+2}(\boldsymbol{\Sigma})}{16|\boldsymbol{\Sigma}|^{1/2}} \right) \\ &\quad \times n^{-4/(d+2r+4)} \\ &= 2^{-(d+r+4)}\pi^{-d/2}|\boldsymbol{\Sigma}|^{-1/2}\mu_{r+2}(\boldsymbol{\Sigma}) \\ &\quad \times \left(\frac{d+2r+4}{d+2r} \right) \left(\frac{4(d+2r)|\boldsymbol{\Sigma}|^{1/2}\mu_r(\mathbf{I}_d)}{\mu_{r+2}(\boldsymbol{\Sigma})} \right)^{4/(d+2r+4)} n^{-4/(d+2r+4)} \\ &= 2^{-(d+r+4)}2^{8/(d+2r+4)}\pi^{-d/2}(d+2r+4)(d+2r)^{-1} \\ &\quad \times \{|\boldsymbol{\Sigma}|^{-(d+2r)/2}\mu_{r+2}(\boldsymbol{\Sigma})^{d+2r}\mu_{r+1}(\mathbf{I}_d)^4\}^{1/(d+2r+4)} n^{-4/(d+2r+4)}, \end{aligned}$$

since $\mu_{r+1}(\mathbf{I}_d) = (d+2r)\mu_r(\mathbf{I}_d)$.

A.3.3. Proof of Corollary 5

Proof. From Corollaries 3 and 4, the ARE is

$$\begin{aligned} \text{ARE}(\mathcal{F} : \mathcal{I}) &= \frac{(d+2r+2)^{(d+2r)/4}|\boldsymbol{\Sigma}|^{-(d+2r+4)/8}\mu_r(\boldsymbol{\Sigma})^{(d+2r+4)/4}}{(d+2r)^{-(d+2r+4)/4}|\boldsymbol{\Sigma}|^{-(d+2r)/8}\mu_{r+2}(\boldsymbol{\Sigma})^{(d+2r)/4}\mu_{r+1}(\mathbf{I}_d)} \\ &= [(d+2r+2)(d+2r)]^{(d+2r)/4}|\boldsymbol{\Sigma}|^{-1/2}\mu_r(\boldsymbol{\Sigma})^{(d+2r+4)/4}\mu_{r+2}(\boldsymbol{\Sigma})^{-(d+2r)/4}\mu_r(\mathbf{I}_d)^{-1} \end{aligned}$$

since $\mu_{r+1}(\mathbf{I}_d) = (d+2r)\mu_r(\mathbf{I}_d)$.

A.3.4. Proof of Corollary 6

Proof. Let the variance be $\Sigma = \sigma^2 \Sigma_\rho$ where $\Sigma_\rho = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$. With this form of the variance, using symmetry arguments, the bandwidth that minimizes the AMISE is of the form $\mathbf{H}_{\text{AMISE}} = c\mathbf{I}_2$, for some positive constant c . So we can apply Corollary 5. Thus

$$\begin{aligned} \text{ARE}(\mathcal{F} : \mathcal{D}) &= \text{ARE}(\mathcal{F} : \mathcal{I}) \\ &= [4(r+1)(r+2)]^{(r+1)/2} |\Sigma|^{-1/2} \mu_r(\Sigma)^{(r+3)/2} \mu_{r+2}(\Sigma)^{-(r+1)/2} \mu_r(\mathbf{I}_2)^{-1} \\ &\quad \times \mu_{r+2}(\Sigma)^{-(r+1)/2} \mu_r(\mathbf{I}_2)^{-1} \\ &= [4(r+1)(r+2)]^{(r+1)/2} \sigma^{-2} (1-\rho^2)^{-1/2} \sigma^{-r(r+3)} (1-\rho^2)^{-r(r+3)/2} Q(r, \rho)^{(r+3)/2} \\ &\quad \times \sigma^{(r+1)(r+2)} (1-\rho^2)^{(r+1)(r+2)/2} Q(r+2, \rho)^{-(r+1)/2} Q(r, 0)^{-1} \\ &= \frac{(1-\rho^2)^{1/2} Q(r, \rho)^{(r+3)/2}}{Q(r, 0) Q(r+2, \rho)^{(r+1)/2}}, \end{aligned}$$

where $Q(r, \rho) = (1-\rho^2)^r \mu_r(\Sigma_\rho) = \sigma^{2r} (1-\rho^2)^r \mu_r(\Sigma)$. Since $\Sigma^{-1} = \sigma^{-2} (1-\rho^2)^{-1} \begin{bmatrix} 1 & -\rho \\ -\rho & 1 \end{bmatrix}$, and with z_1, z_2 independent standard normal random variables,

$$\begin{aligned} Q(r, \rho) &= \mathbb{E}\{(z_1^2 + z_2^2 - 2\rho z_1 z_2)^2\} \\ &= \mathbb{E}\left\{ \sum_{j=0}^r \binom{r}{j} (z_1^2 - 2\rho z_1 z_2)^j (z_2)^{2(r-j)} \right\} \\ &= \mathbb{E}\left\{ \sum_{j=0}^r \binom{r}{j} \sum_{j'=0}^j \binom{j}{j'} z_1^{2j'} (-2\rho z_1 z_2)^{j-j'} (z_2)^{2(r-j)} \right\} \\ &= \mathbb{E}\left\{ \sum_{j=0}^r \binom{r}{j} \sum_{j'=0}^j \binom{j}{j'} (-2\rho)^{j-j'} z_1^{j+j'} z_2^{2r-j-j'} \right\} \\ &= \sum_{j=0}^r \sum_{j'=0}^j \binom{r}{j} \binom{j}{j'} (-2\rho)^{j-j'} m_{j+j'} m_{2r-j-j'}, \end{aligned}$$

where $m_k = (1/2)\{(-1)^k + 1\} \text{OF}(k)$ is the k th central moment of a standard normal variable.

References

Baxandall, P. and Liebeck, H. (1986). *Vector Calculus*. Clarendon Press, Oxford.

- Bhattacharya, P. K. (1967). Estimation of a probability density function and its derivatives. *Sankhyā Ser. A* **29**, 373-382.
- Brinkman, R. R., Gasparetto, M., Lee, S.-J. J., Ribickas, A. J., Perkins, J., Janssen, W., Smiley, R. and Smith, C. (2007). High-content flow cytometry and temporal data analysis for defining a cellular signature of graft-versus-host disease. *Biology of Blood and Marrow Transplantation* **13**, 691-700.
- Chacón, J. E. (2009). Data-driven choice of the smoothing parametrization for kernel density estimators. *Canad. J. Statist.* **37**, 249-265.
- Chacón, J. E. and Duong, T. (2010). Multivariate plug-in bandwidth selection with unconstrained pilot bandwidth matrices. *Test*, **19**, 375-398
- Ćwik, J. and Koronacki, J., (1997). A combined adaptive-mixtures/plug-in estimator of multivariate probability densities. *Comput. Statist. Data Anal.* **26**, 199-218.
- Duong, T., Cowling, A., Koch, I. and Wand, M. P. (2008). Feature significance for multivariate kernel density estimation. *Comput. Statist. Data Anal.* **52**, 4225-4242.
- Godtliebsen, F., Marron, J. S. and Chaudhuri, P. (2002). Significance in scale space for bivariate density estimation. *J. Comput. Graph. Statist.* **11**, 1-21.
- Härdle, W., Marron, J. S. and Wand, M. P. (1990). Bandwidth choice for density derivatives. *J. Roy. Statist. Soc. Ser. B* **52**, 223-232.
- Henderson, H. V. and Searle, S. R. (1979). Vec and vech operators for matrices, with some uses in Jacobians and multivariate statistics. *Canad. J. Statist.* **7**, 65-81.
- Hildenbrand, K. and Hildenbrand, W. (1986). On the mean income effect: a data analysis of the U.K. family expenditure survey. In *Contributions to Mathematical Economics, in Honor of Gerard Debreu* (Edited by W. Hildenbrand and A. Mas-Colell), 247-268. Amsterdam, North-Holland.
- Holmquist, B. (1996a). The d -variate vector Hermite polynomial of order k . *Linear Algebra Appl.* **237/238**, 155-190.
- Holmquist, B. (1996b). Expectations of products of quadratic forms in normal variables. *Stochastic Anal. Appl.* **14**, 149-164.
- Jones, M. C. (1994). On kernel density derivative estimation. *Comm. Statist. Theory Methods* **23**, 2133-2139.
- Le Meur, N., Rossini, A., Gasparetto, M., Smith, C., Brinkman, R. R. and Gentleman, R. (2007). Data quality assessment of ungated flow cytometry data in high throughput experiments. *Cytometry Part A* **71A**, 393-403.
- Magnus, J.R. (1986). The exact moments of a ratio of quadratic forms in normal variables. *Ann. Econom. Statist.* **4**, 95-109.
- Magnus, J. R. and Neudecker, H. (1979). The commutation matrix: some properties and applications. *Ann. Statist.* **7**, 381-394.
- Magnus, J. R. and Neudecker, H. (1999). *Matrix Differential Calculus with Applications in Statistics and Econometrics*. Revised edition. John Wiley, Chichester.
- Naumann, U. and Wand, M. P. (2009). Automation in high-content flow cytometry screening. *Cytometry A* **75A**, 789-797
- Parzen, E. (1962). On estimation of a probability density function and mode. *Ann. Math. Statist.* **33**, 1065-1076.
- Schott, J. R. (2003). Kronecker product permutation matrices and their application to moment matrices of the normal distribution. *J. Multivariate Anal.* **87**, 177-190.

- Schuster, E. F. (1969). Estimation of a probability function and its derivatives. *Ann. Math. Statist.* **40**, 1187-1195.
- Shapiro, H. M. (2003). *Practical Flow Cytometry*. 4th edition. John Wiley, New York.
- Singh, R. S. (1976). Nonparametric estimation of mixed partial derivatives of a multivariate density. *J. Multivariate Anal.* **6**, 111-122.
- Singh, R. S. (1977). Applications of estimators of a density and its derivatives to certain statistical problems *J. Roy. Statist. Soc. Ser. B* **39**, 357-363.
- Singh, R. S. (1979). Mean squared errors of estimates of a density and its derivatives. *Biometrika* **66**, 177-180.
- Singh, R. S. (1987). MISE of kernel estimates of a density and its derivatives. *Statist. Prob. Letters* **5**, 153-159.
- Stone, C. J. (1980). Optimal rates of convergence for nonparametric estimators *Ann. Statist.* **8**, 1348-1360.
- Terrell, G. R. (1990). The maximal smoothing principle in density estimation. *J. Amer. Statist. Assoc.* **85**, 470-477.
- Wand, M. P. (1992). Error analysis for general multivariate kernel estimators. *J. Nonparametr. Stat.* **2**, 2-15.
- Wand, M. P. and Jones, M. C. (1993). Comparison of smoothing parametrizations in bivariate kernel density estimation. *J. Amer. Statist. Assoc.* **88**, 520-528.
- Wand, M. P. and Jones, M. C. (1995). *Kernel Smoothing*. Chapman & Hall, London.
- Wu, T.-J. (1997). Root n bandwidth selectors for kernel estimation of density derivatives. *J. Amer. Statist. Assoc.* **92**, 536-547.
- Wu, T.-J. and Lin, Y. (2000). Information bound for bandwidth selection in kernel estimation of density derivatives. *Statist. Sinica* **10**, 457-473.

Departamento de Matemáticas, Universidad de Extremadura, E-06006 Badajoz, Spain.

E-mail: jechacon@unex.es

Institut Pasteur, Groupe Imagerie et Modélisation; CNRS, URA 2582, F-75015 Paris, France.

Institut Curie, Molecular Mechanisms of Intracellular Transport Laboratory; CNRS, UMR 144, F-75248 Paris, France.

E-mail: tduong@curie.fr

School of Mathematics and Applied Statistics, University of Wollongong, Wollongong, Australia.

E-mail: mwand@uow.edu.au

(Received November 2008; accepted November 2009)