

**Statistica Sinica Preprint No: SS-2018-0484**

<b>Title</b>	Nonparametric density estimation for intentionally corrupted functional data
<b>Manuscript ID</b>	SS-2018-0484
<b>URL</b>	<a href="http://www.stat.sinica.edu.tw/statistica/">http://www.stat.sinica.edu.tw/statistica/</a>
<b>DOI</b>	10.5705/ss.202018.0484
<b>Complete List of Authors</b>	Aurore Delaigle and Alexander Meister
<b>Corresponding Author</b>	Alexander Meister
<b>E-mail</b>	alexander.meister@uni-rostock.de
Notice: Accepted version subject to English editing.	

---

# Nonparametric density estimation for intentionally corrupted functional data

Aurore Delaigle and Alexander Meister

*University of Melbourne, Australia and Universität Rostock, Germany*

*Abstract:* We consider statistical models where functional data are artificially contaminated by independent Wiener processes in order to satisfy privacy constraints. We show that the corrupted observations have a Wiener density which determines the distribution of the original functional random variables, masked near the origin, uniquely, and we construct a nonparametric estimator of that density. We derive an upper bound for its mean integrated squared error which has a polynomial convergence rate, and we establish an asymptotic lower bound on the minimax convergence rates which is close to the rate attained by our estimator. Our estimator requires the choice of a basis and of two smoothing parameters. We propose data-driven ways of choosing them and prove that the asymptotic quality of our estimator is not significantly affected by the empirical parameter selection. We examine the numerical performance of our method via simulated examples.

*Key words and phrases:* classification, convergence rates, differential privacy, infinite-dimensional Gaussian mixtures, Wiener densities.

## 1. Introduction

Data privacy is an important feature of a database, where the collected data are transformed and released so as to make it difficult to identify individuals participating in a study. Various privatisation methods are in use, resulting in different privacy constraints of which a popular one is differential privacy. We refer to Wasserman and Zhou (2010) for a statistical introduction of differential privacy. The privatisation mechanism typically has an impact on the statistical analysis of the data, and one of the research directions in statistical privacy is to find ways of ensuring differential privacy while keeping as much of the information as possible from the original database (see e.g. Hall et al., 2013 in the functional data context and Karwa and Slavkovi, 2016 in the setting of synthetic graphs).

One simple way of ensuring differential privacy is to contaminate the data artificially with additive random noise; see for example Wasserman and Zhou (2010). In the functional data context, Hall et al. (2013) propose a data release mechanism where the observed functional data are contaminated by adding to each function a random Gaussian process (one per functional observation) that is independent of the original data. Proposition 3.3 in Hall et al. (2013) roughly says that the data can be made differentially private whenever the scaling noise factor of the Gaussian pro-

cess is sufficiently large.

In this paper, we show that if the Gaussian process is a Wiener process and the value of the raw data is masked at the origin, then the contaminated data are differentially private, but they also have a density. This contrasts with the usual functional data setting where the assumption that all measures which are admitted to be the true image measure of functional random variables are dominated by a known basic measure seems very hard to justify. There exists no canonical basic measure such as the Lebesgue measure for finite-dimensional Euclidean data or the Haar measure for data in general locally compact groups. As a result, inference and descriptive summaries of functional data are often based on pseudo-densities. See for example Delaigle and Hall (2010) and Ciollaro et al. (2016). Recently, Lin et al. (2018) considered the estimation of densities for functions which lie in a dense subset  $S$  of the Hilbert space  $L_2(D)$ , where  $D$  is a finite interval. There,  $S$  is defined as the (non-closed) linear hull of an orthonormal basis of  $L_2(D)$  and does not contain the functional data contaminated by Wiener processes that we consider. Privacy issues for functional data are also discussed in the recent work of Mirshani et al. (2017). Therein, the authors also deduce the existence of a Gaussian density for fixed functional observations, but nonparametric estimation of that density is not studied.

By contrast, with the privatisation process we propose, the privatised functional data have a Radon-Nikodym derivative (thus a true, non pseudo, density) with respect to the Wiener measure. Exploiting the fact that the contaminating distribution is usually known in this context, we consider statistical inference from such privatised functional data.

To our knowledge, most existing nonparametric approaches for estimating a Wiener density are motivated by diffusion processes. Although these do not include the type of functional data we consider, some of these methods can be applied in our context. See for example Dabo-Niang (2004a), who suggests an orthogonal series estimator, Dabo-Niang (2002, 2004b) and Ferraty and Vieu (2006), who propose a kernel density estimator (see also Prakasa Rao, 2010a, for a generalisation in the case of diffusion processes), and Prakasa Rao (2010b) and Chesneau et al. (2013) who construct a wavelet estimator. See also Baïllo et al. (2011) for a parametric context where the data and the reference measure are Gaussian. However these methods either suffer from slow logarithmic convergence rates, or are derived under abstract assumptions that seem hard to justify in our context, or seem difficult to implement in practice. We propose a fully data-driven estimator which has fast polynomial convergence rates under simple conditions. Although our estimator is motivated by the privacy setting we

consider, our results can be extended to more general cases of functional data which have a Wiener density.

This paper is organised as follows. In Section 2 we introduce our statistical model and show that the Wiener density exists and determines uniquely the image measure of the raw functional random variables masked near zero. Moreover we prove that the privacy constraints are fulfilled when the noise level is sufficiently large. In Section 3 we construct a nonparametric orthonormal series estimator of the Wiener density and propose data-driven procedures for choosing the basis (Section 3.4) and the smoothing parameters (Section 3.5). In Section 4 we derive an explicit upper bound for the mean integrated squared error of our estimator and show that it achieves polynomial convergence rates under intuitive tail restrictions and metric entropy constraints on the measure of the original data. Functional data problems in which such fast rates are available are rare; usually the achievable rates are only logarithmic or sub-polynomial; see e.g. Dabo-Niang (2004a), Mas (2012) and Meister (2016). Finally, we derive a lower bound on the mean integrated square error under our intuitive conditions and we show that choosing the parameters in a data-driven way does not significantly deteriorate the asymptotic performance of our procedure (thus we establish a weak adaptivity result). Numerical simulations are provided in Section 5.

The proofs are deferred to supplement.

## 2. Model, data and applications

### 2.1 Model and data

We observe functional data  $Y_1, \dots, Y_n$  defined on  $[0, 1]$ , without loss of generality, and which, for reasons such as differential privacy constraints discussed in Section 1, have been intentionally contaminated by additive random noise. Specifically, we assume that

$$Y_j = X_j + \sigma W_j, \quad j = 1, \dots, n, \quad (2.1)$$

where the random functions  $X_j$  and  $W_j$ ,  $j = 1, \dots, n$ , are totally independent. Here,  $X_j$  represents the  $j$ th function of interest, which is corrupted by a standard Wiener process  $W_j$  with a deterministic scaling factor  $\sigma > 0$ .

Unlike typical measurement error problems where contamination is due to imprecise measurement or unavoidable perturbation, here the data are contaminated artificially and we can assume that  $\sigma$  is known.

We assume that the  $X_j$ 's take their values in  $C_{0,0}([0, 1])$  where  $C_{0,\ell}([0, 1])$  denotes the set of  $\ell$  times continuously differentiable (or just continuous when  $\ell = 0$ ) functions  $f$  defined on  $[0, 1]$  that are such that  $f(0) = 0$ . The  $X_j$ 's have an unknown probability measure  $P_X$  on the Borel  $\sigma$ -field

## 2.2 Density of contaminated data and differential privacy

---

$\mathfrak{B}(C_{0,0}([0, 1]))$  of  $C_{0,0}([0, 1])$  where we equip the space  $C_{0,0}([0, 1])$  with the supremum norm  $\|\cdot\|_\infty$ . Throughout we use the notation  $V_j = \sigma W_j$ , and we use  $V$ ,  $W$ ,  $X$  and  $Y$  to denote a generic function that has the same distribution as, respectively, the  $V_j$ 's,  $W_j$ 's, the  $X_j$ 's and the  $Y_j$ 's. Critically here, the functional data  $X_j$  are assumed to satisfy  $X_j(0) = 0$ . Indeed, since  $W_j(0) = 0$ , then  $Y_j(0) = X_j(0)$  and if the value of  $X_j$  at zero is not masked, then individuals can be identified from  $Y_j(0)$ . In practice, if the raw data do not satisfy  $X_j(0) = 0$ , they can be pre-masked at zero before the contamination step, for example by replacing  $X_j$  by  $\tilde{X}_j = X_j - X_j(0)$  or  $\tilde{X}_j = X_j w$  where  $w$  is a smooth function such that  $w(0) = 0$  and  $w(1) = 1$ .

## 2.2 Density of contaminated data and differential privacy

In this section, we show that the  $Y_j$ 's have a well defined density with respect to the scaled Wiener measure, and that this density characterises the distribution of the  $X_j$ 's uniquely. Finally, we show that the contamination process ensures differential privacy.

To ensure existence of a density, we need the following assumption, which we assume throughout this work:

### Assumption 1

$X \in C_{0,2}([0, 1])$  a.s..



## 2.2 Density of contaminated data and differential privacy

---

Under Assumption 1, using Girsanov's theorem (Girsanov, 1960), for any Borel measurable mapping  $\varphi$  from  $C_{0,0}([0, 1])$  to  $[0, 1]$ , we have

$$\begin{aligned} E\{\varphi(Y)\} &= E\{\varphi(X + \sigma W)\} \\ &= E\left\{\varphi(\sigma W) \exp\left(\frac{1}{\sigma} \int_0^1 X'(t) dW(t)\right) \exp\left(-\frac{1}{2\sigma^2} \int_0^1 |X'(t)|^2 dt\right)\right\} \\ &= E\left[\varphi(V) E\left\{\exp\left(\frac{1}{\sigma^2} \int_0^1 X'(t) dV(t)\right) \exp\left(-\frac{1}{2\sigma^2} \int_0^1 |X'(t)|^2 dt\right) \middle| V\right\}\right], \end{aligned}$$

so that, by integration by parts, we have, a.s.,

$$\begin{aligned} \frac{dP_Y}{dP_V}(V) &= E\left[\exp\left\{\frac{1}{\sigma^2} \int_0^1 X'(t) dV(t) - \frac{1}{2\sigma^2} \int_0^1 |X'(t)|^2 dt\right\} \middle| V\right] \\ &= \int \exp\left\{\frac{1}{\sigma^2} x'(1)V(1) - \frac{1}{\sigma^2} \int_0^1 x''(t)V(t) dt - \frac{1}{2\sigma^2} \int_0^1 |x'(t)|^2 dt\right\} dP_X(x). \end{aligned} \tag{2.2}$$

Applying the factorization lemma to this conditional expectation, we deduce that there exists a Borel measurable mapping  $f_Y: C_{0,0}([0, 1]) \rightarrow \mathbb{R}$  such that  $f_Y(V)$  is equal to the right hand side of (2.2) almost surely. This implies that  $f_Y$  is the density of  $P_Y$  with respect to  $P_V$ . Thus the contaminated  $Y_j$ 's have a density  $f_Y$ . The next theorem establishes its connection with the measure of the  $X_j$ 's.

**Theorem 1.** *The functional density  $f_Y$  in (2.2) characterises the probability measure  $P_X$  uniquely.*

## 2.2 Density of contaminated data and differential privacy

We deduce from this theorem that inference about  $P_X$  (e.g. goodness-of-fit tests or classification problems; see Section 2.3) can be performed via  $f_Y$ . To use this result in practice, it remains to see whether we can estimate  $f_Y$  nonparametrically using the data  $Y_1, \dots, Y_n$ . This is what we study in Section 3.

Throughout we use the notation  $\langle \cdot, \cdot \rangle$  for the inner product of  $L_2([0, 1])$ ,  $\|\cdot\|_2$  for the corresponding norm, and we make the following assumption:

### Assumption 2

For some constant  $C_{X,1} \in (0, \infty)$ , we have that  $\|X'\|_2 \leq C_{X,1}$  a.s..

The following proposition shows that if the scaling factor  $\sigma$  is large enough, the contaminated data are privatised. For the definition of  $(\alpha, \beta)$ -privacy we refer to Hall et al. (2013); in our setting this criterion means that

$$P[x + \sigma W \in B] \leq \exp(\alpha) \cdot P[\tilde{x} + \sigma W \in B] + \beta, \quad \forall B \in \mathfrak{B}(C_{0,0}([0, 1])),$$

for all  $x, \tilde{x} \in C_{0,2}([0, 1])$  with  $\max\{\|x'\|_2, \|\tilde{x}'\|_2\} \leq C_{X,1}$ .

**Proposition 1.** *For any  $\alpha, \beta > 0$ , choosing  $\sigma > 2C_{X,1}\sqrt{2\log(2/\beta)}/\alpha$  guarantees  $(\alpha, \beta)$ -privacy of the observation of  $Y = X + \sigma W$  under the Assumptions 1 and 2.*

## 2.3 Applications

The existence of a density for contaminated data has important practical applications. One of them is goodness-of-fit testing. Goodness-of-fit tests for functional data have been considered in e.g. Bugni et al. (2009). In our context, using the observed i.i.d. contaminated functional data  $Y_1, \dots, Y_n$ , the problem consists in testing the null hypothesis  $H_0 : X_1 \sim P_X$  versus the alternative  $H_1 : X_1 \not\sim P_X$  for some fixed probability measure  $P_X$  on  $\mathfrak{B}(C_{0,0}([0, 1]))$ . According to Theorem 1,  $H_0$  is equivalent to the claim that  $Y_1$  has the functional density  $f_Y = d(P_X * P_V)/dP_V$ . Using the estimator  $\hat{f}_Y$  of  $f_Y$  that we introduce in Section 3, a testing procedure could be based on

$$T(Y_1, \dots, Y_n) := \begin{cases} 1, & \text{for } \int |\hat{f}_Y(y) - f_Y(y)|^2 dP_V(y) > \rho, \\ 0, & \text{otherwise,} \end{cases}$$

where  $\rho$  is a threshold parameter. In Theorem 2 we derive an upper bound on the mean integrated squared error of our estimator  $\hat{f}_Y$ . Using the Markov inequality, we deduce that the test can attain any given significance level  $\alpha > 0$  if we select  $\rho$  larger or equal to the ratio of this upper bound and  $\alpha$ . While this gives some insights about  $\rho$ , this upper bound does not provide a data-driven rule for selecting  $\rho$  in practice. The latter is a difficult problem; for example, it requires deriving the asymptotic distribution of the fully

data-driven estimator. Another possibility would be to select  $\rho$  using a bootstrap approach. However, before applying such a technique, this too would require careful theoretical considerations that ensure the validity of the bootstrap in this context. While these issues are interesting, they go beyond the scope of this paper and we leave the practical choice of  $\rho$  open for future research.

Another interesting application is classification, which, in our context, can be expressed as follows. We observe training contaminated data pairs  $(Y_i, I_i)$ ,  $i = 1, \dots, n$ , where  $Y_i = X_i + V_i$ , the  $X_i$ 's come from two distinct populations  $\Pi_0$  and  $\Pi_1$  and the class label  $I_i = k$  if  $X_i$  comes from population  $\Pi_k$ , for  $k = 0, 1$ . The  $V_i$ 's are Wiener processes independent of the  $X_i$ 's, and are identically distributed within each population, but the scaling noise parameter  $\sigma$  need not be the same for the two populations. Using these data, the goal is to classify in  $\Pi_0$  or  $\Pi_1$  a new random curve  $Y = X + V$ , where  $X$  comes from either  $\Pi_0$  or  $\Pi_1$ , but whose class label is unknown.

It is well known in general classification problems that the optimal classifier is the Bayes classifier which, adapted to our context, assigns a curve to  $\Pi_1$  if  $E(I|Y = y) > 1/2$  and to  $\Pi_0$  otherwise. In the case where the probability measures  $P_{Y,0}$  and  $P_{Y,1}$  of the  $Y_i$ 's that originate from, respectively,

$\Pi_0$  and  $\Pi_1$ , have well defined densities  $f_{Y,0}$  and  $f_{Y,1}$ , the Bayes classifier can be expressed as: assign  $Y$  to  $\Pi_1$  if  $\pi_1 f_{Y,1}(Y) > \pi_0 f_{Y,0}(Y)$ , and to  $\Pi_0$  otherwise, where  $\pi_k = P(I = k)$ . In the particular Gaussian case, Baïllo et al. (2011) showed that these densities are well defined and showed how to estimate them.

In our case the  $Y_i$ 's are generally not Gaussian but they have functional densities  $f_{Y,k} = dP_{Y,k}/dP_V$ , for  $k = 0, 1$ . Since  $P_{X,0} \neq P_{X,1}$  implies that  $f_{Y,0} \neq f_{Y,1}$  (see Theorem 1), these densities can be used for classification of  $X$  from observations on  $Y$  in the optimal Bayes classifier. There, in practice we classify  $Y$  in  $\Pi_1$  if  $\pi_1 \hat{f}_{Y,1}(Y) \geq \pi_0 \hat{f}_{Y,0}(Y)$  and in  $\Pi_0$  otherwise, where for  $k = 0, 1$ ,  $\hat{f}_{Y,k}$  denotes the estimator of  $f_{Y,k}$  from Section 3 constructed from the training data  $Y_i$  for which  $I_i = k$ .

There exist many other classification procedures for functional data, often based on pseudo-densities or finite dimensional approximations. However, Delaigle and Hall (2012) pointed that, except in the Gaussian case, often such projections do not ensure good finite sample performance. See for example Hall et al. (2001), Ferraty and Vieu (2006), Escabias et al. (2007), Preda et al. (2007) and Shin (2008). See also Dai et al. (2017) for a recent example, where the authors approximate the densities in the two populations by the finite dimensional surrogate densities proposed in Delaigle and

Hall (2010); see Delaigle and Hall (2013) for a related classifier.

### 3. Methodology

In this section we consider the problem of estimating the functional density  $f_Y$  nonparametrically.

#### 3.1 Existing methods

Nonparametric estimation of a density for stochastic processes whose probability measure has a Radon-Nikodym derivative with respect to the Wiener measure has been considered by several authors. In Dabo-Niang (2002, 2004b), the author proposes to use a kernel density estimator; see also Prakasa Rao (2010a). This estimator is simple but it suffers from slow logarithmic convergence rates, which are reflected by its practical performance. A wavelet estimator with polynomial convergence rates was proposed by Prakasa Rao (2010b) and Chesneau et al. (2013), but their conditions are quite technical and it is not clear how their parameters can be chosen in practice. Moreover, their theory is derived under abstract high level conditions which might not be easily satisfied in our context.

A simpler estimator is the orthogonal series estimator of Dabo-Niang (2004a), defined as follows. Let  $\{\varphi_j\}_{j \in \mathbb{N}}$  denote an orthonormal basis of real-valued

### 3.1 Existing methods14

functions of  $[0, 1]$ , where each  $\varphi_j \in L_2([0, 1])$ , and let  $(H_j)_{j \geq 1}$  denote the scaled Hermite polynomials defined by  $H_k(x) = (-1)^k \phi^{(k)}(x) / \{\phi(x) \sqrt{k!}\}$ , for all integer  $k \geq 0$ , where  $\phi(x) = \exp(-x^2/2) / \sqrt{2\pi}$ . Also, for  $x \in C_0([0, 1])$ , let

$$\beta'_{x,\ell} = \int_0^1 \varphi_\ell(t) dx(t). \quad (3.3)$$

Using results from Cameron and Martin (1947), the author notes that, as  $K \rightarrow \infty$ , the Fourier-Hermite series  $(\Psi_{k_1, \dots, k_K})_{0 \leq k_1 \leq K, \dots, 0 \leq k_K \leq K}$ , where, for  $x \in C_0([0, 1])$ ,

$$\Psi_{k_1, \dots, k_K}(x) \equiv H_{k_1, \dots, k_K}(\beta'_{x,1}, \dots, \beta'_{x,K}) \equiv \prod_{\ell=1}^K H_{k_\ell}(\beta'_{x,\ell}), \quad (3.4)$$

forms an orthonormal basis of the Hilbert space of all square-integrable  $C_0([0, 1])$ -valued random variables with respect to the Wiener measure. Motivated by this, the author proposes to estimate the Wiener density  $f_T$  of functional data  $T_1, \dots, T_n$  (that have a Wiener density) by

$$\hat{f}_T^K(x) = \sum_{k_1, \dots, k_K=0}^K \frac{1}{n} \sum_{j=1}^n H_{k_1, \dots, k_K}(\beta'_{T_j,1}, \dots, \beta'_{T_j,K}) \cdot H_{k_1, \dots, k_K}(\beta'_{x,1}, \dots, \beta'_{x,K}), \quad (3.5)$$

where  $K$  is a smoothing parameter. This estimator is very attractive for its simplicity. However, a drawback is that the rates derived by Daboniang (2004a) are logarithmic. In the next two sections, using a two-stage approximation approach (first a sieve approximation of  $f_Y$  and then an es-

### 3.2 Finite-dimensional approximation of $f_Y$

timator of the approximation), we are able to introduce a different regularisation scheme which involves two parameters. This increases the flexibility of the estimator, which, as we shall see, enables us to obtain polynomial convergence rates. Moreover we provide data-driven choices of the basis and the threshold parameters.

#### 3.2 Finite-dimensional approximation of $f_Y$

Recall from (2.2) that for  $V = \sigma W$  with  $W$  a standard Wiener process, we have

$$f_Y(V) = E \left[ \exp \left\{ \frac{1}{\sigma^2} \int_0^1 X'(t) dV(t) - \frac{1}{2\sigma^2} \int_0^1 |X'(t)|^2 dt \right\} \middle| V \right], \quad \text{a.s.},$$

and that our goal is to estimate  $f_Y$  from data  $Y_1, \dots, Y_n$ . Instead of directly expressing  $f_Y$  in the Fourier-Hermite basis at (3.4), we first construct a sieve approximation of  $f_Y$ , and then (see Section 3.3) we express our sieve approximation in the Fourier-Hermite basis.

Using the notation  $\beta'_{x,\ell} = \int_0^1 \varphi_\ell(t) dx(t)$  from Equation (3.3), where  $\{\varphi_j\}_{j \in \mathbb{N}}$  is a real-valued orthonormal basis of  $L_2([0, 1])$ , we can write

$$\int_0^1 X'(t) dV(t) - \frac{1}{2} \int_0^1 |X'(t)|^2 dt = \sum_{j=1}^{\infty} \beta'_{X,j} \cdot \beta'_{V,j} - \frac{1}{2} \sum_{j=1}^{\infty} \beta'^2_{X,j}, \quad (3.6)$$

where the infinite sums should be understood as mean squared limits. Truncating the sums to  $m$  terms, with  $m \geq 1$  an integer, this suggests that we can



### 3.2 Finite-dimensional approximation of $f_Y$ 16

approximate  $f_Y(V)$  by  $f_Y^{[m]}(\beta'_{V,1}, \dots, \beta'_{V,m})$ , where, for all  $s_1, \dots, s_m \in \mathbb{R}$ ,

$$\begin{aligned} f_Y^{[m]}(s_1, \dots, s_m) &= E \left\{ \exp \left( \frac{1}{\sigma^2} \sum_{j=1}^m \beta'_{X,j} \cdot s_j - \frac{1}{2\sigma^2} \sum_{j=1}^m \beta'_{X,j}{}^2 \right) \right\} \\ &= \exp \left( \frac{1}{2\sigma^2} \sum_{j=1}^m s_j^2 \right) \int \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^m (s_j - x_j)^2 \right\} dP_{X,m}(x_1, \dots, x_m), \end{aligned} \quad (3.7)$$

and where  $P_{X,m}$  denotes the measure of  $(\beta'_{X,1}, \dots, \beta'_{X,m})$ .

The following lemma shows that, as long as  $m$  is sufficiently large,  $f_Y^{[m]}(\beta'_{V,1}, \dots, \beta'_{V,m})$  is a good approximation to  $f_Y(V)$ , where  $V$  denotes a generic  $V_i \sim P_V$ .

**Lemma 1.** *Let  $\mathfrak{A}_m$  denote the  $\sigma$ -field generated by  $\beta'_{V_1,1}, \dots, \beta'_{V_1,m}$ . Under Assumptions 1 and 2,*

(a)  $f_Y^{[m]}(\beta'_{V_1,1}, \dots, \beta'_{V_1,m}) = E\{f_Y(V_1)|\mathfrak{A}_m\}$  a.s.,

(b) we have

$$\begin{aligned} E|f_Y^{[m]}(\beta'_{V_1,1}, \dots, \beta'_{V_1,m}) - f_Y(V_1)|^2 \\ \leq \frac{1}{\sigma^2} \cdot \exp(C_{X,1}^2/\sigma^2) \cdot \left( \sum_{j,j'>m} |\langle \varphi_j, \Gamma_X \varphi_{j'} \rangle|^2 \right)^{1/2}, \end{aligned}$$

where the linear operator  $\Gamma_X : L_2([0, 1]) \rightarrow L_2([0, 1])$  is defined by

$$(\Gamma_X f)(t) = E \left\{ X'(t) \int_0^1 X'(s) f(s) ds \right\}, \quad t \in [0, 1], f \in L_2([0, 1]). \quad (3.8)$$

### 3.3 Estimating the sieve approximation of $f_Y$

Since  $\Gamma_X$  is a self-adjoint and positive-semidefinite Hilbert-Schmidt operator, the upper bound in Lemma 1(b) is finite for any orthonormal basis  $\{\varphi_j\}_j$  of  $L_2([0, 1])$ , and converges to zero as  $m \rightarrow \infty$ . Indeed, Assumption 2 guarantees that  $\sum_{j,j'} |\langle \varphi_j, \Gamma_X \varphi_{j'} \rangle|^2 \leq E \|X'_1\|_2^4 \leq C_{X,1}^4 < \infty$ . If  $X$  (and hence  $X'$ ) is centered then  $\Gamma_X$  coincides with the covariance operator of  $X'$ .

### 3.3 Estimating the sieve approximation of $f_Y$

Next we show how to estimate  $f_Y^{[m]}$  using a Fourier-Hermite series. For this, let  $P_{Y,m}$  and  $f_{Y,m}$  denote, respectively, the measure and the  $m$ -dimensional Lebesgue density of the observed random vector  $(\beta'_{Y_j,1}, \dots, \beta'_{Y_j,m})$ , where

$$\beta'_{Y_j,k} = \int_0^1 \varphi_k(t) dY_j(t) = \beta'_{X_j,k} + \beta'_{V_j,k}, \quad j = 1, \dots, n; k = 1, \dots, m.$$

Let  $g_\sigma$  denote the  $N(0, \sigma^2 I_m)$ -density, with  $I_m$  the  $m \times m$ -identity matrix, let  $L_{2,g_\sigma}(\mathbb{R}^m)$  denote the Hilbert space of Borel measurable functions  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  which satisfy  $\|f\|_{g_\sigma}^2 \equiv \int |f(t)|^2 g_\sigma(t) dt < \infty$ , and let  $\langle \cdot, \cdot \rangle_{g_\sigma}$  denote the inner product of  $L_{2,g_\sigma}(\mathbb{R}^m)$ .

It is easy to deduce from (3.7) that

$$f_Y^{[m]}(s_1, \dots, s_m) = f_{Y,m}(s_1, \dots, s_m) / g_\sigma(s_1, \dots, s_m), \quad (3.9)$$

and it can be proved that  $f_Y^{[m]} \in L_{2,g_\sigma}(\mathbb{R}^m)$ . Therefore, if  $\Psi_1, \Psi_2, \dots$  is an

### 3.3 Estimating the sieve approximation of $f_Y$ 18

orthonormal basis of  $L_{2,g_\sigma}(\mathbb{R}^m)$ , we can write

$$f_Y^{[m]} = \sum_{k=1}^{\infty} \alpha_k \Psi_k,$$

$$\alpha_k = \langle \Psi_k, f_Y^{[m]} \rangle_{g_\sigma} = \int \Psi_k(y) f_{Y,m}(y) dy = E\{\Psi_k(\beta'_{Y,1}, \dots, \beta'_{Y,m})\}.$$

Now the sequence  $(H_{k_1, \dots, k_m})_{k_1, \dots, k_m \geq 0}$  of functions  $H_{k_1, \dots, k_m}(x_1, \dots, x_m) = \prod_{j=1}^m H_{k_j}(x_j)$  defined at (3.4) forms an orthonormal basis of  $L_{2,g_1}(\mathbb{R}^m)$ .

Thus we can take  $\Psi_k(\cdot) = H_{k_1, \dots, k_m}(\cdot/\sigma)$ . To estimate  $f_Y^{[m]}$ , we replace  $\alpha_k$  by  $\hat{\alpha}_k = n^{-1} \sum_{j=1}^n \Psi_k(\beta'_{Y_j,1}, \dots, \beta'_{Y_j,m})$ .

Finally, for  $U$  a functional random variable independent of  $Y_1, \dots, Y_n$  which has a density with respect to  $P_V$ , we define our estimator of  $f_Y(U)$  by

$$\begin{aligned} \hat{f}_Y^{[m,K]}(U) &= \sum_{k_1, \dots, k_m \geq 0} \frac{1}{n} \sum_{j=1}^n H_{k_1, \dots, k_m}(\beta'_{Y_j,1}/\sigma, \dots, \beta'_{Y_j,m}/\sigma) H_{k_1, \dots, k_m}(\beta'_{U,1}/\sigma, \dots, \beta'_{U,m}/\sigma) \\ &\quad \times \omega_K(k_1 + \dots + k_m) 1\{k_1 + \dots + k_m \leq K\}, \end{aligned} \tag{3.10}$$

where  $K \geq 0$  is a truncation parameter and  $0 \leq \omega_K(x) \leq 1$  a continuous function defined on  $[0, K]$ . The term  $\omega_K(k_1 + \dots + k_m) 1\{k_1 + \dots + k_m \leq K\}$  prevents the  $k_i$ 's from being too large, which controls the variability of the estimator. Using wavelet terminology, the function  $\omega_K$  dictates whether the

### 3.3 Estimating the sieve approximation of $f_{Y|X}$

---

$k_i$ 's are chosen by a soft or a hard rule. Specifically, a hard rule corresponds to  $\omega_K \equiv 1$ : here all  $k_i$ 's summing to at most  $K$  are given equal weight and as  $K$  increases, new indices appear and play as big a role as older ones. For a soft rule,  $\omega_K(x)$  is taken to be a smooth decreasing function of  $x$ , e.g.  $\omega_K(x) = 1 - x/(K + 1)$ ; as  $K$  increases, new indices start playing a role but have less weight than former ones.

A major difference between (3.10) and Dabo-Niang's (2004a) estimator at (3.5) is our regularisation scheme: because of the two-step construction of our estimator (sieve approximation followed by basis expansion), we do not use all the indices  $(k_1, \dots, k_K) \in \{0, \dots, K\}^K$ . Instead we use  $(k_1, \dots, k_m) \in \{0, \dots, K\}^m$  such that  $k_1 + \dots + k_m \leq K$ , and we assign a weight  $\omega_K(k_1 + \dots + k_m)$  to each group of  $m$  indices. As we will see in the next sections, our use of a second parameter  $m$  and the restriction we put on  $k_1 + \dots + k_m$  drastically improve the quality of the estimator, both theoretically and practically. Moreover, in Section 3.4, we introduce a data-driven way of choosing the basis  $\{\varphi_j\}_{j \in \mathbb{N}}$  used to construct the coefficients  $\beta'_{Y_j, k}$  and  $\beta'_{U, k}$ .

### 3.4 Choosing the $\varphi_j$ 's

To compute our estimator in practice, we need to choose the basis  $\{\varphi_j\}_j$  used in (3.3). Lemma 1(b) implies that if we take the  $\varphi_j$ 's equal to the eigenfunctions of  $\Gamma_X$ , ordered such that the sequence of corresponding eigenvalues  $(\lambda_j)_j$  decreases monotonically, then

$$E|f_Y^{[m]}(\beta'_{V_{1,1}}, \dots, \beta'_{V_{1,m}}) - f_Y(V_1)|^2 \leq \frac{1}{\sigma^2} \cdot \exp(C_{X,1}^2/\sigma^2) \cdot \left(\sum_{j>m} \lambda_j^2\right)^{1/2}. \quad (3.11)$$

This bound decreases monotonically as  $m$  increases, which gives an indication that the first  $m$  terms of the basis capture some of the main characteristics of  $f_Y$ .

Of course, in practice  $\Gamma_X$  is unknown and thus the  $\varphi_j$ 's are unknown. Thus we need to estimate  $\Gamma_X$ , but a priori this does not seem to be an easy task because, up to some mean terms,  $\Gamma_X$  is the covariance function of the first derivative  $X'$  of  $X$ . If we could observe  $X'_1, \dots, X'_n$ , we could use standard covariance estimation techniques such as those in Hall and Hosseini-Nasab (2006), Mas and Ruymgaard (2015) and Jirak (2016). However we only observe the contaminated  $Y_j$ 's. If the  $Y_j$ 's were differentiable, we could take their derivative and estimate  $\Gamma_X$  and its eigenfunctions as in the references just cited. However they are not differentiable and we cannot take such a simple approach.

### 3.4 Choosing the $\varphi_j$ 's<sup>21</sup>

Instead, we propose the following approximation procedure. Let  $\{\psi_j\}_j$  denote an orthonormal basis of  $L_2([0, 1])$ , and recall that  $\varphi_\ell$  denotes the eigenfunction of  $\Gamma_X$  with eigenvalue  $\lambda_\ell$ , where  $\lambda_1 \geq \lambda_2 \geq \dots$ . In the supplement we show that, for all  $k \geq 1$ ,

$$\sum_{j=1}^{\infty} \varphi_{\ell,j} \langle \psi_k, \Gamma_X \psi_j \rangle = \lambda_\ell \varphi_{\ell,k}, \quad (3.12)$$

where  $\varphi_{\ell,j} = \langle \varphi_\ell, \psi_j \rangle$ . If we take the  $\psi_j$ 's to be continuously differentiable and such that  $\psi_j(0) = \psi_j(1) = 0$ , for example if  $\{\psi_j\}_j$  is the Fourier sine basis, then for  $j, k = 1, 2, \dots$ , we have

$$\langle \psi_k, \Gamma_X \psi_j \rangle = \mathcal{M}_{j,k} - \sigma^2 \cdot 1\{j = k\}, \quad (3.13)$$

where  $\mathcal{M}_{j,k} = \int_0^1 \psi_j'(t) \int_0^1 E\{Y(t)Y(s)\} \psi_k'(s) ds dt$  (see the proof in the supplement). We propose to approximate  $\varphi_\ell$  by  $\sum_{j=1}^M \hat{\varphi}_{\ell,j} \psi_j$ , with  $M$  a large positive integer, where  $\hat{\varphi}_{\ell,j}$  denotes an estimator of  $\varphi_{\ell,j}$ . Next we show how to compute  $\hat{\varphi}_{\ell,1}, \dots, \hat{\varphi}_{\ell,M}$  from our data. First, combining (3.12) and (3.13), we have  $\sum_{j=1}^{\infty} \varphi_{\ell,j} (\mathcal{M}_{j,k} - \sigma^2 \cdot 1\{j = k\}) = \lambda_\ell \varphi_{\ell,k}$  so that

$$\sum_{j=1}^M \varphi_{\ell,j} (\mathcal{M}_{j,k} - \sigma^2 \cdot 1\{j = k\}) = \lambda_\ell \varphi_{\ell,k} + R_{k,\ell}, \quad (3.14)$$

where  $R_{k,\ell}$  is a remainder term resulting from the truncation of the sum to  $M$  terms. Let  $I_M$  and  $\mathcal{M}$  denote, respectively, the  $M \times M$ -identity matrix and the  $M \times M$ -matrix whose components are defined by  $\mathcal{M}_{j,k}$ ,

### 3.5 Choosing the parameters $M$ , $m$ and $K$

$j, k = 1, \dots, M$ , and let  $\Phi_\ell = (\varphi_{\ell,1}, \dots, \varphi_{\ell,M})^T$  and  $R_\ell = (R_{1,\ell}, \dots, R_{M,\ell})^T$ .

Then (3.14) implies that  $(\mathcal{M} - \sigma^2 I_M) \Phi_\ell = \lambda \Phi_\ell + R_\ell$ .

Note that  $|R_\ell|$  shrinks to zero as  $M \rightarrow \infty$  since  $|R_\ell|^2 \leq C_{X,1}^4 \sum_{j>M} |\varphi_{\ell,j}|^2$ .

Thus,  $(\mathcal{M} - \sigma^2 I_M) \Phi_\ell \approx \lambda_\ell \Phi_\ell$ , which motivates us to approximate  $\Phi_\ell$  by the

unit eigenvector  $v_\ell$  of the matrix  $\mathcal{M} - \sigma^2 I_M$  corresponding to the  $\ell$ th largest

eigenvalue. Now,  $(\mathcal{M} - \sigma^2 I_M) v_\ell = \lambda_\ell v_\ell$  implies that  $\mathcal{M} v_\ell = (\lambda_\ell + \sigma^2) v_\ell$ .

Thus,  $v_\ell$  is also the eigenvector of  $\mathcal{M}$  corresponding to its  $\ell$ th largest eigen-

value. Of course,  $\mathcal{M}$  is unknown but it can be estimated by

$$\hat{\mathcal{M}} = \frac{1}{n} \sum_{\ell=1}^n \left\{ \int_0^1 \int_0^1 \psi'_j(t) Y_\ell(t) Y_\ell(s) \psi'_k(s) ds dt \right\}_{j,k=1,\dots,M}. \quad (3.15)$$

For  $\ell = 1, \dots, M$ , let  $\hat{v}_\ell$  denote the  $M$  unit eigenvectors of  $\hat{\mathcal{M}}$  (ordered

so that the corresponding eigenvalues decrease monotonically). We propose

to estimate  $\Phi_\ell$  by  $\hat{\Phi}_\ell = (\hat{\varphi}_{\ell,1}, \dots, \hat{\varphi}_{\ell,M})^T = \hat{v}_\ell$ . Finally, we estimate  $\varphi_\ell$  by

$$\hat{\varphi}_\ell = \sum_{j=1}^M \hat{\varphi}_{\ell,j} \psi_j.$$

### 3.5 Choosing the parameters $M$ , $m$ and $K$

To compute the estimator at (3.10) in practice, we need to choose three

parameters:  $M$ , the parameter used in Section 3.4 to construct the basis

functions  $\varphi_j$  employed to compute the projections in (3.3),  $m$ , a parameter

which dictates the dimension of our approximation of  $f_Y$  by  $f_Y^{[m]}$  at (3.7),

and  $K$ , the truncation parameter of our orthogonal series expansion at

### 3.5 Choosing the parameters $M$ , $m$ and $K$

(3.10). Having the  $\hat{\varphi}_j$ 's close to the eigenfunctions of  $\Gamma_X$  is likely to give better practical performance, but it is not necessary for the consistency of our estimator. This suggests that the choice of  $M$  is not crucial and we take  $M = 20$ . By contrast,  $m$  and  $K$  are important smoothing parameters which influence consistency and need to be chosen with care. We suggest choosing  $(m, K)$  by minimising the cross-validation (CV) criterion

$$CV(m, K) = \int |\hat{f}_Y(v)|^2 dP_V(v) - \frac{2}{n} \sum_{i=1}^n \hat{f}_Y^{(-i)}(Y_i), \quad (3.16)$$

with  $\hat{f}_Y^{(-i)}$  defined in the same way as the estimator at (3.10), except that it uses only the data  $Y_1, \dots, Y_{i-1}, Y_{i+1}, \dots, Y_n$ . To compute the integral at (3.16) we generate a large sample (we took a sample of size 10000 in our numerical work) of  $V_j$ 's from  $P_V$  and approximate the integral by the mean of the  $|\hat{f}_Y(V_j)|^2$ 's.

As in standard nonparametric density estimation problems, our cross-validation criterion can have multiple local minima and the global minimum is not necessarily a good choice. In case of multiple local minima, we choose the one that produces the smallest value of  $m + K$ . Moreover, when minimising  $CV(K, m)$  we discard all pairs of values  $(K, m)$  for which more than 50% of the  $\hat{f}_Y^{(-i)}$ 's or of the  $\hat{f}_Y$ 's are negative. For the non discarded  $(K, m)$ 's, we replace each negative  $\hat{f}_Y^{(-i)}(Y_i)$  and  $\hat{f}_Y(V_j)$  by recomputing those estimators by repeatedly replacing  $K$  by  $K - 1$  and  $m$  and  $m - 1$



until the negative estimators become positive.

#### 4. Theoretical properties

In this section we derive theoretical properties of our estimator. For simplicity we derive our results in the case where the weight function  $\omega_K$  in (3.10) is equal to 1. Similar results can be established for a more general weight function, but at the expense of even more technical proofs. In Section 4.1 we derive an upper bound on the mean integrated squared error of our estimator which is valid for all  $n$ . Next, in Section 4.2 we derive asymptotic properties of our estimator.

##### 4.1 Finite sample properties

In the next theorem, we give an upper bound on the mean integrated squared error

$$\mathcal{R}(\hat{f}_Y^{[m,K]}, f_Y) = E \int |\hat{f}_Y^{[m,K]}(v) - f_Y(v)|^2 dP_V(v)$$

of the estimator at (3.10) in the case where the orthonormal basis  $\{\varphi_j\}_j$  and the parameters  $m$  and  $K$  are deterministic. Our result is non asymptotic and is valid for all  $n$ .

**Theorem 2.** *Under Assumptions 1 and 2 and the selection  $\omega_K \equiv 1$ , we*

have  $\mathcal{R}(\hat{f}_Y^{[m,K]}, f_Y) \leq \mathcal{V} + \mathcal{B} + \mathcal{D}$ , where

$$\mathcal{V} = \frac{1}{n} \exp(KC_{X,1}^2/\sigma^2) \cdot \binom{K+m}{K}, \quad \mathcal{B} = \inf_{h \in \mathcal{H}_{m,K}} \|f_Y^{[m]}(\sigma \cdot) - h\|_{g_1}^2,$$

$$\mathcal{D} = \frac{1}{\sigma^2} \cdot \exp(C_{X,1}^2/\sigma^2) \cdot \left( \sum_{j,j' > m} |\langle \varphi_j, \Gamma_X \varphi_{j'} \rangle|^2 \right)^{1/2},$$

and where  $\mathcal{H}_{m,K}$  denotes the linear hull of the  $H_{k_1, \dots, k_m}$ 's for which  $k_1 + \dots + k_m \leq K$ .

In Theorem 2,  $\mathcal{V}$  represents a variance term while  $\mathcal{B}$  represents a bias term which depends on smoothness properties of  $f_Y^{[m]}$ . Both are typical of nonparametric estimators, but the term  $\mathcal{D}$  is of a different type. It reflects the error of the finite-dimensional approximation of the density  $f_Y$  by the function  $f_Y^{[m]}$ .

## 4.2 Asymptotic properties

Next we derive asymptotic properties of our density estimator. For this, we need an additional assumption which will be used when dealing with the term  $\mathcal{D}$  from Theorem 2:

### Assumption 3

There exist constants  $C_{X,2}, C_{X,3} \in (0, \infty)$  and  $\gamma > 0$  such that

$$\sum_{j,j' > m} \left| \int_0^1 \varphi_j(s) (\Gamma_X \varphi_{j'})(s) ds \right|^2 \leq C_{X,2} \cdot \exp(-C_{X,3} m^\gamma), \quad \forall m \in \mathbb{N}.$$

For example, if  $X_1$  is centered and  $\{\varphi_j\}_j$  is the principal component basis with eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots$  discussed in Section 3.4, then Assumption 3 is satisfied as soon as  $\sum_{j=1}^{\infty} \exp(C'_{X,3} j^\gamma) \cdot \lambda_j^2 < \infty$  for some  $C'_{X,3} > C_{X,3}$ . In this case, Assumption 3 can be interpreted as an exponential decay of the eigenvalues of  $\Gamma_X$ ; concretely Assumption 3 is satisfied if there exist some  $C''_{X,3} > C'_{X,3} > C_{X,3}$  and some  $C'''_{X,3} > 0$  such that  $\lambda_j \leq C'''_{X,3} \exp(-C''_{X,3} j^\gamma / 2)$  for all integer  $j \geq 1$ .

The next theorem establishes an upper bound on the convergence rates of the mean integrated squared error of our estimator  $\hat{f}_Y^{[m,K]}$  as the sample size  $n$  tends to infinity. We establish the upper bound uniformly over the class  $\mathcal{F}_X = \mathcal{F}_X(C_{X,1}, C_{X,2}, C_{X,3}, \gamma, \{\varphi_j\}_j)$  of all admitted image measures of  $X_1$  such that Assumptions 1 to 3 are satisfied for some deterministic orthonormal basis  $\{\varphi_j\}_j$  of  $L_2([0, 1])$ . The next three theorems consider functions in this class, which implies that they are derived under Assumptions 1 to 3.

**Theorem 3.** *Assume that  $\gamma \in (0, 1)$  and select the weight function  $\omega_K \equiv 1$  and the parameters  $K$  and  $m$  such that  $K = K_n = \lfloor \gamma(\log n) / \log(\log n) \rfloor$ ,  $m = m_n = \lfloor (C_M \cdot \log n)^{1/\gamma} \rfloor$ , for some finite constant  $C_M > 2/C_{X,3}$ . Then*

our estimator  $\hat{f}^{[m,K]}$  satisfies

$$\limsup_{n \rightarrow \infty} \sup_{P_X \in \mathcal{F}_X} \log \{ \mathcal{R}(\hat{f}_Y^{[m,K]}, f_Y) \} / \log n \leq -\gamma.$$

Theorem 3 shows that the risk of our estimator converges to zero faster than  $\mathcal{O}(n^{-\gamma'})$  for any  $\gamma' < \gamma < 1$ . In particular, our estimator achieves polynomial convergence rates, which is usually impossible in problems of nonparametric functional regression or density estimation. In standard problems of that type where the data range over an infinite-dimensional space, only logarithmic or sub-algebraic rates can usually be achieved (see e.g. Mas, 2012, Chagny and Roche, 2014 and Meister, 2016). In our case the dimension of the data is infinite as well; however the density  $f_Y$  forms an infinite-dimensional Gaussian mixture and its smoothness degree is sufficiently high to overcome the difficulty caused by high dimensionality.

The next theorem provides an asymptotic lower bound for the problem of estimating  $f_Y$  nonparametrically. For simplicity we restrict to the case where  $C_{X,1} = 1$ .

**Theorem 4.** *Assume that  $\gamma \in (0, 1)$  and let  $C_{X,1} = 1$  in Assumption 2. Moreover, assume that the orthonormal basis  $\{\varphi_j\}_j$  of  $L_2([0, 1])$  is such that all  $\varphi_j$ 's are continuously differentiable. Then, for any sequence  $(\hat{f}_n)_n$*

of estimators of  $f_Y$  computed from the data  $Y_1, \dots, Y_n$ , we have

$$\liminf_{n \rightarrow \infty} \sup_{P_X \in \mathcal{F}_X} \log \{ \mathcal{R}(\hat{f}_n, f_Y) \} / \log n \geq -\gamma + (\gamma - 1)^2 / (\gamma - 2).$$

We learn from the theorem that, in this problem, no nonparametric estimator can reach the parametric squared convergence rate  $n^{-1}$ . This is significantly different from the simpler problem of nonparametric estimation of one-dimensional Gaussian mixtures, where the parametric rates are achievable up to a logarithmic factor (see Kim, 2014). Note that the upper bound in Theorem 3 is usually larger than the lower bound in Theorem 4, although the two bounds are very close to each other for  $\gamma$  close to 1. Rather than our estimator being suboptimal, we suspect that our lower bound is not sharp enough. Deriving the exact minimax rates seems a very challenging open problem for future research.

As is standard in nonparametric estimation problems requiring the choice of smoothing parameters, Theorem 3 was derived under a deterministic choice of  $m$  and  $K$ . Next we establish an asymptotic result in the case where  $(\hat{m}, \hat{K})$  is chosen by cross-validation as at (3.16), where minimisation is performed over the mesh

$$G = \{ \lfloor \log n \rfloor, \dots, \lfloor (\log n)^{1/\gamma_0} \rfloor \} \times \{ 1, \dots, \lfloor (\log n) / \log(\log n) \rfloor \}, \quad (4.17)$$

for some constant  $\gamma_0 \in (0, \gamma)$ . The following theorem shows that the con-

vergence rates from Theorem 3 can be maintained at least in a weak sense.

**Theorem 5.** *Our estimator  $\hat{f}_Y^{[\hat{m}, \hat{K}]}$ , where  $\omega_K \equiv 1$  and  $(\hat{m}, \hat{K})$  is selected by cross-validation over the mesh  $G$  at (4.17), satisfies*

$$\lim_{n \rightarrow \infty} \sup_{P_X \in \mathcal{F}_X} P \left\{ n^\gamma \int | \hat{f}_Y^{[\hat{m}, \hat{K}]}(x) - f_Y(x) |^2 dP_V(x) \geq n^d \right\} = 0,$$

for all  $\gamma \in [\gamma_0, 1)$  and  $d > 0$ .

## 5. Simulation results

To illustrate the performance of our density estimation procedure, we performed simulations in different settings. For a grid of  $T = 101$  points  $0 = t_0 < t_1 < \dots < t_T = 1$  equispaced by  $\Delta t = 1/(T - 1)$ , we generated data  $Y_i(t_k) = \sum_{j=1}^J \sqrt{\lambda_j} Z_{ik} \phi_j(t_k) + \sigma W_i(t_k)$ , where the  $Z_{ik}$ 's are i.i.d., each  $Z_{ik}$  is the average of the two independent  $U[-.1, .1]$  random variables,  $W_i(t_0) = 0$  and, for  $k = 1, \dots, T$ ,  $W_i(t_k) = W_i(t_{k-1}) + \epsilon_{ik}$ , where the  $\epsilon_{ik}$ 's are i.i.d.  $\sim N(0, \Delta t)$ . We considered five settings: (i)  $J = 20$ ,  $\sigma = 0.1$ ,  $\lambda_j = \exp(-j)$  and  $\phi_j(t) = \sqrt{2} \sin(\pi t j)$ ; (ii) same as (i) but with  $J = 40$ ; (iii) same as (ii) but with  $\sigma = 0.075$ ; (iv) same as (i) but with  $\sigma = 0.075$ ,  $\phi_j(t) = \sqrt{2} \cos(\pi t j) \kappa(t)$ ,  $\kappa(t) = 2 \exp(10t) / \{1 + \exp(10t)\} - 1$ ; (v) same as (i) but with  $\sigma = 0.075$ ,  $\phi_j(t) = \sqrt{2} \sin(\pi t j) \kappa(t)$ .

In each case we generated  $B = 200$  samples of  $Y_i(t_k)$ 's, of sizes  $n = 500$ ,

1000, 2000 and 5000. Then, for  $b = 1, \dots, B$ , using the  $b$ th sample of  $Y_i(t_k)$ 's, we computed our density estimator  $\hat{f}_Y^{[m,K]}(V)$  at (3.10) for  $10^4$  functions  $V$  generated from the same distribution as  $\sigma W$ , where  $m$  and  $K$  were chosen by cross-validation by minimisation of (3.16) and where we took the weight function  $\omega_K(x) = 1 - x/(K + 1)$ . The basis functions  $\varphi_j$  were computed as in Section 3.4 with  $M = 20$  and  $\psi_j(t) = \sqrt{2} \sin(\pi t j)$ ; we denote by DM the resulting estimator. Each time the  $m$  and  $K$  selected by CV produced a negative estimator  $\hat{f}_Y(v)$  for a new data curve  $v$ , for that curve  $v$  we repeatedly replaced,  $K$  by  $K - 1$  and  $m$  and  $m - 1$  until the resulting value of  $(m, K)$  was such that  $\hat{f}_Y(v) > 0$ .

In each case we also computed the estimator of Dabo-Niang (2004a) with our adaptive basis of  $\varphi_j$ 's, which we denote by DN. We chose  $K$  by minimisation of the cross-validation criterion at (3.16), replacing there our estimator by this estimator and  $(m, K)$  by  $K$ . As for our estimator, each time the selected value of  $K$  produced a negative estimator for a new curve  $v$ , we replaced, for that curve  $v$ ,  $K$  by the largest value smaller or equal to  $K$  which produced a positive estimator.

We also considered the kernel density estimator of Dabo-Niang (2004b), which requires the choice of a bandwidth. To choose it in practice we considered several versions of cross-validation and a nearest-neighbour bandwidth

Table 1: Simulation results for density estimation:  $10^4 \times$  median [first quartile, second quartile] of  $2 \times 10^6$  values of the SE.

Model	Method	$n = 500$	$n = 1000$	$n = 2000$	$n = 5000$
(i)	DM	635[145,2242]	492[120,1660]	395[103,1252]	316[86,953]
	DN	891[171,4122]	800[166,3439]	664[125,2970]	527[100,2271]
(ii)	DM	683[152,2427]	506[123,1732]	409[108,1293]	343[94,1051]
	DN	911[179,4133]	823[168,3568]	659[124,2990]	544[101,2420]
(iii)	DM	1134[237,4538]	898[188,3529]	813[175,3237]	784[165,3197]
	DN	1375[209,8046]	1200[186,7325]	1081[174,6611]	1025[177,5574]
(iv)	DM	908[194,3788]	801[172,3158]	744[154,3135]	590[124,2399]
	DN	1468[232,8351]	1151[183,6878]	1097[190,6514]	1052[196,5460]
(v)	DM	849[187,3287]	751[163,2812]	654[143,2500]	565[122,2273]
	DN	1097[170,6389]	1024[172,5817]	914[160,5133]	865[160,4309]

version of the estimator. However we encountered major numerical issues with denominators getting too close to zero and did not manage to obtain reasonable results. Therefore we do not consider this estimator in our numerical work.

The results of our simulations are summarised in Table 1 where, for each case and each sample size  $n$  we present  $10^4$  times the median and the first and third quartiles of the squared error  $SE = \{\hat{f}_Y(V) - f_Y(V)\}^2$



Table 2: Average computational time (in seconds) for computing one density estimator (including the CV choice of smoothing parameters).

Model	Method	$n = 500$	$n = 1000$	$n = 2000$	$n = 5000$
(i)	DM	94	114	130	198
	DN	42	46	54	77
(ii)	DM	95	113	135	200
	DN	49	55	68	96
(iii)	DM	102	116	138	218
	DN	50	53	71	97
(iv)	DM	104	110	127	191
	DN	46	47	59	82
(v)	DM	91	130	125	182
	DN	41	47	65	100

computed for the  $200 \times 10^4$   $V$  values. As expected by the theory, both estimators improved as sample size increased, and overall our estimator worked significantly better than Dabo-Niang's (2004a) estimator. In Table 2, for our estimator and that of Dabo-Niang (2004a), we also show the average time (in seconds and averaged over 10 simulated examples) required to compute one density estimator and its associated data-driven smoothing parameters on a Windows computer with Intel Xeon processor E5-2643 v4

and 32GB memory. Recall that our estimator requires the choice by CV of two smoothing parameters  $m$  and  $K$  whereas that of Dabo-Niang (2004a) requires to choose one smoothing parameter  $K$ . It is unsurprising then that our estimator requires longer computational time: this is the price to pay for the additional accuracy brought by choosing, in a data-driven way, two parameters instead of one.

### Acknowledgments

Delaigle's research was supported by a grant and a fellowship from the Australian Research Council. The authors are grateful to the editors and two referees for their helpful and valuable comments.

### References

- Appell, P. (1988). Sur une classe de polynômes. *Annales scientifiques de l'École Normale Supérieure 2me série* **9**, 119–144.
- Baïllo, A., Cuevas, A. and Cuesta-Albertos, J.A. (2011). Supervised classification for a family of Gaussian functional models. *Scand. J. Statist.* **38**, 480–498.

- Bugni, F.A., Hall, P., Horowitz, J.L. and Neumann, G.R. (2009). Goodness-of-fit tests for functional data. *Econometrics J.* **12**, 1–18.
- Cameron, R. H. and Martin, W. T. (1947). The orthogonal development of nonlinear functionals in series FourierHermite functionals. *Ann. Math.*, **48**, 385–392.
- Chagny, G. and Roche, A. (2014). Adaptive and minimax estimation of the cumulative distribution function given a functional covariate. *Elec. J. Statist.* **8**, 2352–2404.
- Chesneau, C., Kachour, M. and Maillot, B. (2013). Nonparametric estimation for functional data by wavelet thresholding. *REVSTAT – Statistical Journal* **11**, 211–230.
- Ciollaro, M., Genovese, Ch.R. and Wang, D. (2016). Nonparametric clustering of functional data using pseudo-densities. *Elec. J. Statist.* **10**, 2922-2972.
- Dabo-Niang, S. (2002). Estimation de la densité en dimension infinie: Application aux processus de type diffusion. *C.R. Acad. Sci. Paris I*, **334**, 213-216.

- Dabo-Niang, S. (2003). Density estimation in a separable metric space. *Pub. Inst. Stat. Univ. Paris*, **47**, fasc. 1-2, 3–21.
- Dabo-Niang, S. (2004a). Density estimation by orthogonal series in an infinite dimensional space: application to processes of diffusion type I. *J. Nonparametric Statistics*, **16**, 171–186.
- Dabo-Niang, S. (2004b). Kernel density estimator in an infinite-dimensional space with a rate of convergence in the case of diffusion process. *Applied Mathematics Letters* **17**, 381–386.
- Dabo-Niang, S. and Yao, A.-F. (2013). Kernel spatial density estimation in infinite dimension space. *Metrika* **76**, 19–52.
- Dai, X., Müller, H.-G. and Yao, F. (2017). Optimal Bayes classifiers for functional data and density ratios. *Biometrika* **104**, 545–560.
- Delaigle, A. and Hall, P. (2010). Defining probability density for a distribution of random functions. *Ann. Statist.* **38**, 1171–1193.
- Delaigle, A. and Hall, P. (2012). Achieving near-perfect classification for functional data. *J. Roy. Statist. Soc., Ser. B* **74**, 267–286.
- Delaigle, A., Hall, P. and Bathia, N. (2012). Componentwise classification and clustering of functional data. *Biometrika*, **99**, 299–313.

- Delaigle, A. and Hall, P. (2013). Classification using censored functional data. *J. Amer. Statist. Assoc.* **108**, 1269–1283.
- Escabias, M., Aguilera, A.M. and Valderrama, M.J. (2007). Functional PLS logit regression model. *Comput. Statist. Data Anal.* **51**, 4891–4902.
- Ferraty, F. and Vieu, P. (2006). *Nonparametric functional data analysis: theory and practice*, Springer.
- Girsanov, I. V. (1960). On transforming a certain class of stochastic processes by absolutely continuous substitution of measures. *Theo. Probab. Appl.* **5**, 285–301.
- Hall, P. and Hosseini-Nasab, M. (2006). On properties of functional principal components analysis. *J. Roy. Statist. Soc., Ser B* **68**, 109–126.
- Hall, P., Poskitt, D. and Presnell, B. (2001). A functional data-analytic approach to signal discrimination. *Technometrics* **43**, 1–9.
- Hall, R., Rinaldo, A. and Wasserman, L. (2013). Differential privacy for functions and functional data. *J. Mach. Learn. Research* **14**, 703–727.
- Jirak, M. (2016). Optimal eigen expansions and uniform bounds. *Probab. Theo. Rel. Fields* **166**, 753–799.

- Karwa, V. and Slavkovi, A. (2016). Inference using noisy degrees: differentially private  $\beta$ -model and synthetic graphs. *Ann. Statist.* **44**, 87–112.
- Kim, A.K.H. (2014). Minimax bounds for estimation of normal mixtures. *Bernoulli* **20**, 1802–1818.
- Lin, Z., Müller, H.-G. and Yao, F. (2018). Mixture inner product spaces and their application to functional data analysis. *Ann. Statist.* **46**, 370–400.
- Mas, A. (2012). Lower bound in regression for functional data by small ball probability representation in Hilbert space. *Elec. J. Statist.* **6**, 1745–1778.
- Mas A. and Ruymgaart, F. (2015). High dimensional principal projections. *Complex Anal. Operator Theo.* **9**, 35–63.
- Meister, A. (2016). Optimal classification and nonparametric regression for functional data. *Bernoulli* **22**, 1729–1744.
- Mirshani, A., Reimherr, M. and Slavkovic, A. (2017). Establishing statistical privacy for functional data via functional densities. *arXiv:1711.06660*.
- Prakasa Rao, B.L.S. (2010a). Nonparametric density estimation for functional data by delta sequences. *Braz. J. Probab. Stat.* **24**, 468–478.
- Prakasa Rao, B.L.S. (2010b). Nonparametric density estimation for functional data via wavelets. *Comm. Stat. – Theo. Meth.* **39**, 1608–1618.

Preda, C., Saporta, G. and Leveder, C. (2007). PLS classification of functional data. *Comput. Statist.* **22**, 223–235.

Shin, H. (2008). An extension of Fisher’s discriminant analysis for stochastic processes. *J. Mult. Anal.* **99**, 1191–1216.

Wasserman, L., and Zhou, S. (2010). A statistical framework for differential privacy. *J. Amer. Statist. Assoc.* **105**, 375–389.

Zhang, X.L., Begleiter, H., Porjesz, B., Wang, W. and Litke, A. (1995). Event related potentials during object recognition tasks. *Brain Research Bulletin* **38**, 531–538.

Aurore Delaigle

School of Mathematics and Statistics and Australian Research Council Centre of Excellence for Mathematical and Statistical Frontiers,

University of Melbourne, Australia.

E-mail: aurored@unimelb.edu.au

Alexander Meister

Institut für Mathematik,

Universität Rostock,

D-18051 Rostock, Germany.

E-mail: [alexander.meister@uni-rostock.de](mailto:alexander.meister@uni-rostock.de)

Statistica Sinica