

OPTIMAL ESTIMATION OF THE MEAN FUNCTION BASED ON DISCRETELY SAMPLED FUNCTIONAL DATA: PHASE TRANSITION

BY T. TONY CAI¹ AND MING YUAN²

University of Pennsylvania and Georgia Institute of Technology

The problem of estimating the mean of random functions based on discretely sampled data arises naturally in functional data analysis. In this paper, we study optimal estimation of the mean function under both common and independent designs. Minimax rates of convergence are established and easily implementable rate-optimal estimators are introduced. The analysis reveals interesting and different phase transition phenomena in the two cases. Under the common design, the sampling frequency solely determines the optimal rate of convergence when it is relatively small and the sampling frequency has no effect on the optimal rate when it is large. On the other hand, under the independent design, the optimal rate of convergence is determined jointly by the sampling frequency and the number of curves when the sampling frequency is relatively small. When it is large, the sampling frequency has no effect on the optimal rate. Another interesting contrast between the two settings is that smoothing is necessary under the independent design, while, somewhat surprisingly, it is not essential under the common design.

1. Introduction. Estimating the mean function based on discretely sampled noisy observations is one of the most basic problems in functional data analysis. Much progress has been made on developing estimation methodologies. The two monographs by Ramsay and Silverman (2002, 2005) provide comprehensive discussions on the methods and applications. See also Ferraty and Vieu (2006).

Let $X(\cdot)$ be a random function defined on the unit interval $\mathcal{T} = [0, 1]$ and X_1, \dots, X_n be a sample of n independent copies of X . The goal is to estimate the mean function $g_0(\cdot) := \mathbb{E}(X(\cdot))$ based on noisy observations from discrete locations on these curves:

$$(1.1) \quad Y_{ij} = X_i(T_{ij}) + \varepsilon_{ij}, \quad j = 1, 2, \dots, m_i \text{ and } i = 1, 2, \dots, n,$$

where T_{ij} are sampling points, and ε_{ij} are independent random noise variables with $\mathbb{E}\varepsilon_{ij} = 0$ and finite second moment $\mathbb{E}\varepsilon_{ij}^2 = \sigma_0^2 < +\infty$. The sample path of X

Received May 2010; revised May 2011.

¹Supported in part by NSF FRG Grant DMS-08-54973.

²Supported in part by NSF Career Award DMS-08-46234.

Key words and phrases. Functional data, mean function, minimax, rate of convergence, phase transition, reproducing kernel Hilbert space, smoothing splines, Sobolev space.

is assumed to be smooth in that it belongs to the usual Sobolev–Hilbert spaces of order r almost surely, such that

$$(1.2) \quad \mathbb{E} \left(\int_{\mathcal{T}} [X^{(r)}(t)]^2 dt \right) < +\infty.$$

Such problems naturally arise in a variety of applications and are typical in functional data analysis [see, e.g., Ramsay and Silverman (2005), Ferraty and Vieu (2006)]. Various methods have been proposed. However, little is known about their theoretical properties.

In the present paper, we study optimal estimation of the mean function in two different settings. One is when the observations are sampled at the same locations across curves, that is, $T_{1j} = T_{2j} = \dots = T_{mj} =: T_j$ for all $j = 1, \dots, m$. We shall refer to this setting as *common design* because the sampling locations are common to all curves. Another setting is when the T_{ij} are independently sampled from \mathcal{T} , which we shall refer to as *independent design*. We establish the optimal rates of convergence for estimating the mean function in both settings. Our analysis reveals interesting and different phase transition phenomena in the two cases. Another interesting contrast between the two settings is that smoothing is necessary under the independent design, while, somewhat surprisingly, it is not essential under the common design. We remark that under the independent design, the number of sampling points oftentimes varies from curve to curve and may even be random itself. However, for ease of presentation and better illustration of similarities and differences between the two types of designs, we shall assume an equal number of sampling points on each curve in the discussions given in this section.

Earlier studies of nonparametric estimation of the mean function g_0 from a collection of discretely sampled curves can be traced back to at least Hart and Wehrly (1986) and Rice and Silverman (1991) in the case of common design. In this setting, ignoring the temporal nature of $\{T_j : 1 \leq j \leq m\}$, the problem of estimating g_0 can be translated into estimating the mean vector $(g_0(T_1), \dots, g_0(T_m))'$, a typical problem in multivariate analysis. Such notions are often quickly discarded because they essentially lead to estimating $g_0(T_j)$ by its sample mean

$$(1.3) \quad \bar{Y}_{.j} = \frac{1}{n} \sum_{i=1}^n Y_{ij},$$

based on the standard Gauss–Markov theory [see, e.g., Rice and Silverman (1991)].

Note that $\mathbb{E}(Y_{ij}|T) = g_0(T_{ij})$ and that the smoothness of X implies that g_0 is also smooth. It is therefore plausible to assume that smoothing is essential for optimal estimation of g_0 . For example, a natural approach for estimating g_0 is to regress Y_{ij} on T_{ij} nonparametrically via kernel or spline smoothing. Various methods have been introduced along this vein [see, e.g., Rice and Silverman (1991)]. However, not much is known about their theoretical properties. It is noteworthy

that this setting differs from the usual nonparametric smoothing in that the observations from the same curve are highly correlated. Nonparametric smoothing with certain correlated errors has been previously studied by Hall and Hart (1990), Wang (1996) and Johnstone and Silverman (1997), among others. Interested readers are referred to Opsomer, Wang and Yang (2001) for a recent survey of existing results. But neither of these earlier developments can be applied to account for the dependency induced by the functional nature in our setting. To comprehend the effectiveness of smoothing in the current context, we establish minimax bounds on the convergence rate of the integrated squared error for estimating g_0 .

Under the common design, it is shown that the minimax rate is of the order $m^{-2r} + n^{-1}$ where the two terms can be attributed to discretization and stochastic error, respectively. This rate is fundamentally different from the usual nonparametric rate of $(nm)^{-2r/(2r+1)}$ when observations are obtained at nm distinct locations in order to recover an r times differentiable function [see, e.g., Stone (1982)]. The rate obtained here is jointly determined by the sampling frequency m and the number of curves n rather than the total number of observations mn . A distinct feature of the rate is the phase transition which occurs when m is of the order $n^{1/2r}$. When the functions are sparsely sampled, that is, $m = O(n^{1/2r})$, the optimal rate is of the order m^{-2r} , solely determined by the sampling frequency. On the other hand, when the sampling frequency is high, that is, $m \gg n^{1/2r}$, the optimal rate remains $1/n$ regardless of m . Moreover, our development uncovers a surprising fact that *interpolation* of $\{(T_j, \bar{Y}_j) : j = 1, \dots, m\}$, that is, estimating $g_0(T_j)$ by \bar{Y}_j , is rate optimal. In other words, contrary to the conventional wisdom, smoothing does not result in improved convergence rates.

In addition to the common design, another popular sampling scheme is the independent design where the T_{ij} are independently sampled from \mathcal{T} . A natural approach is to smooth observations from each curve separately and then average over all smoothed estimates. However, the success of this two-step procedure hinges upon the availability of a reasonable estimate for each individual curve. In contrast to the case of common design, we show that under the independent design, the minimax rate for estimating g_0 is $(nm)^{-2r/(2r+1)} + n^{-1}$, which can be attained by smoothing $\{(T_{ij}, Y_{ij}) : 1 \leq i \leq n, 1 \leq j \leq m\}$ altogether. This implies that in the extreme case of $m = 1$, the optimal rate of estimating g_0 is $n^{-2r/(2r+1)}$, which also suggests the sub-optimality of the aforementioned two-step procedure because it is impossible to smooth a curve with only a single observation. Similar to the common design, there is a phase transition phenomenon in the optimal rate of convergence with a boundary at $m = n^{1/2r}$. When the sampling frequency m is small, that is, $m = O(n^{1/2r})$, the optimal rate is of the order $(nm)^{-2r/(2r+1)}$ which depends jointly on the values of both m and n . In the case of high sampling frequency with $m \gg n^{1/2r}$, the optimal rate is always $1/n$ and does not depend on m .

It is interesting to compare the minimax rates of convergence in the two settings. The phase transition boundary for both designs occurs at the same value,

$m = n^{1/2r}$. When m is above the boundary, that is, $m \geq n^{1/2r}$, there is no difference between the common and independent designs, and both have the optimal rate of n^{-1} . When m is below the boundary, that is, $m \ll n^{1/2r}$, the independent design is always superior to the common design in that it offers a faster rate of convergence.

Our results connect with several observations made earlier in the literature on longitudinal and functional data analysis. Many longitudinal studies follow the independent design, and the number of sampling points on each curve is typically small. In such settings, it is widely recognized that one needs to pool the data to obtain good estimates, and the two-step procedure of averaging the smooth curves may be suboptimal. Our analysis here provides a rigorous justification for such empirical observations by pinpointing to what extent the two-step procedure is suboptimal. The phase transition observed here also relates to the earlier work by Hall, Müller and Wang (2006) on estimating eigenfunctions of the covariance kernel when the number of sampling points is either fixed or of larger than $n^{1/4+\delta}$ for some $\delta > 0$. It was shown that the eigenfunctions can be estimated at the rate of $n^{-4/5}$ in the former case and $1/n$ in the latter. We show here that estimating the mean function has similar behavior. Furthermore, we characterize the exact nature of such transition behavior as the sampling frequency changes.

The rest of the paper is organized as follows. In Section 2 the optimal rate of convergence under the common design is established. We first derive a minimax lower bound and then show that the lower bound is in fact rate sharp. This is accomplished by constructing a rate-optimal smoothing splines estimator. The minimax upper bound is obtained separately for the common fixed design and common random design. Section 3 considers the independent design and establishes the optimal rate of convergence in this case. The rate-optimal estimators are easily implementable. Numerical studies are carried out in Section 4 to demonstrate the theoretical results. Section 5 discusses connections and differences of our results with other related work. All proofs are relegated to Section 6.

2. Optimal rate of convergence under common design. In this section we consider the common design where each curve is observed at the same set of locations $\{T_j : 1 \leq j \leq m\}$. We first derive a minimax lower bound and then show that this lower bound is sharp by constructing a smoothing splines estimator that attains the same rate of convergence as the lower bound.

2.1. Minimax lower bound. Let $\mathcal{P}(r; M_0)$ be the collection of probability measures for a random function X such that its sample path is r times differentiable almost surely and

$$(2.1) \quad \mathbb{E} \int_{\mathcal{T}} [X^{(r)}(t)]^2 dt \leq M_0$$

for some constant $M_0 > 0$. Our first main result establishes the minimax lower bound for estimating the mean function over $\mathcal{P}(r; M_0)$ under the common design.

THEOREM 2.1. *Suppose the sampling locations are common in model (1.1). Then there exists a constant $d > 0$ depending only on M_0 and the variance σ_0^2 of measurement error ε_{ij} such that for any estimate \tilde{g} based on observations $\{(T_j, Y_{ij}) : 1 \leq i \leq n, 1 \leq j \leq m\}$,*

$$(2.2) \quad \limsup_{n \rightarrow \infty} \sup_{\mathcal{L}(X) \in \mathcal{P}(r; M_0)} P(\|\tilde{g} - g_0\|_{\mathcal{L}_2}^2 > d(m^{-2r} + n^{-1})) > 0.$$

The lower bound established in Theorem 2.1 holds true for both common fixed design where T_j 's are deterministic, and common random design where T_j 's are also random. The term m^{-2r} in the lower bound is due to the deterministic approximation error, and the term n^{-1} is attributed to the stochastic error. It is clear that neither can be further improved. To see this, first consider the situation where there is no stochastic variation and the mean function g_0 is observed exactly at the points T_j , $j = 1, \dots, m$. It is well known [see, e.g., DeVore and Lorentz (1993)] that due to discretization, it is not possible to recover g_0 at a rate faster than m^{-2r} for all g_0 such that $\int [g_0^{(r)}]^2 \leq M_0$. On the other hand, the second term n^{-1} is inevitable since the mean function g_0 cannot be estimated at a faster rate even if the whole random functions X_1, \dots, X_n are observed completely. We shall show later in this section that the rate given in the lower bound is optimal in that it is attainable by a smoothing splines estimator.

It is interesting to notice the phase transition phenomenon in the minimax bound. When the sampling frequency m is large, it has no effect on the rate of convergence, and g_0 can be estimated at the rate of $1/n$, the best possible rate when the whole functions were observed. More surprisingly, such saturation occurs when m is rather small, that is, of the order $n^{1/2r}$. On the other hand, when the functions are sparsely sampled, that is, $m = O(n^{1/2r})$, the rate is determined only by the sampling frequency m . Moreover, the rate m^{-2r} is in fact also the optimal interpolation rate. In other words, when the functions are sparsely sampled, the mean function g_0 can be estimated as well as if it is observed directly without noise.

The rate is to be contrasted with the usual nonparametric regression with nm observations at arbitrary locations. In such a setting, it is well known [see, e.g., Tsybakov (2009)] that the optimal rate for estimating g_0 is $(mn)^{-2r/(2r+1)}$, and typically stochastic error and approximation error are of the same order to balance the bias-variance trade-off.

2.2. Minimax upper bound: Smoothing splines estimate. We now consider the upper bound for the minimax risk and construct specific rate optimal estimators under the common design. These upper bounds show that the rate of convergence given in the lower bound established in Theorem 2.1 is sharp. More specifically, it is shown that a smoothing splines estimator attains the optimal rate of convergence over the parameter space $\mathcal{P}(r; M_0)$.

We shall consider a smoothing splines type of estimate suggested by Rice and Silverman (1991). Observe that $f \mapsto \int [f^{(r)}]^2$ is a squared semi-norm and therefore convex. By Jensen’s inequality,

$$(2.3) \quad \int_{\mathcal{T}} [g_0^{(r)}(t)]^2 dt \leq \mathbb{E} \int_{\mathcal{T}} [X^{(r)}(t)]^2 dt < \infty,$$

which implies that g_0 belongs to the r th order Sobolev–Hilbert space,

$$\mathcal{W}_2^r([0, 1]) = \{g : [0, 1] \rightarrow \mathbf{R} \mid g, g^{(1)}, \dots, g^{(r-1)} \text{ are absolutely continuous and } g^{(r)} \in \mathcal{L}_2([0, 1])\}.$$

Taking this into account, the following smoothing splines estimate can be employed to estimate g_0 :

$$(2.4) \quad \hat{g}_\lambda = \arg \min_{g \in \mathcal{W}_2^r} \left\{ \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m (Y_{ij} - g(T_j))^2 + \lambda \int_{\mathcal{T}} [g^{(r)}(t)]^2 dt \right\},$$

where $\lambda > 0$ is a tuning parameter that balances the fidelity to the data and the smoothness of the estimate.

Similarly to the smoothing splines for the usual nonparametric regression, \hat{g}_λ can be conveniently computed, although the minimization is taken over an infinitely-dimensional functional space. First observe that \hat{g}_λ can be equivalently rewritten as

$$(2.5) \quad \hat{g}_\lambda = \arg \min_{g \in \mathcal{W}_2^r} \left\{ \frac{1}{m} \sum_{j=1}^m (\bar{Y}_{\cdot j} - g(T_j))^2 + \lambda \int_{\mathcal{T}} [g^{(r)}(t)]^2 dt \right\}.$$

Appealing to the so-called representer theorem [see, e.g., Wahba (1990)], the solution of the minimization problem can be expressed as

$$(2.6) \quad \hat{g}_\lambda(t) = \sum_{k=0}^{r-1} d_k t^k + \sum_{j=1}^m c_j K(t, T_j)$$

for some coefficients $d_0, \dots, d_{r-1}, c_1, \dots, c_m$, where

$$(2.7) \quad K(s, t) = \frac{1}{(r!)^2} B_r(s) B_r(t) - \frac{1}{(2r)!} B_{2r}(|s - t|),$$

where $B_m(\cdot)$ is the m th Bernoulli polynomial. Plugging (2.6) back into (2.5), the coefficients and subsequently \hat{g}_λ can be solved in a straightforward way. This observation makes the smoothing splines procedure easily implementable. The readers are referred to Wahba (1990) for further details.

Despite the similarity between \hat{g}_λ and the smoothing splines estimate in the usual nonparametric regression, they have very different asymptotic properties. It is shown in the following that \hat{g}_λ achieves the lower bound established in Theorem 2.1.

The analyses for the common fixed design and the common random design are similar, and we shall focus on the fixed design where the common sampling locations T_1, \dots, T_m are deterministic. In this case, we assume without loss of generality that $T_1 \leq T_2 \leq \dots \leq T_m$. The following theorem shows that the lower bound established in Theorem 2.1 is attained by the smoothing splines estimate \hat{g}_λ .

THEOREM 2.2. *Consider the common fixed design and assume that*

$$(2.8) \quad \max_{0 \leq j \leq m} |T_{j+1} - T_j| \leq C_0 m^{-1}$$

for some constant $C_0 > 0$ where we follow the convention that $T_0 = 0$ and $T_{m+1} = 1$. Then

$$(2.9) \quad \lim_{D \rightarrow \infty} \limsup_{n \rightarrow \infty} \sup_{\mathcal{L}(X) \in \mathcal{P}(r; M_0)} P(\|\hat{g}_\lambda - g_0\|_{\mathcal{L}_2}^2 > D(m^{-2r} + n^{-1})) = 0$$

for any $\lambda = O(m^{-2r} + n^{-1})$.

Together with Theorem 2.1, Theorem 2.2 shows that \hat{g}_λ is minimax rate optimal if the tuning parameter λ is set to be of the order $O(m^{-2r} + n^{-1})$. We note the necessity of the condition given by (2.8). It is clearly satisfied when the design is equidistant, that is, $T_j = 2j/(2m + 1)$. The condition ensures that the random functions are observed on a sufficiently regular grid.

It is of conceptual importance to compare the rate of \hat{g}_λ with those generally achieved in the usual nonparametric regression setting. Defined by (2.4), \hat{g}_λ essentially regresses Y_{ij} on T_j . Similarly to the usual nonparametric regression, the validity of the estimate is driven by $\mathbb{E}(Y_{ij}|T_j) = g_0(T_j)$. The difference, however, is that Y_{i1}, \dots, Y_{im} are highly correlated because they are observed from the same random function $X_i(\cdot)$. When all the Y_{ij} 's are independently sampled at T_j 's, it can be derived that the optimal rate for estimating g_0 is $m^{-2r} + (mn)^{-2r/(2r+1)}$. As we show here, the dependency induced by the functional nature of our problem leads to the different rate $m^{-2r} + n^{-1}$.

A distinct feature of the behavior of \hat{g}_λ is in the choice of the tuning parameter λ . Tuning parameter selection plays a paramount role in the usual nonparametric regression, as it balances the the tradeoff between bias and variance. Optimal choice of λ is of the order $(mn)^{-2r/(2r+1)}$ in the usual nonparametric regression. In contrast, in our setting, more flexibility is allowed in the choice of the tuning parameter in that \hat{g}_λ is rate optimal so long as λ is sufficiently small. In particular, taking $\lambda \rightarrow 0^+$, \hat{g}_λ reduces to the splines interpolation, that is, the solution to

$$(2.10) \quad \min_{g \in \mathcal{W}_2^r} \int_{\mathcal{T}} [g^{(r)}(t)]^2 \quad \text{subject to } g(T_j) = \bar{Y}_{.j}, \quad j = 1, \dots, m.$$

This amounts to, in particular, estimating $g_0(T_j)$ by $\bar{Y}_{.j}$. In other words, there is no benefit from smoothing in terms of the convergence rate. However, as we will see in Section 4, smoothing can lead to improved finite sample performance.

REMARK. More general statements can also be made without the condition on the spacing of sampling points. More specifically, denote by

$$R(T_1, \dots, T_m) = \max_j |T_{j+1} - T_j|$$

the discretization resolution. Using the same argument, one can show that the optimal convergence rate in the minimax sense is $R^{2r} + n^{-1}$ and \hat{g}_λ is rate optimal so long as $\lambda = O(R^{2r} + n^{-1})$.

REMARK. Although we have focused here on the case when the sampling points are deterministic, a similar statement can also be made for the setting where the sampling points are random. In particular, assuming that T_j are independent and identically distributed with a density function η such that $\inf_{t \in \mathcal{T}} \eta(t) \geq c_0 > 0$ and $g_0 \in \mathcal{W}_\infty^r$, it can be shown that the smoothing splines estimator \hat{g}_λ satisfies

$$(2.11) \quad \lim_{D \rightarrow \infty} \limsup_{n \rightarrow \infty} \sup_{\mathcal{L}(X) \in \mathcal{P}(r; M_0)} P(\|\hat{g}_\lambda - g_0\|_{\mathcal{L}_2}^2 > D(m^{-2r} + n^{-1})) = 0$$

for any $\lambda = O(m^{-2r} + n^{-1})$. In other words, \hat{g}_λ remains rate optimal.

3. Optimal rate of convergence under independent design. In many applications, the random functions X_i are not observed at common locations. Instead, each curve is discretely observed at a different set of points [see, e.g., James and Hastie (2001), Rice and Wu (2001), Diggle et al. (2002), Yao, Müller and Wang (2005)]. In these settings, it is more appropriate to model the sampling points T_{ij} as independently sampled from a common distribution. In this section we shall consider optimal estimation of the mean function under the independent design.

Interestingly, the behavior of the estimation problem is drastically different between the common design and the independent design. To keep our treatment general, we allow the number of sampling points to vary. Let m be the harmonic mean of m_1, \dots, m_n , that is,

$$m := \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \right)^{-1}.$$

Denote by $\mathcal{M}(m)$ the collection of sampling frequencies (m_1, \dots, m_n) whose harmonic mean is m . In parallel to Theorem 2.1, we have the following minimax lower bound for estimating g_0 under the independent design.

THEOREM 3.1. *Suppose T_{ij} are independent and identically distributed with a density function η such that $\inf_{t \in \mathcal{T}} \eta(t) \geq c_0 > 0$. Then there exists a constant $d > 0$ depending only on M_0 and σ_0^2 such that for any estimate \tilde{g} based on observations $\{(T_{ij}, Y_{ij}) : 1 \leq i \leq n, 1 \leq j \leq m\}$,*

$$(3.1) \quad \limsup_{n \rightarrow \infty} \sup_{\substack{\mathcal{L}(X) \in \mathcal{P}(r; M_0) \\ (m_1, \dots, m_n) \in \mathcal{M}(m)}}} P(\|\tilde{g} - g_0\|_{\mathcal{L}_2}^2 > d((nm)^{-2r/(2r+1)} + n^{-1})) > 0.$$

The minimax lower bound given in Theorem 3.1 can also be achieved using the smoothing splines type of estimate. To account for the different sampling frequency for different curves, we consider the following estimate of g_0 :

$$(3.2) \quad \hat{g}_\lambda = \arg \min_{g \in \mathcal{W}_2^r} \left\{ \frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \sum_{j=1}^{m_i} (Y_{ij} - g(T_{ij}))^2 + \lambda \int_{\mathcal{T}} [g^{(r)}(t)]^2 dt \right\}.$$

THEOREM 3.2. *Under the conditions of Theorem 3.1, if $\lambda \asymp (nm)^{-2r/(2r+1)}$, then the smoothing splines estimator \hat{g}_λ satisfies*

$$(3.3) \quad \lim_{D \rightarrow \infty} \limsup_{n \rightarrow \infty} \sup_{\substack{\mathcal{L}(X) \in \mathcal{P}(r; M_0) \\ (m_1, \dots, m_n) \in \mathcal{M}(m)}}} P(\|\hat{g} - g_0\|_{\mathcal{L}_2}^2 > D((nm)^{-2r/(2r+1)} + n^{-1})) = 0.$$

In other words, \hat{g}_λ is rate optimal.

Theorems 3.1 and 3.2 demonstrate both similarities and significant differences between the two types of designs in terms of the convergence rate. For either the common design or the independent design, the sampling frequency only plays a role in determining the convergence rate when the functions are sparsely sampled, that is, $m = O(n^{1/2r})$. But how the sampling frequency affects the convergence rate when each curve is sparsely sampled differs between the two designs. For the independent design, the total number of observations mn , whereas for the common design m alone, determines the minimax rate. It is also noteworthy that when $m = O(n^{1/2r})$, the optimal rate under the independent design, $(mn)^{-2r/(2r+1)}$, is the same as if all the observations are independently observed. In other words, the dependency among Y_{i1}, \dots, Y_{im} does not affect the convergence rate in this case.

REMARK. We emphasize that Theorems 3.1 and 3.2 apply to both deterministic and random sampling frequencies. In particular for random sampling frequencies, together with the law of large numbers, the same minimax bound holds when we replace the harmonic mean by

$$\left(\frac{1}{n} \sum_{i=1}^n \mathbb{E}(1/m_i) \right)^{-1},$$

when assuming that m_i 's are independent.

3.1. Comparison with two-stage estimate. A popular strategy to handle discretely sampled functional data in practice is a two-stage procedure. In the first

step, nonparametric regression is run for data from each curve to obtain estimate \tilde{X}_i of $X_i, i = 1, 2, \dots, n$. For example, they can be obtained by smoothing splines

$$(3.4) \quad \tilde{X}_{i,\lambda} = \arg \min_{f \in \mathcal{W}_2^r} \left\{ \frac{1}{m} \sum_{j=1}^m (Y_{ij} - g(T_{ij}))^2 + \lambda \int_{\mathcal{T}} [f^{(r)}(t)]^2 dt \right\}.$$

Any subsequent inference can be carried out using the $\tilde{X}_{i,\lambda}$ as if they were the original true random functions. In particular, the mean function g_0 can be estimated by the simple average

$$(3.5) \quad \tilde{g}_\lambda = \frac{1}{n} \sum_{i=1}^n \tilde{X}_{i,\lambda}.$$

Although formulated differently, it is worth pointing out that this procedure is equivalent to the smoothing splines estimate \hat{g}_λ under the common design.

PROPOSITION 3.3. *Under the common design, that is, $T_{ij} = T_j$ for $1 \leq i \leq n$ and $1 \leq j \leq m$. The estimate \tilde{g}_λ from the two-stage procedure is equivalent to the smoothing splines estimate \hat{g}_λ : $\hat{g}_\lambda = \tilde{g}_\lambda$.*

In light of Theorems 2.1 and 2.2, the two-step procedure is also rate optimal under the common design. But it is of great practical importance to note that in order to achieve the optimality, it is critical that in the first step we *undersmooth* each curve by using a sufficiently small tuning parameter.

Under independent design, however, the equivalence no longer holds. The success of the two-step estimate \tilde{g}_λ depends upon getting a good estimate of each curve, which is not possible when m is very small. In the extreme case of $m = 1$, the procedure is no longer applicable, but Theorem 3.2 indicates that smoothing splines estimate \hat{g}_λ can still achieve the optimal convergence rate of $n^{-2r/(2r+1)}$. The readers are also referred to Hall, Müller and Wang (2006) for discussions on the pros and cons of similar two-step procedures in the context of estimating the functional principal components.

4. Numerical experiments. The smoothing splines estimators are easy to implement. To demonstrate the practical implications of our theoretical results, we carried out a set of simulation studies. The true mean function g_0 is fixed as

$$(4.1) \quad g_0 = \sum_{k=1}^{50} 4(-1)^{k+1} k^{-2} \phi_k,$$

where $\phi_1(t) = 1$ and $\phi_{k+1}(t) = \sqrt{2} \cos(k\pi t)$ for $k \geq 1$. The random function X was generated as

$$(4.2) \quad X = g_0 + \sum_{k=1}^{50} \zeta_k Z_k \phi_k,$$

where Z_k are independently sampled from the uniform distribution on $[-\sqrt{3}, \sqrt{3}]$, and ζ_k are deterministic. It is not hard to see that ζ_k^2 are the eigenvalues of the covariance function of X and therefore determine the smoothness of a sample curve. In particular, we take $\zeta_k = (-1)^{k+1}k^{-1.1/2}$. It is clear that the sample path of X belongs to the second order Sobolev space ($r = 2$).

We begin with a set of simulations designed to demonstrate the effect of interpolation and smoothing under common design. A data set of fifty curves were first simulated according to the aforementioned scheme. For each curve, ten noisy observations were taken at equidistant locations on each curve following model (1.1) with $\sigma_0^2 = 0.5^2$. The observations, together with g_0 (grey line), are given in the right panel of Figure 1. Smoothing splines estimate \hat{g}_λ is also computed with a variety of values for λ . The integrated squared error, $\|\hat{g}_\lambda - g_0\|_{\mathcal{L}_2}$, as a function of the tuning parameter λ is given in the left panel. For λ smaller than 0.1, the smoothing splines estimate essentially reduces to the spline interpolation. To contrast the effect of interpolation and smoothing, the right panel also includes the interpolation estimate (solid black line) and \hat{g}_λ (red dashed line) with the tuning parameter chosen to minimize the integrated squared error. We observe from the

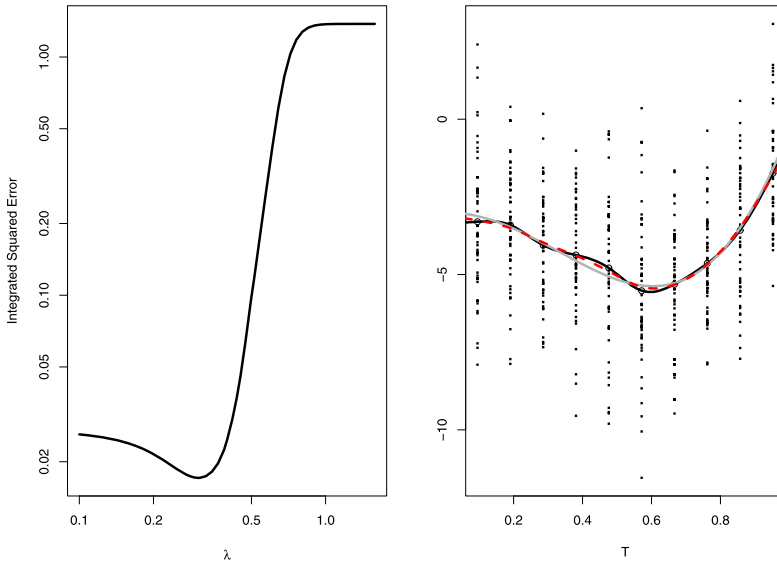


FIG. 1. *Effect of smoothing under common design: for a typical data set with fifty curves, ten observations were taken on each curve. The observations and g_0 (solid grey line) are given in the right panel together with the spline interpolation estimate (solid black line) and smoothing splines estimate (red dashed line) with the tuning parameter chosen to yield the smallest integrated squared error. The left panel gives the integrated squared error of the smoothing splines estimate as a function of the tuning parameter. It is noted that the smoothing splines estimate essentially reduces to the spline interpolation for λ smaller than 0.1.*

figure that smoothing does lead to slightly improved finite sample performance although it does not affect the convergence rate as shown in Section 2.

The next numerical experiment intends to demonstrate the effect of sample size n , sampling frequency m as well as design. To this end, we simulated n curves, and from each curve, m discrete observations were taken following model (1.1) with $\sigma_0^2 = 0.5^2$. The sampling locations are either fixed at $T_j = (2j)/(2m + 1)$, $j = 1, \dots, m$, for common design or randomly sampled from the uniform distribution on $[0, 1]$. The smoothing splines estimate \hat{g}_λ for each simulated data set, and the tuning parameter is set to yield the smallest integrated squared error and therefore reflect the best performance of the estimating procedure for each data set. We repeat the experiment with varying combinations of $n = 25, 50$ or 200 , $m = 1, 5, 10$ or 50 . For the common design, we restrict to $m = 10$ or 50 to give more meaningful comparison. The true function g_0 as well as its estimates obtained in each of the settings are given in Figure 2.

Figure 2 agrees pretty well with our theoretical results. For instance, increasing either m or n leads to improved estimates, whereas such improvement is more visible for small values of m . Moreover, for the same value of m and n , independent designs tend to yield better estimates.

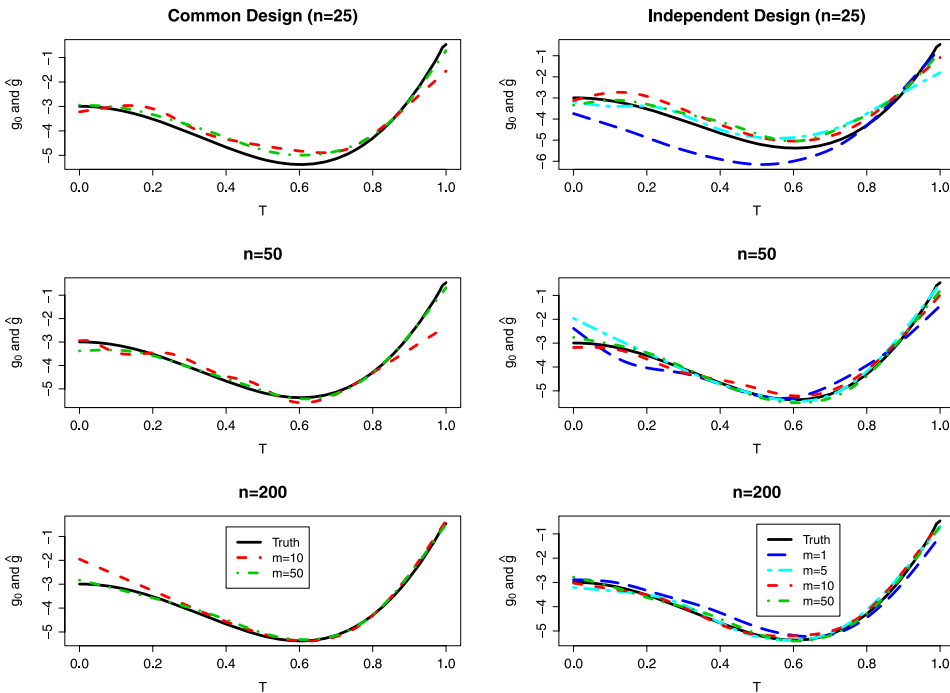


FIG. 2. Effect of m, n and type of design on estimating g_0 : smoothing splines estimates obtained under various combinations are plotted together with g_0 .

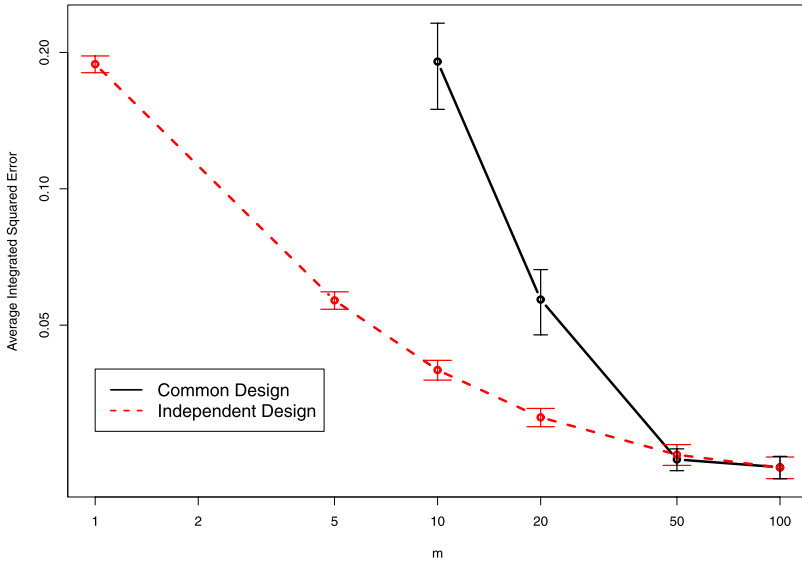


FIG. 3. Effect of design type and sampling frequency on estimating g_0 : the black solid line and circles correspond to common design whereas the red dashed lines and circles correspond to independent design. The error bars correspond to the average \pm one standard errors based on two hundred repetitions. Note that both axes are in log scale to yield better comparison.

To further contrast the two types of designs and the effect of sampling frequency on estimating g_0 , we now fix the number of curves at $n = 100$. For the common design, we consider $m = 10, 20, 50$ or 100 . For the independent design, we let $m = 1, 5, 10, 20, 50$ or 100 . For each combination of (n, m) , two hundred data sets were simulated following the same mechanism as before. Figure 3 gives the estimation error averaged over the one hundred data sets for each combination of (n, m) . It clearly shows that independent design is preferable over common design when m is small; and the two types of designs are similar when m is large. Both phenomena are in agreement with our theoretical results developed in the earlier sections.

5. Discussions. We have established the optimal rates of convergence for estimating the mean function under both the common design and independent design. The results reveal several significant differences in the behavior of the minimax estimation problem between the two designs. These revelations have important theoretical and practical implications. In particular, for sparsely sampled functions, the independent design leads to a faster rate of convergence when compared to the common design and thus should be preferred in practice.

The optimal rates of convergence for estimating the mean function based on discretely sampled random functions behave in a fundamentally different way from the minimax rate of convergence in the conventional nonparametric regression

problems. The optimal rates in the mean function estimation are jointly determined by the sampling frequency m and the number of curves n rather than the total number of observations mn .

The observation that one can estimate the mean function as well as if the whole curves are available when $m \gg n^{1/2r}$ bears some similarity to some recent findings on estimating the covariance kernel and its eigenfunction under independent design. Assuming that X is twice differentiable (i.e., $r = 2$), Hall, Müller and Wang (2006) showed that when $m \gg n^{1/4+\delta}$ for some $\delta > 0$, the covariance kernel and its eigenfunctions can be estimated at the rate of $1/n$ when using a two-step procedure. More recently, the cutoff point is further improved to $n^{1/2r} \log n$ with general r by Cai and Yuan (2010) using an alternative method. Intuitively one may expect estimating the covariance kernel to be more difficult than estimating the mean function, which suggests that these results may not be improved much further for estimating the covariance kernel or its eigenfunctions.

We have also shown that the particular smoothing splines type of estimate discussed earlier by Rice and Silverman (1991) attains the optimal convergence rates under both designs with appropriate tuning. The smoothing splines estimator is well suited for nonparametric estimation over Sobolev spaces. We note, however, other nonparametric techniques such as kernel or local polynomial estimators can also be used. We expect that kernel smoothing or other methods with proper choice of the tuning parameters can also achieve the optimal rate of convergence. Further study in this direction is beyond the scope of the current paper, and we leave it for future research.

Finally, we emphasize that although we have focused on the univariate Sobolev space for simplicity, the phenomena observed and techniques developed apply to more general functional spaces. Consider, for example, the multivariate setting where $\mathcal{T} = [0, 1]^d$. Following the same arguments, it can be shown that the minimax rate for estimating an r -times differentiable function is $m^{-2r/d} + n^{-1}$ under the common design and $(nm)^{-2r/(2r+d)} + n^{-1}$ under the independent design. The phase transition phenomena thus remain in the multidimensional setting under both designs with a transition boundary of $m = n^{d/2r}$.

6. Proofs.

PROOF OF THEOREM 2.1. Let \mathcal{D} be the collection all measurable functions of $\{(T_{ij}, Y_{ij}) : 1 \leq i \leq n, 1 \leq j \leq m\}$. First note that it is straightforward to show that

$$\limsup_{n \rightarrow \infty} \inf_{\tilde{g} \in \mathcal{D}} \sup_{\mathcal{L}(X) \in \mathcal{P}(r; M_0)} P(\|\tilde{g} - g_0\|_{\mathcal{L}_2}^2 > dn^{-1}) > 0$$

by considering X as an unknown constant function where the problem essentially becomes estimating the mean from n i.i.d. observations, and $1/n$ is known as the

optimal rate. It now suffices to show that

$$\limsup_{n \rightarrow \infty} \inf_{\tilde{g} \in \mathcal{D}} \sup_{\mathcal{L}(X) \in \mathcal{P}(r; M_0)} P(\|\tilde{g} - g_0\|_{\mathcal{L}_2}^2 > dm^{-2r}) > 0.$$

Let $\varphi_1, \dots, \varphi_{2m}$ be $2m$ functions from \mathcal{W}_2^r with distinct support, that is,

$$\varphi_k(\cdot) = h^r K\left(\frac{\cdot - t_k}{h}\right), \quad k = 1, \dots, 2m,$$

where $h = 1/(2m)$, $t_k = (k - 1)/2m + 1/4m$, and $K : \mathbf{R} \rightarrow [0, \infty)$ is an r times differentiable function with support $[-1/2, 1/2]$. See [Tsybakov \(2009\)](#) for explicit construction of such functions.

For each $b = (b_1, \dots, b_{2m}) \in \{0, 1\}^{2m}$, define

$$g_b(\cdot) = \sum_{k=1}^{2m} b_k \varphi_k(\cdot).$$

It is clear that

$$\min_{H(b, b') \geq 1} \frac{\|g_b - g_{b'}\|_{\mathcal{L}_2}^2}{H(b, b')} = \|\varphi_k\|_{\mathcal{L}_2}^2 = (2m)^{-(2r+1)} \|K\|_{\mathcal{L}_2}^2.$$

The claim then follows from an application of Assouad’s lemma [[Assouad \(1983\)](#)]. □

PROOF OF THEOREM 2.2. It is well known [see, e.g., [Green and Silverman \(1994\)](#)] that \hat{g}_λ can be characterized as the solution to the following:

$$\min_{g \in \mathcal{W}_2^r} \int_{\mathcal{T}} [g^{(r)}(t)]^2 dt \quad \text{subject to } g(T_j) = \hat{g}_\lambda(T_j), \quad j = 1, \dots, m.$$

Write

$$\delta_j = \hat{g}_\lambda(T_{ij}) - g_0(T_{ij}),$$

and let h be the linear interpolation of $\{(T_j, \delta_j) : 1 \leq j \leq m\}$, that is,

$$(6.1) \quad h(t) = \begin{cases} \delta_1, & 0 \leq t \leq T_1, \\ \delta_j \frac{T_{j+1} - t}{T_{j+1} - T_j} + \delta_{j+1} \frac{t - T_j}{T_{j+1} - T_j}, & T_j \leq t \leq T_{j+1}, \\ \delta_m, & T_m \leq t \leq 1. \end{cases}$$

Then $\hat{g}_\lambda = Q_T(g_0 + h)$ where Q_T be the operator associated with the r th order spline interpolation, that is, $Q_T(f)$ is the solution to

$$\min_{g \in \mathcal{W}_2^r} \int_{\mathcal{T}} [g^{(r)}(t)]^2 dt \quad \text{subject to } g(T_j) = f(T_j), \quad j = 1, \dots, m.$$

Recall that Q_T is a linear operator in that $Q_T(f_1 + f_2) = Q_T(f_1) + Q_T(f_2)$ [see, e.g., DeVore and Lorentz (1993)]. Therefore, $\hat{g}_\lambda = Q_T(g_0) + Q_T(h)$. By the triangular inequality,

$$(6.2) \quad \|\hat{g} - g_0\|_{\mathcal{L}_2} \leq \|Q_T(g_0) - g_0\|_{\mathcal{L}_2} + \|Q_T(h)\|_{\mathcal{L}_2}.$$

The first term on the right-hand side represents the approximation error of spline interpolation for g_0 , and it is well known that it can be bounded by [see, e.g., DeVore and Lorentz (1993)]

$$(6.3) \quad \begin{aligned} & \|Q_T(g_0) - g_0\|_{\mathcal{L}_2}^2 \\ & \leq c_0 \left(\max_{0 \leq j \leq m} |T_{j+1} - T_j|^{2r} \right) \int_T [g_0^{(r)}(t)]^2 dt \\ & \leq c_0 M_0 m^{-2r}. \end{aligned}$$

Hereafter, we shall use $c_0 > 0$ as a generic constant which may take different values at different appearance.

It now remains to bound $\|Q_T(h)\|_{\mathcal{L}_2}$. We appeal to the relationship between spline interpolation and the best local polynomial approximation. Let

$$I_j = [T_{j-r+1}, T_{j+r}]$$

with the convention that $T_j = 0$ for $j < 1$ and $T_j = 1$ for $j > m$. Denote by P_j the best approximation error that can be achieved on I_j by a polynomial of order less than r , that is,

$$(6.4) \quad P_j(f) = \min_{a_k : k < r} \int_{I_j} \left[\sum_{k=0}^{r-1} a_k t^k - f(t) \right]^2 dt.$$

It can be shown [see, e.g., Theorem 4.5 on page 147 of DeVore and Lorentz (1993)] that

$$\|f - Q_T(f)\|_{\mathcal{L}_2}^2 \leq c_0 \left(\sum_{j=1}^m P_j(f)^2 \right).$$

Then

$$\|Q_T(h)\|_{\mathcal{L}_2} \leq \|h\|_{\mathcal{L}_2} + c_0 \left(\sum_{j=1}^m P_j(f)^2 \right)^{1/2}.$$

Together with the fact that

$$P_j(h)^2 \leq \int_{I_j} h(t)^2 dt,$$

we have

$$\begin{aligned} \|\mathcal{Q}_T(h)\|_{\mathcal{L}_2}^2 &\leq c_0 \|h\|_{\mathcal{L}_2}^2 \\ &\leq c_0 \sum_{j=1}^m \delta_j^2 (T_{j+1} - T_{j-1}) \leq c_0 m^{-1} \sum_{j=1}^m \delta_j^2 \\ &= c_0 m^{-1} \sum_{j=1}^m [\hat{g}_\lambda(T_j) - g_0(T_j)]^2 \\ &\leq c_0 m^{-1} \sum_{j=1}^m ([\bar{Y}_{\cdot,j} - \hat{g}_\lambda(T_j)]^2 + [\bar{Y}_{\cdot,j} - g_0(T_j)]^2). \end{aligned}$$

Observe that

$$\mathbb{E} \left(m^{-1} \sum_{j=1}^m [\bar{Y}_{\cdot,j} - g_0(T_j)]^2 \right) = c_0 \sigma_0^2 n^{-1}.$$

It suffices to show that

$$m^{-1} \sum_{j=1}^m [\bar{Y}_{\cdot,j} - \hat{g}_\lambda(T_j)]^2 = O_p(m^{-2r} + n^{-1}).$$

To this end, note that by the definition of \hat{g}_λ ,

$$\begin{aligned} \frac{1}{m} \sum_{j=1}^m (\bar{Y}_{\cdot,j} - \hat{g}_\lambda(T_j))^2 &\leq \frac{1}{m} \sum_{j=1}^m (\bar{Y}_{\cdot,j} - \hat{g}_\lambda(T_j))^2 + \lambda \int_T [\hat{g}_\lambda^{(r)}(t)]^2 dt \\ &\leq \frac{1}{m} \sum_{j=1}^m (\bar{Y}_{\cdot,j} - g_0(T_j))^2 + \lambda \int_T [g_0^{(r)}(t)]^2 dt \\ &\leq O_p(m^{-2r} + n^{-1}), \end{aligned}$$

because $\lambda = O(m^{-2r} + n^{-1})$. The proof is now complete. \square

PROOF OF THEOREM 3.1. Note that any lower bound for a specific case yields immediately a lower bound for the general case. It therefore suffices to consider the case when X is a Gaussian process and $m_1 = m_2 = \dots = m_n =: m$. Denote by $N = c(nm)^{1/(2r+1)}$ where $c > 0$ is a constant to be specified later. Let $b = (b_1, \dots, b_N) \in \{0, 1\}^N$ be a binary sequence, and write

$$g_b(\cdot) = M_0^{1/2} \pi^{-r} \sum_{k=N+1}^{2N} N^{-1/2} k^{-r} b_{k-N} \varphi_k(\cdot),$$

where $\varphi_k(t) = \sqrt{2} \cos(\pi kt)$. It is not hard to see that

$$\begin{aligned} \int_{\mathcal{T}} [g_b^{(r)}(t)]^2 dt &= M_0 \pi^{-2r} \sum_{k \geq N+1}^{2N} (\pi k)^{2r} (N^{-1/2} k^{-r} b_{k-N})^2 \\ &= M_0 N^{-1} \sum_{k=N+1}^{2N} b_{k-N} \leq M_0. \end{aligned}$$

Furthermore,

$$\begin{aligned} \|g_b - g_{b'}\|_{\mathcal{L}_2}^2 &= M_0 \pi^{-2r} N^{-1} \sum_{k=N+1}^{2N} k^{-2r} (b_{k-N} - b'_{k-N})^2 \\ &\geq M_0 \pi^{-2r} (2N)^{-(2r+1)} \sum_{k=N+1}^{2N} (b_{k-m} - b'_{k-m})^2 \\ &= c_0 N^{-(2r+1)} H(b, b') \end{aligned}$$

for some constant $c_0 > 0$. By the Varshamov–Gilbert bound [see, e.g., Tsybakov (2009)], there exists a collection of binary sequences $\{b^{(1)}, \dots, b^{(M)}\} \subset \{0, 1\}^N$ such that $M \geq 2^{N/8}$, and

$$H(b^{(j)}, b^{(k)}) \geq N/8 \quad \forall 1 \leq j < k \leq M.$$

Then

$$\|g_{b^{(j)}} - g_{b^{(k)}}\|_{\mathcal{L}_2} \geq c_0 N^{-r}.$$

Assume that X is a Gaussian process with mean g_b , T follows a uniform distribution on \mathcal{T} and the measurement error $\varepsilon \sim N(0, \sigma_0^2)$. Conditional on $\{T_{ij} : j = 1, \dots, m\}$, $Z_i = (Z_{i1}, \dots, Z_{im})'$ follows a multivariate normal distribution with mean $\mu_b = (g_b(T_{i1}), \dots, g_b(T_{im}))'$ and covariance matrix $\Sigma(T) = (C_0(T_{ij}, T_{ik}))_{1 \leq j, k \leq m} + \sigma_0^2 I$. Therefore, the Kullback–Leibler distance from probability measure $\Pi_{g_{b^{(j)}}}$ to $\Pi_{g_{b^{(k)}}}$ can be bounded by

$$\begin{aligned} \text{KL}(\Pi_{g_{b^{(j)}}} | \Pi_{g_{b^{(k)}}}) &= n \mathbb{E}_T [(\mu_{g_{b^{(j)}}} - \mu_{g_{b^{(k)}}})' \Sigma^{-1}(T) (\mu_{g_{b^{(j)}}} - \mu_{g_{b^{(k)}}})] \\ &\leq n \sigma_0^{-2} \mathbb{E}_T \|\mu_{g_{b^{(j)}}} - \mu_{g_{b^{(k)}}}\|^2 \\ &= nm \sigma_0^{-2} \|g_{b^{(j)}} - g_{b^{(k)}}\|_{\mathcal{L}_2}^2 \\ &\leq c_1 nm \sigma_0^{-2} N^{-2r}. \end{aligned}$$

An application of Fano’s lemma now yields

$$\begin{aligned} \max_{1 \leq j \leq M} \mathbb{E}_{g_{b^{(j)}}} \|\tilde{g} - g_{b^{(j)}}\|_{\mathcal{L}_2} &\geq c_0 N^{-r} \left(1 - \frac{\log(c_1 nm \sigma_0^{-2} N^{-2r}) + \log 2}{\log M} \right) \\ &\asymp (nm)^{-r/(2r+1)} \end{aligned}$$

with an appropriate choice of c , for any estimate \tilde{g} . This in turn implies that

$$\limsup_{n \rightarrow \infty} \inf_{\tilde{g} \in \mathcal{D}} \sup_{\mathcal{L}(X) \in \mathcal{P}(r; M_0)} P(\|\tilde{g} - g_0\|_{\mathcal{L}_2}^2 > d(nm)^{-2r/(2r+1)}) > 0.$$

The proof can then be completed by considering X as an unknown constant function. \square

PROOF OF THEOREM 3.2. For brevity, in what follows, we treat the sampling frequencies m_1, \dots, m_n as deterministic. All the arguments, however, also apply to the situation when they are random by treating all the expectations and probabilities as conditional on m_1, \dots, m_n . Similarly, we shall also assume that T_j 's follow uniform distribution. The argument can be easily applied to handle more general distributions.

It is well known that \mathcal{W}_2^r , endowed with the norm

$$(6.5) \quad \|f\|_{\mathcal{W}_2^r}^2 = \int f^2 + \int (f^{(r)})^2,$$

forms a reproducing kernel Hilbert space [Aronszajn (1950)]. Let \mathcal{H}_0 be the collection of all polynomials of order less than r and \mathcal{H}_1 be its orthogonal complement in \mathcal{W}_2^r . Let $\{\phi_k : 1 \leq k \leq r\}$ be a set of orthonormal basis functions of \mathcal{H}_0 , and $\{\phi_k : k > r\}$ an orthonormal basis of \mathcal{H}_1 such that any $f \in \mathcal{W}_2^r$ admits the representation

$$f = \sum_{v \geq 1} f_v \phi_v.$$

Furthermore,

$$\|f\|_{\mathcal{L}_2}^2 = \sum_{v \geq 1} f_v^2 \quad \text{and} \quad \|f\|_{\mathcal{W}_2^r}^2 = \sum_{v \geq 1} (1 + \rho_v^{-1}) f_v^2,$$

where $\rho_1 = \dots = \rho_r = +\infty$ and $\rho_v \asymp v^{-2r}$.

Recall that

$$\hat{g} = \arg \min_{g \in \mathcal{H}(K)} \left\{ \ell_{mn}(g) + \lambda \int [g^{(r)}]^2 \right\},$$

where

$$\ell_{mn}(g) = \frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \sum_{j=1}^{m_i} (Y_{ij} - g(T_{ij}))^2.$$

For brevity, we shall abbreviate the subscript of \hat{g} hereafter when no confusion occurs. Write

$$\begin{aligned} \ell_\infty(g) &= \mathbb{E} \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \sum_{j=1}^{m_i} [Y_{ij} - g(T_{ij})]^2 \right) \\ &= \mathbb{E}([Y_{11} - g_0(T_{11})]^2) + \int_{\mathcal{T}} [g(s) - g_0(s)]^2 ds. \end{aligned}$$

Let

$$\bar{g} = \arg \min_{g \in \mathcal{H}(K)} \left\{ \ell_\infty(g) + \lambda \int [g^{(r)}]^2 \right\}.$$

Denote

$$\ell_{mn,\lambda}(g) = \ell_{mn}(g) + \lambda \int [g^{(r)}]^2; \quad \ell_{\infty,\lambda}(g) = \ell_\infty(g) + \lambda \int [g^{(r)}]^2.$$

Let

$$\tilde{g} = \bar{g} - \frac{1}{2} G_\lambda^{-1} D \ell_{mn,\lambda}(\bar{g}),$$

where $G_\lambda = (1/2) D^2 \ell_{\infty,\lambda}(\bar{g})$ and D stands for the Fréchet derivative. It is clear that

$$\hat{g} - g_0 = (\bar{g} - g_0) + (\hat{g} - \tilde{g}) + (\tilde{g} - \bar{g}).$$

We proceed by bounding the three terms on the right-hand side separately. In particular, it can be shown that

$$(6.6) \quad \|\bar{g} - g_0\|_{\mathcal{L}_2}^2 \leq c_0 \lambda \int [g_0^{(r)}]^2$$

and

$$(6.7) \quad \|\tilde{g} - \bar{g}\|_{\mathcal{L}_2}^2 = O_p(n^{-1} + (nm)^{-1} \lambda^{-1/(2r)}).$$

Furthermore, if

$$nm \lambda^{1/(2r)} \rightarrow \infty,$$

then

$$(6.8) \quad \|\hat{g} - \tilde{g}\|_{\mathcal{L}_2}^2 = o_p(n^{-1} + (nm)^{-1} \lambda^{-1/(2r)}).$$

Therefore,

$$\|\hat{g} - g_0\|_{\mathcal{L}_2}^2 = O_p(\lambda + n^{-1} + (nm)^{-1} \lambda^{-1/(2r)}).$$

Taking

$$\lambda \asymp (nm)^{-2r/(2r+1)}$$

yields

$$\|\hat{g} - g_0\|_{\mathcal{L}_2}^2 = O_p(n^{-1} + (nm)^{-2r/(2r+1)}).$$

We now set to establish bounds (6.6)–(6.8). For brevity, we shall assume in what follows that all expectations are taken conditionally on m_1, \dots, m_n unless otherwise indicated. Define

$$(6.9) \quad \|g\|_\alpha^2 = \sum_{v \geq 1} (1 + \rho_v^{-1})^\alpha g_v^2,$$

where $0 \leq \alpha \leq 1$.

We begin with $\bar{g} - g_0$. Write

$$(6.10) \quad g_0(\cdot) = \sum_{k \geq 1} a_k \phi_k(\cdot), \quad g(\cdot) = \sum_{k \geq 1} b_k \phi_k(\cdot).$$

Then

$$(6.11) \quad \ell_\infty(g) = \mathbb{E}([Y_{11} - g_0(T_{11})]^2) + \sum_{k \geq 1} (b_k - a_k)^2.$$

It is not hard to see

$$(6.12) \quad \bar{b}_k := \langle \bar{g}, \phi_k \rangle_{\mathcal{L}_2} = \arg \min\{(b_k - a_k)^2 + \lambda \rho_k^{-1} b_k^2\} = \frac{a_k}{1 + \lambda \rho_k^{-1}}.$$

Hence,

$$\begin{aligned} \|\bar{g} - g_0\|_\alpha^2 &= \sum_{k \geq 1} (1 + \rho_k^{-1})^\alpha (\bar{b}_k - a_k)^2 \\ &= \sum_{k \geq 1} (1 + \rho_k^{-1})^\alpha \left(\frac{\lambda \rho_k^{-1}}{1 + \lambda \rho_k^{-1}} \right)^2 a_k^2 \\ &\leq c_0 \lambda^2 \sup_{k \geq 1} \frac{\rho_k^{-(1+\alpha)}}{(1 + \lambda \rho_k^{-1})^2} \sum_{k=1}^\infty \rho_k^{-1} a_k^2 \\ &\leq c_0 \lambda^{1-\alpha} \int [g_0^{(r)}]^2. \end{aligned}$$

Next, we consider $\tilde{g} - \bar{g}$. Notice that $D\ell_{mn,\lambda}(\tilde{g}) = D\ell_{mn,\lambda}(\bar{g}) - D\ell_{\infty,\lambda}(\bar{g}) = D\ell_{mn}(\bar{g}) - D\ell_\infty(\bar{g})$. Therefore

$$\begin{aligned} \mathbb{E}[D\ell_{mn,\lambda}(\tilde{g})f]^2 &= \mathbb{E}[D\ell_{mn}(\bar{g})f - D\ell_\infty(\bar{g})f]^2 \\ &= \frac{4}{n^2} \sum_{i=1}^n \frac{1}{m_i^2} \text{Var} \left[\sum_{j=1}^{m_i} ([Y_{ij} - \bar{g}(T_{ij})]f(T_{ij})) \right]. \end{aligned}$$

Note that

$$\begin{aligned} &\text{Var} \left[\sum_{j=1}^{m_i} ([Y_{ij} - \bar{g}(T_{ij})]f(T_{ij})) \right] \\ &= \text{Var} \left[\mathbb{E} \left(\sum_{j=1}^{m_i} [Y_{ij} - \bar{g}(T_{ij})]f(T_{ij}) | T \right) \right] + \mathbb{E} \left[\text{Var} \left(\sum_{j=1}^{m_i} Y_{ij} f(T_{ij}) | T \right) \right] \\ &= \text{Var} \left[\sum_{j=1}^{m_i} ([g_0(T_{ij}) - \bar{g}(T_{ij})]f(T_{ij})) \right] + \mathbb{E} \left[\text{Var} \left(\sum_{j=1}^{m_i} Y_{ij} f(T_{ij}) | m_i, T \right) \right]. \end{aligned}$$

The first term on the rightmost-hand side can be bounded by

$$\begin{aligned} & \mathbb{V}\text{ar} \left[\sum_{j=1}^{m_i} ([g_0(T_{ij}) - \bar{g}(T_{ij})] f(T_{ij})) \right] \\ &= m_i \mathbb{V}\text{ar}([g_0(T_{i1}) - \bar{g}(T_{i1})] f(T_{i1})) \\ &\leq m_i \mathbb{E}([g_0(T_{i1}) - \bar{g}(T_{i1})] f(T_{i1}))^2 \\ &= m_i \int_{\mathcal{T}} ([g_0(t) - \bar{g}(t)] f(t))^2 dt \\ &\leq m_i \int_{\mathcal{T}} [g_0(t) - \bar{g}(t)]^2 dt \int_{\mathcal{T}} f^2(t) dt, \end{aligned}$$

where the last inequality follows from the Cauchy–Schwarz inequality. Together with (6.6), we get

$$(6.13) \quad \mathbb{V}\text{ar} \left[\sum_{j=1}^{m_i} ([g_0(T_{ij}) - \bar{g}(T_{ij})] f(T_{ij})) \right] \leq c_0 m_i \|f\|_{\mathcal{L}_2}^2 \lambda.$$

We now set out to compute the the second term. First observe that

$$\begin{aligned} & \mathbb{V}\text{ar} \left(\sum_{j=1}^{m_i} Y_{ij} f(T_{ij}) \middle| T \right) \\ &= \sum_{j,k=1}^{m_i} f(T_{ij}) f(T_{ik}) (C_0(T_{ij}, T_{ik}) + \sigma_0^2 \delta_{jk}), \end{aligned}$$

where δ_{jk} is Kronecker’s delta. Therefore,

$$\begin{aligned} & \mathbb{E} \left[\mathbb{V}\text{ar} \left(\sum_{j=1}^m Y_{1j} f(T_{1j}) \middle| T \right) \middle| m_i \right] \\ &= m_i (m_i - 1) \int_{\mathcal{T} \times \mathcal{T}} f(s) C_0(s, t) f(t) ds dt \\ &\quad + m_i \sigma_0^2 \|f\|_{\mathcal{L}_2}^2 + m_i \int_{\mathcal{T}} f^2(s) C(s, s) ds. \end{aligned}$$

Summing up, we have

$$(6.14) \quad \mathbb{E}[D\ell_{mn,\lambda}(\bar{g})\phi_k]^2 \leq \frac{c_0}{n^2} \left(\sum_{i=1}^n \frac{1}{m_i} \right) + \frac{4c_k}{n},$$

where

$$(6.15) \quad c_k = \int_{\mathcal{T} \times \mathcal{T}} \phi_k(s) C_0(s, t) \phi_k(t) ds dt.$$

Therefore,

$$\begin{aligned} \mathbb{E}\|\tilde{g} - \bar{g}\|_\alpha^2 &= \mathbb{E}\left\|\frac{1}{2}G_\lambda^{-1}D\ell_{nm,\lambda}(\bar{g})\right\|_\alpha^2 \\ &= \frac{1}{4}\mathbb{E}\left[\sum_{k\geq 1}(1 + \rho_k^{-1})^\alpha(1 + \lambda\rho_k^{-1})^{-2}(D\ell_{nm,\lambda}(\bar{g})\phi_k)^2\right] \\ &\leq \frac{c_0}{n^2}\left(\sum_{i=1}^n\frac{1}{m_i}\right)\sum_{k\geq 1}(1 + \rho_k^{-1})^\alpha(1 + \lambda\rho_k^{-1})^{-2} \\ &\quad + \frac{1}{n}\sum_{k\geq 1}(1 + \rho_k^{-1})^\alpha(1 + \lambda\rho_k^{-1})^{-2}c_k. \end{aligned}$$

Observe that

$$\sum_{k\geq 1}(1 + \rho_k^{-1})^\alpha(1 + \lambda\rho_k^{-1})^{-2} \leq c_0\lambda^{-\alpha-1/(2r)}$$

and

$$\sum_{k\geq 1}(1 + \rho_k^{-1})^\alpha(1 + \lambda\rho_k^{-1})^{-2}c_k \leq \sum_{k\geq 1}(1 + \rho_k^{-1})c_k = \mathbb{E}\|X\|_{\mathcal{W}_2^r}^2 < \infty.$$

Thus,

$$\mathbb{E}\|\tilde{g} - \bar{g}\|_\alpha^2 \leq c_0\left[\frac{1}{n^2}\left(\sum_{i=1}^n\frac{1}{m_i}\right)\lambda^{-\alpha-1/(2r)} + \frac{1}{n}\right].$$

It remains to bound $\hat{g} - \tilde{g}$. It can be easily verified that

$$(6.16) \quad \hat{g} - \tilde{g} = \frac{1}{2}G_\lambda^{-1}[D^2\ell_\infty(\bar{g})(\hat{g} - \bar{g}) - D^2\ell_{mn}(\bar{g})(\hat{g} - \bar{g})].$$

Then

$$\begin{aligned} \|\hat{g} - \tilde{g}\|_\alpha^2 &= \sum_{k\geq 1}(1 + \rho_k^{-1})^\alpha(1 + \lambda\rho_k^{-1})^{-2} \\ &\quad \times \left[\frac{1}{n}\sum_{i=1}^n\frac{1}{m_i}\sum_{j=1}^{m_i}(\hat{g}(T_{ij}) - \bar{g}(T_{ij}))\phi_k(T_{ij})\right. \\ &\quad \left. - \int_{\mathcal{T}}(\hat{g}(s) - \bar{g}(s))\phi_k(s) ds\right]^2. \end{aligned}$$

Clearly, $(\hat{g} - \bar{g})\phi_k \in \mathcal{H}(K)$. Write

$$(6.17) \quad (\hat{g} - \bar{g})\phi_k = \sum_{j\geq 1}h_j\phi_j.$$

Then, by the Cauchy–Schwarz inequality,

$$\begin{aligned} & \left[\frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \sum_{j=1}^{m_i} (\hat{g}(T_{ij}) - \bar{g}(T_{ij})) \phi_k(T_{ij}) - \int_{\mathcal{T}} (\hat{g}(s) - \bar{g}(s)) \phi_k(s) ds \right]^2 \\ &= \left[\sum_{k_1 \geq 1} h_{k_1} \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \sum_{j=1}^{m_i} \phi_{k_1}(T_{ij}) - \int_{\mathcal{T}} \phi_{k_1}(s) ds \right) \right]^2 \\ &\leq \left[\sum_{k_1 \geq 1} (1 + \rho_{k_1}^{-1})^\gamma h_{k_1}^2 \right] \\ &\quad \times \left[\sum_{k_1 \geq 1} (1 + \rho_{k_1}^{-1})^{-\gamma} \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \sum_{j=1}^{m_i} \phi_{k_1}(T_{ij}) - \int_{\mathcal{T}} \phi_{k_1}(s) ds \right)^2 \right] \\ &\leq \|\hat{g} - \bar{g}\|_\gamma^2 (1 + \rho_k^{-1})^\gamma \\ &\quad \times \left[\sum_{k_1 \geq 1} (1 + \rho_{k_1}^{-1})^{-\gamma} \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \sum_{j=1}^{m_i} \phi_{k_1}(T_{ij}) - \int_{\mathcal{T}} \phi_{k_1}(s) ds \right)^2 \right] \end{aligned}$$

for any $0 \leq \gamma \leq 1$, where in the last inequality, we used the fact that

$$(6.18) \quad \|(\hat{g} - \bar{g})\phi_k\|_\gamma \leq \|\hat{g} - \bar{g}\|_\gamma \|\phi_k\|_\gamma = \|\hat{g} - \bar{g}\|_\gamma (1 + \rho_k^{-1})^{\gamma/2}.$$

Following a similar calculation as before, it can be shown that

$$\begin{aligned} & \mathbb{E} \left[\sum_{k_1 \geq 1} (1 + \rho_{k_1}^{-1})^{-\gamma} \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \sum_{j=1}^{m_i} \phi_{k_1}(T_{ij}) - \int_{\mathcal{T}} \phi_{k_1}(s) ds \right)^2 \right] \\ & \leq \frac{1}{n^2} \left(\sum_{i=1}^n \frac{1}{m_i} \right) \sum_{k_1 \geq 1} (1 + \rho_{k_1}^{-1})^{-\gamma}, \end{aligned}$$

which is finite whenever $\gamma > 1/2r$. Recall that m is the harmonic mean of m_1, \dots, m_n . Therefore,

$$(6.19) \quad \|\hat{g} - \bar{g}\|_\alpha^2 \leq O_p \left(\frac{1}{nm\lambda^{\alpha+\gamma+1/(2r)}} \right) \|\hat{g} - \bar{g}\|_\gamma^2.$$

If $\alpha > 1/2r$, then taking $\gamma = \alpha$ yields

$$(6.20) \quad \|\hat{g} - \bar{g}\|_\alpha^2 = O_p \left(\frac{1}{nm\lambda^{2\alpha+1/(2r)}} \right) \|\hat{g} - \bar{g}\|_\alpha^2 = o_p(\|\hat{g} - \bar{g}\|_\alpha),$$

assuming that

$$(6.21) \quad nm\lambda^{2\alpha+1/(2r)} \rightarrow \infty.$$

Together with the triangular inequality

$$(6.22) \quad \|\tilde{g} - \bar{g}\|_\alpha \geq \|\hat{g} - \bar{g}\|_\alpha - \|\hat{g} - \tilde{g}\|_\alpha = (1 - o_p(1))\|\hat{g} - \bar{g}\|_\alpha.$$

Therefore,

$$(6.23) \quad \|\hat{g} - \bar{g}\|_\alpha^2 = O_p(\|\tilde{g} - \bar{g}\|_\alpha^2) = O_p(n^{-1} + (nm)^{-1}\lambda^{-\alpha-1/(2r)}).$$

Together with (6.19),

$$\begin{aligned} \|\hat{g} - \tilde{g}\|_{\mathcal{L}_2}^2 &= O_p\left(\frac{1}{nm\lambda^{\alpha+1/(2r)}}(n^{-1} + (nm)^{-1}\lambda^{-\alpha-1/(2r)})\right) \\ &= o_p(n^{-1}\lambda^\alpha + (nm)^{-1}\lambda^{-1/(2r)}). \end{aligned}$$

We conclude by noting that in the case when m_1, \dots, m_n are random, m can also be replaced with the expectation of the harmonic mean thanks to the law of large numbers. \square

PROOF OF PROPOSITION 3.3. Let $Q_{T,\lambda}$ be the smoothing spline operator, that is, $Q_{T,\lambda}(f_1, \dots, f_m)$ is the solution to

$$\min_{g \in \mathcal{W}_2^r} \left\{ \frac{1}{m} \sum_{j=1}^m (f_j - g(T_j))^2 + \lambda \int_T [g^{(r)}(t)]^2 dt \right\}.$$

It is clear that

$$\hat{g}_\lambda = Q_{T,\lambda}(\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_m)$$

and

$$\tilde{X}_{i,\lambda} = Q_{T,\lambda}(Y_{i1}, Y_{i2}, \dots, Y_{im}).$$

Because $Q_{T,\lambda}$ is a linear operator [see, e.g., Wahba (1990)], we have

$$\begin{aligned} \tilde{g}_\lambda &= \frac{1}{n} \sum_{i=1}^n \tilde{X}_{i,\lambda} = \frac{1}{n} \sum_{i=1}^n Q_{T,\lambda}(Y_{i1}, Y_{i2}, \dots, Y_{im}) \\ &= Q_{T,\lambda}(\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_m) = \hat{g}_\lambda. \end{aligned} \quad \square$$

Acknowledgments. We thank an Associate Editor and two referees for their constructive comments which have helped to improve the presentation of the paper.

REFERENCES

ARONSZAJN, N. (1950). Theory of reproducing kernels. *Trans. Amer. Math. Soc.* **68** 337–404. MR0051437
 ASSOUD, P. (1983). Deux remarques sur l'estimation. *C. R. Acad. Sci. Paris Sér. I Math.* **296** 1021–1024. MR0777600

- CAI, T. and YUAN, M. (2010). Nonparametric covariance function estimation for functional and longitudinal data. Technical report, Georgia Institute of Technology, Atlanta, GA.
- DEVORE, R. A. and LORENTZ, G. G. (1993). *Constructive Approximation. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]* **303**. Springer, Berlin. [MR1261635](#)
- DIGGLE, P., HEAGERTY, P., LIANG, K. and ZEGER, S. (2002). *Analysis of Longitudinal Data*, 2nd ed. Oxford Univ. Press, Oxford. [MR2049007](#)
- FERRATY, F. and VIEU, P. (2006). *Nonparametric Functional Data Analysis: Theory and Practice*. Springer, New York. [MR2229687](#)
- GREEN, P. J. and SILVERMAN, B. W. (1994). *Nonparametric Regression and Generalized Linear Models: A Roughness Penalty Approach. Monographs on Statistics and Applied Probability* **58**. Chapman and Hall, London. [MR1270012](#)
- HALL, P. and HART, J. D. (1990). Nonparametric regression with long-range dependence. *Stochastic Process. Appl.* **36** 339–351. [MR1084984](#)
- HALL, P., MÜLLER, H.-G. and WANG, J.-L. (2006). Properties of principal component methods for functional and longitudinal data analysis. *Ann. Statist.* **34** 1493–1517. [MR2278365](#)
- HART, J. D. and WEHRLY, T. E. (1986). Kernel regression estimation using repeated measurements data. *J. Amer. Statist. Assoc.* **81** 1080–1088. [MR0867635](#)
- JAMES, G. M. and HASTIE, T. J. (2001). Functional linear discriminant analysis for irregularly sampled curves. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **63** 533–550. [MR1858401](#)
- JOHNSTONE, I. M. and SILVERMAN, B. W. (1997). Wavelet threshold estimators for data with correlated noise. *J. Roy. Statist. Soc. Ser. B* **59** 319–351. [MR1440585](#)
- OPSOMER, J., WANG, Y. and YANG, Y. (2001). Nonparametric regression with correlated errors. *Statist. Sci.* **16** 134–153. [MR1861070](#)
- RAMSAY, J. O. and SILVERMAN, B. W. (2002). *Applied Functional Data Analysis: Methods and Case Studies*. Springer, New York. [MR1910407](#)
- RAMSAY, J. O. and SILVERMAN, B. W. (2005). *Functional Data Analysis*, 2nd ed. Springer, New York. [MR2168993](#)
- RICE, J. A. and SILVERMAN, B. W. (1991). Estimating the mean and covariance structure nonparametrically when the data are curves. *J. Roy. Statist. Soc. Ser. B* **53** 233–243. [MR1094283](#)
- RICE, J. A. and WU, C. O. (2001). Nonparametric mixed effects models for unequally sampled noisy curves. *Biometrics* **57** 253–259. [MR1833314](#)
- STONE, C. J. (1982). Optimal global rates of convergence for nonparametric regression. *Ann. Statist.* **10** 1040–1053. [MR0673642](#)
- TSYBAKOV, A. B. (2009). *Introduction to Nonparametric Estimation*. Springer, New York. [MR2724359](#)
- WAHBA, G. (1990). *Spline Models for Observational Data. CBMS-NSF Regional Conference Series in Applied Mathematics* **59**. SIAM, Philadelphia, PA. [MR1045442](#)
- WANG, Y. (1996). Function estimation via wavelet shrinkage for long-memory data. *Ann. Statist.* **24** 466–484. [MR1394972](#)
- YAO, F., MÜLLER, H.-G. and WANG, J.-L. (2005). Functional data analysis for sparse longitudinal data. *J. Amer. Statist. Assoc.* **100** 577–590. [MR2160561](#)

DEPARTMENT OF STATISTICS
THE WHARTON SCHOOL
UNIVERSITY OF PENNSYLVANIA
PHILADELPHIA, PENNSYLVANIA 19104
USA
E-MAIL: tcai@wharton.upenn.edu

SCHOOL OF INDUSTRIAL AND SYSTEMS ENGINEERING
GEORGIA INSTITUTE OF TECHNOLOGY
ATLANTA, GEORGIA 30332
USA
E-MAIL: myuan@isye.gatech.edu