

We propose a nonparametric method to perform functional principal components analysis for the case of sparse longitudinal data. The method aims at irregularly spaced longitudinal data, where the number of repeated measurements available per subject is small. In contrast, classical functional data analysis requires a large number of regularly spaced measurements per subject. We assume that the repeated measurements are located randomly with a random number of repetitions for each subject and are determined by an underlying smooth random (subject-specific) trajectory plus measurement errors. Basic elements of our approach are the parsimonious estimation of the covariance structure and mean function of the trajectories, and the estimation of the variance of the measurement errors. The eigenfunction basis is estimated from the data, and functional principal components score estimates are obtained by a conditioning step. This conditional estimation method is conceptually simple and straightforward to implement. A key step is the derivation of asymptotic consistency and distribution results under mild conditions, using tools from functional analysis. Functional data analysis for sparse longitudinal data enables prediction of individual smooth trajectories even if only one or few measurements are available for a subject. Asymptotic pointwise and simultaneous confidence bands are obtained for predicted individual trajectories, based on asymptotic distributions, for simultaneous bands under the assumption of a finite number of components. Model selection techniques, such as the Akaike information criterion, are used to choose the model dimension corresponding to the number of eigenfunctions in the model. The methods are illustrated with a simulation study, longitudinal CD4 data for a sample of AIDS patients, and time-course gene expression data for the yeast cell cycle.

KEY WORDS: Asymptotics; Conditioning; Confidence band; Measurement error; Principal components; Simultaneous inference; Smoothing.

1. INTRODUCTION

We develop a version of functional principal components (FPC) analysis, in which the FPC scores are framed as conditional expectations. We demonstrate that this extends the applicability of FPC analysis to situations in longitudinal data analysis, where only few repeated and sufficiently irregularly spaced measurements are available per subject, and refer to this approach as principal components analysis through conditional expectation (PACE) for longitudinal data.

When the observed data are in the form of random curves rather than scalars or vectors, dimension reduction is mandatory, and FPC analysis has become a common tool to achieve this, by reducing random trajectories to a set of FPC scores. However, this method encounters difficulties when applied to longitudinal data with only few repeated observations per subject.

Beyond dimension reduction, FPC analysis attempts to characterize the dominant modes of variation of a sample of random trajectories around an overall mean trend function. There exists an extensive literature on FPC analysis when individuals are measured at a dense grid of regularly spaced time points. The method was introduced by Rao (1958) for growth curves, and the basic principle has been studied by Besse and Ramsay (1986), Castro, Lawton, and Sylvestre (1986), and Berkey, Laird, Valadian, and Gardner (1991). Rice and Silverman (1991) discussed smoothing and smoothing parameter choice in this context, whereas Jones and Rice (1992) emphasized applications. Various theoretical properties have been studied by Silverman (1996), Boente and Fraiman (2000), and Kneip and Utikal (2001). (For an introduction and summary, see

Ramsay and Silverman 1997.) Staniswalis and Lee (1998) proposed kernel-based functional principal components analysis for repeated measurements with an irregular grid of time points. The case of irregular grids was also studied by Besse, Cardot, and Ferraty (1997) and Boularan, Ferré, and Vieu (1993). However, when the time points vary widely across subjects and are sparse, down to one or two measurements, the FPC scores defined through the Karhunen–Loève expansion are not well approximated by the usual integration method.

Shi, Weiss, and Taylor (1996), Rice and Wu (2000), James, Hastie, and Sugar (2001), and James and Sugar (2003) proposed B-splines to model the individual curves with random coefficients through mixed effects models. James et al. (2001) and James and Sugar (2003) emphasized the case of sparse data, postulating a reduced-rank mixed-effects model through a B-spline basis for the underlying random trajectories. In contrast, we represent the trajectories directly through the Karhunen–Loève expansion, determining the eigenfunctions from the data. Perhaps owing to the complexity of their modeling approach, James et al. (2001) did not investigate the asymptotic properties of the estimated components in relation to the true components, such as the behavior of the estimated covariance structure, eigenvalues, and eigenfunctions, especially for the sparse situation. Instead, they constructed pointwise confidence intervals for the individual curves using bootstrap. With our simpler and more direct approach, we are able to derive asymptotic properties, using tools from functional analysis. We can also derive both pointwise and simultaneous bands for predicted individual trajectories. This requires first obtaining the uniform convergence results for nonparametric function and surface estimates under dependence structure that follows from the longitudinal nature of the data. The dependence is a consequence of the assumed random nature of the observed sample of trajectories, which sets our work apart from previous results

Fang Yao is Assistant Professor, Department of Statistics, Colorado State University, Fort Collins, CO 80523 (E-mail: yao@stat.colostate.edu). Hans-Georg Müller is Professor (E-mail: muller@math.ucdavis.edu) and Jane-Ling Wang is Professor (E-mail: jwang@math.ucdavis.edu), Department of Statistics, University of California, Davis, CA 95616. This research was supported in part by National Science Foundation grants DMS-98-03637, DMS-99-71602, DMS-02-04869, DMS-03-54448, and DMS-04-06430. The authors thank the associate editor and two referees for insightful comments on a previous version of this article that led to many improvements.

where either the observed functions are nonrandom with independent measurements (Kneip 1994), are random vectors of large but fixed dimensions (Ferré 1995), or are random trajectories sampled on dense and regular grids (Cardot, Ferraty, and Sarda 1999).

The contributions of this article are as follows. First, we provide a new technique, PACE, for longitudinal and functional data, a method designed to handle sparse and irregular longitudinal data for which the pooled time points are sufficiently dense. Second, the presence of additional measurement errors is taken into account, extending previous approaches of Staniswalis and Lee (1998) and Yao et al. (2003). Third, an emphasis is on the derivation of asymptotic consistency properties, by first establishing uniform convergence for smoothed estimates of the mean and covariance functions under mild assumptions. These uniform consistency results are developed for smoothers in the situation where repeated, and thus dependent, measurements are obtained for the same subject. Then we couple these results with the theory of eigenfunctions and eigenvalues of compact linear operators, to obtain uniform convergence of estimated eigenfunctions and eigenvalues. To our knowledge, there exist only few published asymptotic results for FPC (Dauxois, Pousse, and Romain 1982; Bosq 1991; Silverman 1996), and none for functional data analysis in the sparse situation. Fourth, we derive the asymptotic distribution needed to obtain pointwise confidence intervals for individual trajectories, and obtain asymptotic simultaneous bands for these trajectories.

The main novelty of our work is that we establish the conditional method for the case of sparse and irregular data, show that this provides a straightforward and simple tool for the modeling of longitudinal data, and derive asymptotic results for this method. Under Gaussian assumptions, the proposed estimation of individual FPC scores in PACE corresponds to the best prediction, combining the data from the individual subject to be predicted with data from the entire collection of subjects. In the non-Gaussian case, it provides an estimate for the best linear prediction. The proposed PACE method extends to the case of sparse and irregular data, provided that as the number of subjects increases, the pooled time points from the entire sample become dense in the domain of the data. We suggest one-curve-leave-out cross-validation for choosing auxiliary parameters, such as the degree of smoothing and the model dimension, corresponding to the number of eigenfunctions to be included, similar to the approach of Rice and Silverman (1991). For faster computing, we also consider the Akaike information criterion (AIC) to select the number of eigenfunctions.

The remainder of the article is organized as follows. In Section 2 we introduce the PACE approach, that is, the proposed conditional estimates for the FPC scores. We present asymptotic results for the proposed method in Section 3, with proofs deferred to the Appendix. We discuss simulation results that illustrate the usefulness of the methodology in Section 4. Applications of PACE to longitudinal CD4 data and time-course gene expression data for yeast cell cycle genes are the theme of Section 5, followed by concluding remarks in Section 6 and proofs and theoretical results in the Appendix.

2. FUNCTIONAL PRINCIPAL COMPONENT ANALYSIS FOR PEAR DATA

2.1 MODELING MEASUREMENT ERRORS

We model sparse functional data as noisy sampled points from a collection of trajectories that are assumed to be independent realizations of a smooth random function, with unknown mean function $\mu(\cdot) = \mu(\cdot)$ and covariance function $\text{cov}(\cdot, \cdot) = \text{cov}(\cdot, \cdot)$. The domain of $\cdot(\cdot)$ typically is a bounded and closed time interval \mathcal{I} . Although we refer to the index variable as time, it could also be a spatial variable, such as in image or geoscience applications. We assume that there is an orthogonal expansion (in the L^2 sense) of $\mu(\cdot)$ in terms of eigenfunctions $\phi_k(\cdot)$ and nonincreasing eigenvalues λ_k : $\mu(\cdot) = \sum_{k=1}^{\infty} \lambda_k \phi_k(\cdot)$, $\lambda_k \geq 0$. In classical FPC analysis, it is assumed that the i th random curve can be expressed as $\mu(\cdot) + \sum_{k=1}^{\infty} \lambda_k^{-1/2} \epsilon_{ik} \phi_k(\cdot)$, $\epsilon_{ik} \in \mathbb{R}$, where the ϵ_{ik} are uncorrelated random variables with mean 0 and variance $\sigma^2 = \lambda_k$, where $\sum_{k=1}^{\infty} \lambda_k < \infty$, $\lambda_1 \geq \lambda_2 \geq \dots$.

We consider an extended version of the model that incorporates uncorrelated measurement errors with mean 0 and constant variance σ^2 to reflect additive measurement errors (see also Rice and Wu 2000). Let y_{ij} be the i th observation of the random function $\mu(\cdot)$, made at a random time t_{ij} , and let ϵ_{ij} be the additional measurement errors that are assumed to be iid and independent of the random coefficients ϵ_{ik} , where $i = 1, \dots, n$, $j = 1, \dots, m$, $k = 1, 2, \dots$. Then the model that we consider is

$$y_{ij} = \mu(t_{ij}) + \sum_{k=1}^{\infty} \lambda_k^{-1/2} \epsilon_{ik} \phi_k(t_{ij}) + \epsilon_{ij}, \quad \epsilon_{ij} \in \mathbb{R}, \quad (1)$$

where $\epsilon_{ij} = 0$, $\text{var}(\epsilon_{ij}) = \sigma^2$, and the number of measurements made on the i th subject is considered random, reflecting sparse and irregular designs. The random variables ϵ_{ij} are assumed to be iid and independent of all other random variables.

2.2 ESTIMATING MEAN, COVARIANCE, AND EIGENFUNCTIONS

We assume that mean, covariance, and eigenfunctions are smooth. We use local linear smoothers (Fan and Gijbels 1996) for function and surface estimation, fitting local lines in one dimension and local planes in two dimensions by weighted least squares. In a first step, we estimate the mean function μ based on the pooled data from all individuals. The formula for this local linear smoother is in (A.1) in the Appendix. Data-adaptive methods for bandwidth choice are available (see Müller and Prewitt 1993 for surface smoothing and Rice and Silverman 1991 for one-curve-leave-out cross-validation); subjective choices are often adequate. (For issues of smoothing dependent data, see Lin and Carroll 2000.) Adapting to estimated correlations when estimating the mean function did not lead to improvements (simulations not reported); therefore, we do not incorporate such adjustments.

Note that in model (1), $\text{cov}(y_i | t_i, y_j | t_j) = \text{cov}(\mu(t_i), \mu(t_j)) + \sigma^2 \delta_{ij}$, where δ_{ij} is 1 if $i = j$ and 0 otherwise. Let $\tilde{C}(t_i, t_j) = (y_i - \hat{\mu}(t_i))(y_j - \hat{\mu}(t_j))$ be the “raw” covariances, where $\hat{\mu}(\cdot)$ is the estimated mean function obtained from the previous step. It is easy to see that $[\tilde{C}(t_i, t_j)]_{i,j} \approx \text{cov}(\mu(t_i), \mu(t_j)) + \sigma^2 \delta_{ij}$. Therefore, the diagonal of the raw

covariances should be removed; that is, only (t, t) , $t \neq s$, should be included as input data for the covariance surface smoothing step (as previously observed in Staniswalis and Lee 1998). We use one-curve-leave-out cross-validation to choose the smoothing parameter for this surface smoothing step.

The variance $\hat{\sigma}^2$ of the measurement errors is of interest in model (1). Let $\hat{f}(t, s)$ be a smooth surface estimate [see (A.2) in the App.] of $\text{cov}(Y(t), Y(s))$. Following Yao et al. (2003), because the covariance of $Y(t)$ is maximal along the diagonal, a local quadratic rather than a local linear fit is expected to better approximate the shape of the surface in the direction orthogonal to the diagonal. We thus fit a local quadratic component along the direction perpendicular to the diagonal and a local linear component in the direction of the diagonal; implementation of this local smoother is achieved by rotating the coordinates by 45 degrees and then minimizing weighted least squares [similar to (A.2)] in rotated coordinates with local quadratic and linear components, see (A.3) in the Appendix.

Denote the diagonal of the resulting surface estimate by $\tilde{f}(t)$ and a local linear smoother focusing on diagonal values $\{Y(t) + \tilde{f}(t)\}$ by $\hat{f}(t)$, obtained by (A.1) with $\{Y(t), \tilde{f}(t)\}$ as input. To mitigate boundary effects, we cut off the two ends of the interval to get a more stable estimate, following a suggestion of Staniswalis and Lee (1998). Let $|I|$ denote the length of I , and let I_1 be the interval $I_1 = [\inf\{t \in I : t \geq |I|/4\}, \sup\{t \in I : t \leq |I|/4\}]$. The proposed estimate of $\hat{\sigma}^2$ is

$$\hat{\sigma}^2 = \frac{2}{|I|} \int_{I_1} \{\hat{f}(t) - \tilde{f}(t)\}^2 dt \quad (2)$$

if $\hat{\sigma}^2 > 0$ and $\hat{\sigma}^2 = 0$ otherwise.

The estimates of eigenfunctions and eigenvalues correspond to the solutions \hat{f}_k and $\hat{\lambda}_k$ of the eigenequations,

$$\int_{\mathcal{T}} \hat{f}_k(t, s) \hat{f}_k(t) dt = \hat{\lambda}_k \hat{f}_k(t), \quad (3)$$

where the \hat{f}_k are subject to $\int_{\mathcal{T}} \hat{f}_k(t)^2 dt = 1$ and $\int_{\mathcal{T}} \hat{f}_k(t) \times \hat{f}_l(t) dt = 0$ for $k \neq l$. We estimate the eigenfunctions by discretizing the smoothed covariance, as previously described by Rice and Silverman (1991) and Capra and Müller (1997).

2.3 Functional Principal Component Analysis

The FPC scores $\alpha_k = \int_{\mathcal{T}} (Y(t) - \mu(t)) \hat{f}_k(t) dt$ have traditionally been estimated by numerical integration, which works well when the density of the grid of measurements for each subject is sufficiently large. Because in our model the $Y(t)$ are available only at discrete random times t_i , reflecting the sparseness of the data, the integrals in the definition of the FPC scores accordingly would be approximated by sums, substituting $Y(t_i)$ as defined in (1) for $Y(t)$ and estimates $\hat{\mu}(t)$ for $\mu(t)$ and $\hat{f}_k(t)$ for $f_k(t)$, leading to $\hat{\alpha}_k = \sum_{i=1}^n (Y(t_i) - \hat{\mu}(t_i)) \hat{f}_k(t_i) (t_i - t_{i-1})$, setting $t_0 = 0$. For sparse functional data, $\hat{\alpha}_k$ will not provide reasonable approximations to α_k , for example, when one has only two observations per subject. Moreover, when the measurements are contaminated

with errors, the underlying random process cannot be directly observed. Substituting $Y(t_i)$ for $Y(t)$ then leads to biased FPC scores. These considerations motivate the alternative PACE method to obtain the FPC scores.

Assume that in (1), $Y(t)$ and $\mu(t)$ are jointly Gaussian. In all of what follows, the results pertaining to expectations are always conditional on the observation times $t_i, i = 1, \dots, n$, $i = 1, \dots, n$. For simplicity, the dependence on t_i is suppressed. Write $\tilde{\mathbf{X}} = (X(t_1), \dots, X(t_n))$, $\tilde{\mathbf{Y}} = (Y(t_1), \dots, Y(t_n))$, $\boldsymbol{\mu} = (\mu(t_1), \dots, \mu(t_n))$, and $\boldsymbol{\nu} = (\nu(t_1), \dots, \nu(t_n))$. The best prediction of the FPC scores for the k th subject, given the data from that individual, is the conditional expectation, which, under Gaussian assumptions [also given in (A5) in Sec. 3], is found to be (see, e.g., thm. 3.2.4 in Mardia, Kent, and Bibby 1979)

$$\hat{\alpha}_k = [\boldsymbol{\nu} | \tilde{\mathbf{Y}}] = \boldsymbol{\nu} + \boldsymbol{\Sigma}_Y^{-1} (\tilde{\mathbf{Y}} - \boldsymbol{\mu}), \quad (4)$$

where $\boldsymbol{\Sigma}_Y = \text{cov}(\tilde{\mathbf{Y}}, \tilde{\mathbf{Y}}) = \text{cov}(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}) + \boldsymbol{\nu} \boldsymbol{\nu}^T$; that is, the (i, j) entry of the $n \times n$ matrix $\boldsymbol{\Sigma}_Y$ is $(\boldsymbol{\Sigma}_Y)_{ij}$

the linear functions of $\tilde{\mathbf{Y}}$, we have $[\tilde{\mu}(\cdot), \tilde{\sigma}(\cdot)] = [\tilde{\mu}(\cdot), \tilde{\sigma}(\cdot)]$, that is, $\text{var}(\tilde{\mu}(\cdot), \tilde{\sigma}(\cdot)) = \text{var}(\mu(\cdot), \sigma(\cdot)) - \text{var}(\hat{\mu}(\cdot), \hat{\sigma}(\cdot)) = \Sigma$, where $\Sigma = \Lambda - \mathbf{H}\Sigma_{\tilde{\mathbf{Y}}}^{-1}\mathbf{H}$ and $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_p\}$. Then, under Gaussian assumptions, $(\tilde{\mu}(\cdot), \tilde{\sigma}(\cdot)) \sim \mathcal{N}(\mathbf{0}, \Sigma)$.

We construct asymptotic pointwise confidence intervals for individual trajectories as follows. Let $\hat{\mu}(\cdot) = \hat{\Lambda} - \hat{\mathbf{H}}\hat{\Sigma}_{\tilde{\mathbf{Y}}}^{-1}\hat{\mathbf{H}}$, where $\hat{\Lambda} = \text{diag}\{\hat{\lambda}_1, \dots, \hat{\lambda}_p\}$ and $\hat{\mathbf{H}} = (\hat{\mu}_1(\cdot), \dots, \hat{\mu}_p(\cdot))$. For $t \in \mathcal{T}$, let $\hat{\mu}_1(t) = (\hat{\mu}_1(t), \dots, \hat{\mu}_p(t))$, $\hat{\sigma}_1(t) = (\hat{\sigma}_1(t), \dots, \hat{\sigma}_p(t))$, and $\hat{\mu}(t) = \hat{\mu}(\cdot) + \hat{\sigma}_1(t)$. Theorem 4 establishes that the distribution of $\{\hat{\mu}(t) - \mu(t)\}$ may be asymptotically approximated by $\mathcal{N}(\mathbf{0}, \hat{\Sigma}_1(t))$. Because we assume that $\mu(\cdot)$ can be approximated sufficiently well by the first eigenfunctions, we may construct the $(1 - \alpha)$ asymptotic pointwise interval for $\mu(t)$,

$$\hat{\mu}(t) \pm \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \sqrt{\hat{\Sigma}_1(t)}, \tag{7}$$

where Φ is the standard Gaussian cdf. These confidence intervals are constructed by ignoring the bias that results from the truncation at $\hat{\mu}$.

Next, consider the construction of asymptotic simultaneous confidence bands. Let $\hat{\mu}(t) = \mu(t) + \sum_{j=1}^p \hat{\sigma}_j(t)$. Theorem 5 provides the asymptotic simultaneous band for $\{\hat{\mu}(t) - \mu(t)\}$, for a given fixed t . The Karhunen–Loève theorem implies that $\sup_{t \in \mathcal{T}} [\hat{\mu}(t) - \mu(t)]^2$ is small for fixed t and sufficiently large n . Therefore, ignoring a remaining approximation error that may be interpreted as a bias, we may construct $(1 - \alpha)$ asymptotic simultaneous bands for $\mu(t)$ through

$$\hat{\mu}(t) \pm \sqrt{\chi^2_{2, 1-\alpha} \hat{\Sigma}_1(t)}, \tag{8}$$

where $\chi^2_{2, 1-\alpha}$ is the $100(1 - \alpha)$ th percentile of the chi-squared distribution with 2 degrees of freedom. Because $\sqrt{\chi^2_{2, 1-\alpha}} > \Phi^{-1}(1 - \alpha/2)$ for all $\alpha \geq 1$, the asymptotic simultaneous band is always wider than the corresponding asymptotic pointwise confidence intervals.

We obtain simultaneous intervals for all linear combinations of the FPC scores analogously. Given \mathcal{R} , let $\mathcal{L} \subseteq \mathcal{R}$ be a linear space with dimension $\leq p$. Then, asymptotically, it follows from the uniform result in Corollary 2 in Section 3 that for all linear combinations \mathbf{I}_μ , simultaneously, where $\mathbf{I} \in \mathcal{L}$,

$$\mathbf{I}_\mu \in \mathbf{I}_\mu \hat{\mu} \pm \sqrt{\chi^2_{2, 1-\alpha} \mathbf{I}_\mu \hat{\Sigma}_1 \mathbf{I}_\mu}, \tag{9}$$

with approximate probability $(1 - \alpha)$.

2.5. Choosing the number of eigenfunctions

To choose the number of eigenfunctions that provides a reasonable approximation to the infinite-dimensional process, we may use the cross-validation score based on the one-curve-leave-out prediction error (Rice and Silverman 1991). Let $\hat{\mu}^{(-)}$ and $\hat{\sigma}^{(-)}$ be the estimated mean and eigenfunctions after removing the data for the i th subject. Then we choose p so as to minimize the cross-validation score based on the squared prediction error,

$$\text{CV}(p) = \sum_{i=1}^n \sum_{j=1}^n \{ \hat{\mu}^{(-)}(t_j) - \hat{\mu}(t_j) \}^2, \tag{10}$$

where $\hat{\mu}^{(-)}$ is the predicted curve for the i th subject, computed after removing the data for this subject, that is, $\hat{\mu}^{(-)}(t) = \hat{\mu}(t) + \sum_{j=1}^n \hat{\sigma}_j^{(-)}(t)$, where $\hat{\sigma}_j^{(-)}$ is obtained by (5).

One can also adapt AIC-type criteria (Shibata 1981) to this situation. In simulations not reported here, we found that AIC is computationally more efficient while the results are similar to those obtained by cross-validation. A pseudo-Gaussian log-likelihood, summing the contributions from all subjects, conditional on the estimated FPC scores $\hat{\mu}$ (5), is given by

$$\begin{aligned} \text{AIC} = -\log \hat{L}(\hat{\mu}) = & \sum_{i=1}^n \left\{ -\frac{1}{2} \log |\Sigma_i| - \frac{1}{2} \log \hat{\Sigma}_i^{-2} \right. \\ & \left. - \frac{1}{2\hat{\Sigma}_i^2} \left(\tilde{\mathbf{Y}}_i - \hat{\mu}_i - \sum_{j=1}^p \hat{\sigma}_j \hat{\sigma}_j^T \right) \right. \\ & \left. \times \left(\tilde{\mathbf{Y}}_i - \hat{\mu}_i - \sum_{j=1}^p \hat{\sigma}_j \hat{\sigma}_j^T \right) \right\}, \tag{11} \end{aligned}$$

where we define $\text{AIC} = -\log \hat{L} + \frac{p}{n}$.

3. AIC PROPER USE

We derive consistency and distribution results demonstrating the consistency of the estimated FPC scores $\hat{\mu}$ in (5) for the true conditional expectations $\tilde{\mu}$ in (4). Uniform convergence of the local linear estimators of mean and covariance functions on bounded intervals plays a central role in obtaining these results and thus is established first (Thm. 1). Proofs are deferred to the Appendix.

The data $(t_i, y_i) = (t_i, y_i)$, $i = 1, \dots, n$, coming from model (1), are assumed to have the same distribution as (t, y) , with joint density $f(t, y)$. Assume that the observation times are iid with marginal density $f(t)$, but that dependence is allowed between observations t_i and t_j , coming from the same subject or cluster. The following assumptions pertain to the number of observations n_i made on the i th subject or cluster:

- (A1.1) The number of observations n_i made for the i th subject or cluster is a random variable with $n_i \stackrel{\text{iid}}{\sim} \text{Pois}(\lambda_i)$, where $\lambda_i > 0$ is a positive discrete random variable, with $\lambda_i < \infty$ and $\{\lambda_i > 1\} > 0$.

The observation times and measurements are assumed to be independent of the number of measurements, that is, for any subset $\mathcal{I} \subseteq \{1, \dots, n\}$ and for all $i = 1, \dots, n$,

- (A1.2) $(\{t_i : i \in \mathcal{I}\}, \{y_i : i \in \mathcal{I}\})$ is independent of n_i .

Writing $\tilde{\mathbf{T}} = (t_1, \dots, t_n)$,

definitions of these smoothers.] Kernel $\phi_1(\cdot)$ is also used for obtaining the estimate $\hat{\mu}$ for $\{(\cdot, \cdot) + \mathcal{I}^2\}$ with the local linear smoother. Let μ_1 , μ_2 , and μ_3 be the bandwidths for estimating $\hat{\mu}$, $\hat{\mu}_1$, and $\hat{\mu}_2$. Assume that ϕ_1 and ϕ_2 are compactly supported densities with properties (B2.1a) and (B2.2a) and (B2.1b) and (B2.2b). We develop asymptotics as the number of subjects $n \rightarrow \infty$, and require the following:

- (A2.1) $\mu_1 \rightarrow 0$, $\frac{4}{\mu_1} \rightarrow \infty$, and $\frac{6}{\mu_1} < \infty$.
- (A2.2) $\mu_2 \rightarrow 0$, $\frac{4}{\mu_2} \rightarrow \infty$, and $\frac{6}{\mu_2} < \infty$.
- (A2.3) $\mu_3 \rightarrow 0$, $\frac{4}{\mu_3} \rightarrow \infty$, and $\frac{6}{\mu_3} < \infty$.

Define the Fourier transforms of $\phi_1(\cdot)$ and $\phi_2(\cdot, \cdot)$ by $\phi_1(\omega) = \int \phi_1(\cdot) e^{-i\omega \cdot} d\cdot$ and $\phi_2(\omega, \nu) = \int \phi_2(\cdot, \cdot) e^{-i(\omega \cdot + \nu \cdot)} d(\cdot, \cdot)$. They satisfy the following:

- (A3.1) $\phi_1(\cdot)$ is absolutely integrable, that is, $\int |\phi_1(\cdot)| < \infty$.
- (A3.2) $\phi_2(\cdot, \cdot)$ is absolutely integrable, that is, $\iint |\phi_2(\cdot, \cdot)| < \infty$.

Assume that the fourth moment of ϕ_1 centered at $\mu(\cdot)$ is finite, that is,

(A4) $\int |\phi_1(\cdot - \mu(\cdot))|^4 < \infty$.

Then we obtain uniform convergence rates for local linear estimators $\hat{\mu}(\cdot)$ of $\mu(\cdot)$ and $\hat{\mu}_1(\cdot, \cdot)$ of $\mu_1(\cdot, \cdot)$ on compact sets and \mathcal{I}^2 .

Under (A1.1)–(A4) and (B1.1)–(B2.2b) with $\mu = 0$, $\mu_1 = 2$ in (B2.2a) and $\mu_2 = (0, 0)$, $\mu_3 = 2$ in (B2.2b),

$$\sup_{\mathcal{I}} |\hat{\mu}(\cdot) - \mu(\cdot)| = \left(\frac{1}{\sqrt{\mu}} \right) \tag{12}$$

and

$$\sup_{\mathcal{I}} |\hat{\mu}_1(\cdot, \cdot) - \mu_1(\cdot, \cdot)| = \left(\frac{1}{\sqrt{\mu}} \right). \tag{13}$$

The consistency of $\hat{\mu}_2(\cdot)$ is obtained as a consequence.

Under (A1.1)–(A4) and (B1.1)–(B2.2b) with $\mu = 0$, $\mu_1 = 2$ in (B2.2a) and $\mu_2 = (0, 0)$, $\mu_3 = 2$ in (B2.2b),

$$|\hat{\mu}_2 - \mu_2| = \left(\frac{1}{\sqrt{\mu}} \left(\frac{1}{2} + \frac{1}{\mu} \right) \right). \tag{14}$$

We note that the rates of convergence provided in (12) and (13) are slower than the optimal ones known for the case of smoothing functions or surfaces from sufficiently densely spaced independent measurements. These rates would be of order $(\sqrt{\log n} / (\mu))$ for function estimates and $(\sqrt{\log n} / (\mu^2))$ for surface estimates. It is an interesting question whether these rates remain optimal for the present dependent-data setting and whether they can be attained in the situation of dependent and sparse data that we are dealing with.

Next, consider the real separable Hilbert space $\mathcal{H}^2(\cdot) \equiv$ endowed with inner product $\langle \cdot, \cdot \rangle = \int_{\mathcal{I}} (\cdot) (\cdot)$ and norm $\| \cdot \| = \sqrt{\langle \cdot, \cdot \rangle}$ (Courant and Hilbert 1953). Let \mathcal{I} denote the set of indices of the eigenfunctions corresponding to eigenvalues of multiplicity 1. We obtain the consistency of the $\hat{\mu}$ in (3) for \mathcal{I} , the consistency of $\hat{\mu}_1$ in (3) for \mathcal{I} in the \mathcal{H}^2 norm $\| \cdot \|$, by choosing $\hat{\mu}$ appropriately when \mathcal{I} is of multiplicity 1, and furthermore the uniform consistency of $\hat{\mu}$ for \mathcal{I} on the bounded interval \mathcal{I} .

Under (A1.1)–(A4) and (B1.1)–(B2.2b) with $\mu = 0$, $\mu_1 = 2$ in (B2.2a) and $\mu_2 = (0, 0)$, $\mu_3 = 2$ in (B2.2b),

$$|\hat{\mu} - \mu| = \left(\frac{1}{\sqrt{\mu}} \right); \tag{15}$$

$$\| \hat{\mu}_1 - \mu_1 \| = \left(\frac{1}{\sqrt{\mu}} \right), \quad \mu \in \mathcal{I}; \tag{16}$$

and

$$\sup_{\mathcal{I}} |\hat{\mu}_2(\cdot) - \mu_2(\cdot)| = \left(\frac{1}{\sqrt{\mu}} \right), \quad \mu \in \mathcal{I}. \tag{17}$$

We remark that the rates (15)–(17) are direct consequences of the rates (12) and (13), as is evident from the proofs. If the rates in (12) and (13) are both (\cdot) , then the rates in (15)–(17) will also be (\cdot) .

For the following results, we require Gaussian assumptions:

- (A5) The FPC scores $\epsilon_{i,j}$ and measurement errors $\epsilon_{i,j}$ in (1) are jointly Gaussian.

We also assume that the data asymptotically follow a linear scheme:

- (A6) The number, location, and values of measurements for a given subject or cluster remain unaltered as $n \rightarrow \infty$.

The target trajectories that we aim to predict are

$$\tilde{\mu}(\cdot) = \mu(\cdot) + \sum_{i=1}^{\infty} \tilde{\epsilon}_i(\cdot), \quad i = 1, \dots, \infty, \tag{18}$$

with $\tilde{\epsilon}_i$ as defined in (4). We note that $\tilde{\mu}$ may be defined as a limit of random functions $\tilde{\mu}(\cdot) = \mu(\cdot) + \sum_{i=1}^{\infty} \tilde{\epsilon}_i(\cdot)$, as $\sup_{\mathcal{I}} [\tilde{\mu}(\cdot) - \tilde{\mu}(\cdot)]^2 \rightarrow 0$ (see Lemma A.3 in the App.). For any $i \geq 1$, the target curve $\tilde{\mu}(\cdot)$ is then estimated by

$$\hat{\mu}(\cdot) = \hat{\mu}(\cdot) + \sum_{i=1}^{\infty} \hat{\epsilon}_i(\cdot), \tag{19}$$

with $\hat{\epsilon}_i$ as in (5).

Assume (A1.1)–(A6) and (B1.1)–(B2.2b) with $\mu = 0$, $\mu_1 = 2$ in (B2.2a) and $\mu_2 = (0, 0)$, $\mu_3 = 2$ in (B2.2b). Then

$$\lim_{n \rightarrow \infty} \hat{\mu} = \tilde{\mu} \text{ in probability,} \tag{20}$$

and for all $\epsilon \in \mathcal{I}$,

$$\lim_{n \rightarrow \infty} \hat{\mu}(\epsilon) = \tilde{\mu}(\epsilon)$$

Applying Theorems 1 and 2, the estimate $\hat{\mu}(\cdot)$ is consistent for $\mu(\cdot)$ for all $\cdot \in \mathcal{I}$; that is, $\lim_{n \rightarrow \infty} \lim_{T \rightarrow \infty} \hat{\mu}(\cdot) = \mu(\cdot)$ in probability.

Assume (A1.1)–(A7) and (B1.1)–(B2.2b) with $\beta = 0$, $\nu = 2$ in (B2.2a) and $\beta = (0, 0)$, $\nu = 2$ in (B2.2b). For all $\cdot \in \mathcal{I}$ and $\cdot \in \mathfrak{R}$,

$$\lim_{n \rightarrow \infty} \lim_{T \rightarrow \infty} \left\{ \frac{\hat{\mu}(\cdot) - \mu(\cdot)}{\sqrt{\hat{\mu}(\cdot)}} \leq \cdot \right\} = \Phi(\cdot), \quad (22)$$

where Φ is the standard Gaussian cdf.

The number of random components and eigenfunctions that are needed in Theorems 3 and 4 to approximate the trajectory $\tilde{\mu}(\cdot)$ depends primarily on the complexity of the covariance structure $\Sigma(\cdot, \cdot)$ and on the number and location of the measurements observed for a given subject. It also depends on the sample size n , through the eigenfunction and covariance estimates. Although data-based choices for \cdot are available through (10) and (11) and are successful in practical applications, results (21) and (22) indicate that for large n , the number of components \cdot needs to be increased to obtain consistency, but these results do not provide further guidance as to how \cdot should be chosen in dependence on n .

We next establish $(1 - \alpha)$ asymptotic simultaneous inference for $\{\hat{\mu}(\cdot) - \mu(\cdot)\}$ on the domain \mathcal{I} , where $\mu(\cdot) = \mu(\cdot) + \sum_{i=1}^{\cdot} \mu_i(\cdot)$. For these results, we are providing not functional asymptotics, but instead finite-dimensional asymptotics, because the number of included components \cdot is considered fixed, whereas the sample size $n \rightarrow \infty$ as before. If \cdot is chosen such that only trajectories truncated at the first \cdot components $\mu_i(\cdot)$ of their expansion are of interest, then the following two results provide simultaneous confidence bands, as well as simultaneous confidence sets for the first \cdot random effects. Corollary 2 is a variation of Scheffé’s method.

Under (A1.1)–(A7) and (B1.1)–(B2.2b) with $\beta = 0$, $\nu = 2$ in (B2.2a) and $\beta = (0, 0)$, $\nu = 2$ in (B2.2b), for fixed number of components \cdot ,

$$\lim_{n \rightarrow \infty} \left\{ \sup_{\cdot \in \mathcal{I}} \frac{|\hat{\mu}(\cdot) - \mu(\cdot)|}{\sqrt{\hat{\mu}(\cdot)}} \leq \sqrt{\chi^2_{2, 1-\alpha}} \right\} \geq 1 - \alpha, \quad (23)$$

where $\chi^2_{2, 1-\alpha}$ is the $(1 - \alpha)$ th percentile of the chi-squared distribution with 2 degrees of freedom.

Assuming \cdot components, let $\mathcal{L} \subseteq \mathfrak{R}$ be a linear space with dimension $\leq \cdot$. By arguments analogous to the proof of Theorem 5, we obtain the asymptotic simultaneous $(1 - \alpha)$ confidence region for all linear combinations $\mathbf{1}^T \hat{\mu}(\cdot)$, where $\mathbf{1} \in \mathcal{L}$.

Under the assumptions of Theorem 5,

$$\lim_{n \rightarrow \infty} \left\{ \sup_{\mathbf{1} \in \mathcal{L}} \frac{|\mathbf{1}^T (\hat{\mu}(\cdot) - \mu(\cdot))|}{\sqrt{\mathbf{1}^T \hat{\mu}(\cdot) \mathbf{1}}} \leq \sqrt{\chi^2_{2, 1-\alpha}} \right\} \geq 1 - \alpha, \quad (24)$$

where $\chi^2_{2, 1-\alpha}$ is the $(1 - \alpha)$ th percentile of the chi-squared distribution with 2 degrees of freedom.

4. IM LA ION DIE

To illustrate the implementation of sparse FPC analysis by PACE, we construct 100 iid normal and 100 iid nonnormal samples each consisting of $n = 100$ random trajectories. The simulated processes have mean function $\mu(\cdot) = \sin(\cdot)$ and covariance function derived from two eigenfunctions, $\mu_1(\cdot) = -\cos(\cdot/10)/\sqrt{5}$ and $\mu_2(\cdot) = \sin(\cdot/10)/\sqrt{5}$, $0 \leq \cdot \leq 10$. We chose $\lambda_1 = 4$, $\lambda_2 = 1$, and $\lambda_i = 0$, $i \geq 3$, as eigenvalues and $\sigma^2 = .25$ as the variance of the additional measurement errors in (1), which are assumed to be normal with mean 0. For the smoothing steps, we use univariate and bivariate Epanechnikov kernel functions, that is, $K_1(\cdot) = 3/4(1 - \cdot^2)\mathbb{1}_{[-1,1]}(\cdot)$ and $K_2(\cdot, \cdot) = 9/16(1 - \cdot^2)(1 - \cdot^2)\mathbb{1}_{[-1,1]}(\cdot)\mathbb{1}_{[-1,1]}(\cdot)$, where $\mathbb{1}(\cdot) = 1$ if $\cdot \in [-1, 1]$ and 0 otherwise for any set \cdot . For an equally spaced grid $\{0, \dots, 50\}$ on $[0, 10]$ with $t_0 = 0$ and $t_{50} = 10$, let $\cdot = t_j + \epsilon_j$, where ϵ_j are iid with $\epsilon_j \sim (0, .1)$, $\epsilon_j = 0$ if $t_j < 0$, and $\epsilon_j = 10$ if $t_j > 10$, allowing for nonequidistant “jittered” designs. Each curve was sampled at a random number of points, chosen from a discrete uniform distribution on $\{1, \dots, 4\}$, and the locations of the measurements were randomly chosen from $\{1, \dots, 49\}$ without replacement. For the 100 normal samples, the FPC scores were generated from $(0, \cdot)$, whereas the \cdot for the nonnormal samples were generated from a mixture of two normals, $(\sqrt{2}, \cdot/2)$ with probability 1/2 and $(-\sqrt{2}, \cdot/2)$ with probability 1/2.

To demonstrate the superior performance of the conditional method, Table 1 reports the average mean squared error (MSE) for the true curves μ , $MSE = \sum_{i=1}^2 \int_0^{10} \{\mu_i(\cdot) - \hat{\mu}_i(\cdot)\}^2 / \cdot$, where $\hat{\mu}(\cdot) = \hat{\mu}(\cdot) + \sum_{i=1}^2 \hat{\mu}_i(\cdot)$ and the $\hat{\mu}_i$ ’s were obtained using either the proposed PACE method (5) or the integration method. The number of eigenfunctions \cdot in each run was chosen by the AIC (11). In each simulation consisting of 100 Monte Carlo runs (for a total of 400 runs, normal/mixture and sparse/nonsparse), there were always more than 95 runs in which two eigenfunctions were chosen.

Another outcome measure of interest is the average squared error (ASE) for the two FPC scores, $ASE(\cdot) = \sum_{i=1}^2 \int_0^{10} \{\hat{\mu}_i(\cdot) - \mu_i(\cdot)\}^2 / \cdot$, $\cdot = 1, 2$, also listed in Table 1. We also compared the two methods for irregular but nonsparse simulated data, where the number of observations for each curve was randomly chosen from $\{30, \dots, 40\}$, with results given in Table 1. We find that the gains in the sparse situation are dramatic when switching from the traditional method to the PACE method. For the case of an underlying normal distribution, the MSE was reduced by 43% using the PACE method (5) as compared with

	Method	FPC A					
		MSE	ASE(1)	ASE(2)	MSE	ASE(1)	ASE(2)
N = 100	Integration	1.33	.762	.453	1.30	.737	.453
	PACE	2.32	1.58	.622	2.25	1.53	.631
N = 30	Integration	.259	.127	.110	.256	.132	.105
	PACE	.286	.159	.115	.286	.168	.114

the traditional method; the ASE() were reduced by 52%/27% (= 1, 2). For the mixture distribution case, the decreases were still 42% for MSE and 52%/28% for ASE() (= 1, 2). In nonsparse situations, the traditional estimates provide reasonable approximations to the underlying integrals, but nevertheless PACE still produces better estimates, with improvements of 10%/10% for MSE and 20%/21%, 5%/8% for ASE(), = 1, 2, for normal/nonnormal samples. We conclude that the gains obtainable using PACE are substantial for sparse data and also extend to the case of dense and non-Gaussian data.

5. APPLICATION

5.1 L . . . CD4 C . . .

Because CD4 counts constitute a critical assessment of the status of the immune system and are used as an important marker in describing the progress to AIDS in adults, CD4 cell counts and CD4 percentages (i.e., CD4 counts divided by the total number of lymphocytes) are commonly used markers for the health status of human immunodeficiency virus (HIV) infected persons. The dataset considered here is from the Multicenter AIDS Cohort Study, which includes repeated measurements of physical exams, laboratory results, and CD4 percentages for 283 homosexual men who became HIV-positive between 1984 and 1991. All individuals were scheduled to have their measurements made at semiannual visits. However, because many individuals missed scheduled visits and the HIV infections happened randomly during the study, the data are sparse, with unequal numbers of repeated measurements per subject and different measurement times, , per individual. The number of observations per subject ranged from 1 to 14, with a median of 6. The trajectories in their entirety are assembled in Figure 1(a).

That the data from such a classical longitudinal study, with measurements intended to be spaced at regular 6-month intervals, are quite well suited for analysis by PACE is illustrated by Figure 2. As this figure shows, the assembled pairs (,) measurements

54.8(t)-2820-256dense-096002 91 T4 5cau9592 cm /Im1 Do Q

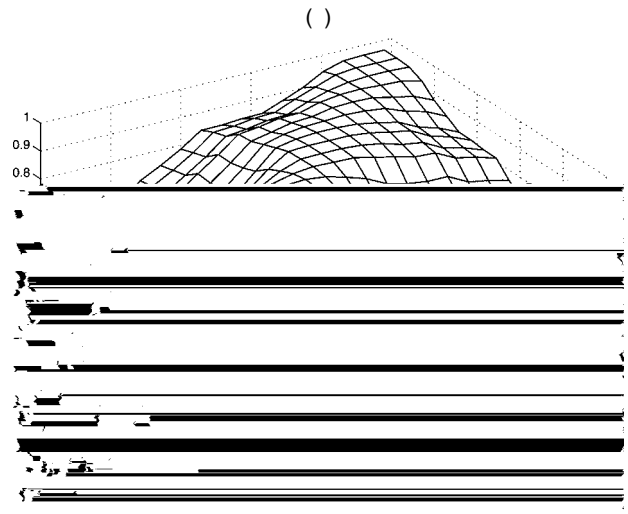
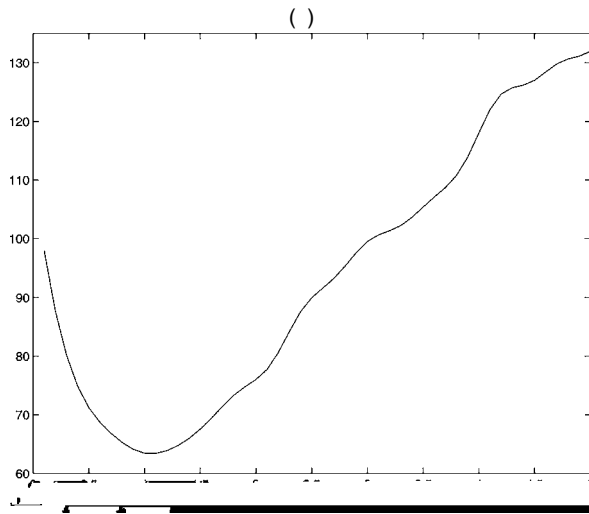


Figure 3. Estimated mean function of CD4 percentage. Figure 4. Estimated covariance surface of CD4 percentage.

early times, decreasing until about 1 year and then increasing again. Measurements made on the same subject are strongly correlated, irrespective of the time difference. However, the correlation between very early and late counts dies off relatively rapidly, whereas for middle and later times, the dependence patterns persist more strongly. These features would be difficult to anticipate in a traditional parametric model; they would not be produced, by, for example, linear random-effects models.

Next, consider the eigenfunction decomposition of the estimated covariance surface. Three eigenfunctions shown in the upper panels of Figure 4 are used to approximate the infinite-dimensional process. The choice $k = 3$ emerges as a reasonable choice, supported both by the AIC (11) and one-curve-leave-out cross-validation. The first eigenfunction is somewhat similar to the mean function, the second corresponds to a contrast be-

tween very early and very late times, and the third corresponds to a contrast between the early and the medium plus later times. These eigenfunctions account for 76.9%, 12.3%, and 8.1% of the total variation. Most of the variability is thus in the direction of overall CD4 percentage level. In exploring such data, extreme individual cases are difficult to detect by visual examination due to irregular sampling and substantial noise. One way to explore the variability in the sample and to single out extreme cases is to identify cases that exhibit large principal component scores in the directions of a few leading eigenfunctions (Jones and Rice 1992). Three such cases, corresponding to the largest absolute values of the projections on the first three eigenfunctions, are shown in the lower panels of Figure 4.

The predicted curves and 95% pointwise and simultaneous confidence bands for four randomly chosen individuals are displayed in Figure 5, where the principal component scores of

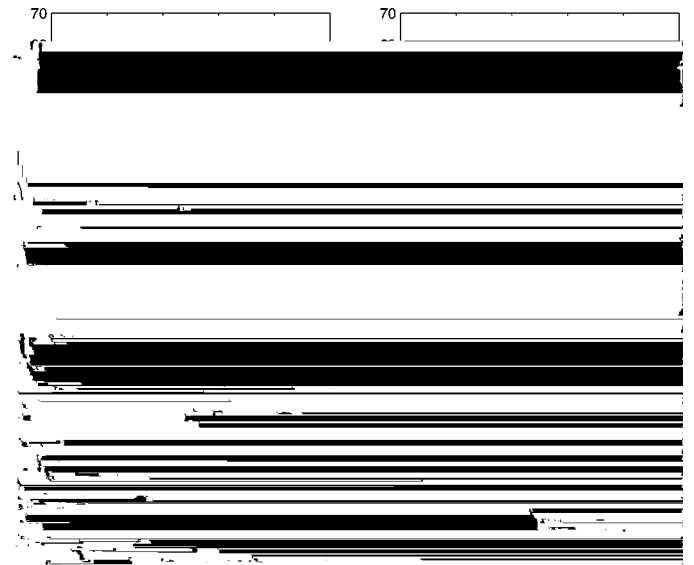
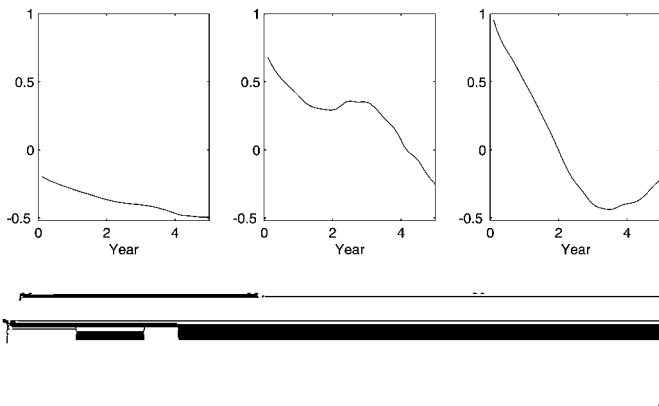


Figure 4. The first three eigenfunctions of the estimated covariance surface. The three cases shown in the lower panels correspond to the largest absolute values of the projections on the first three eigenfunctions.

Figure 5. Predicted curves and 95% pointwise and simultaneous confidence bands for four randomly chosen individuals.

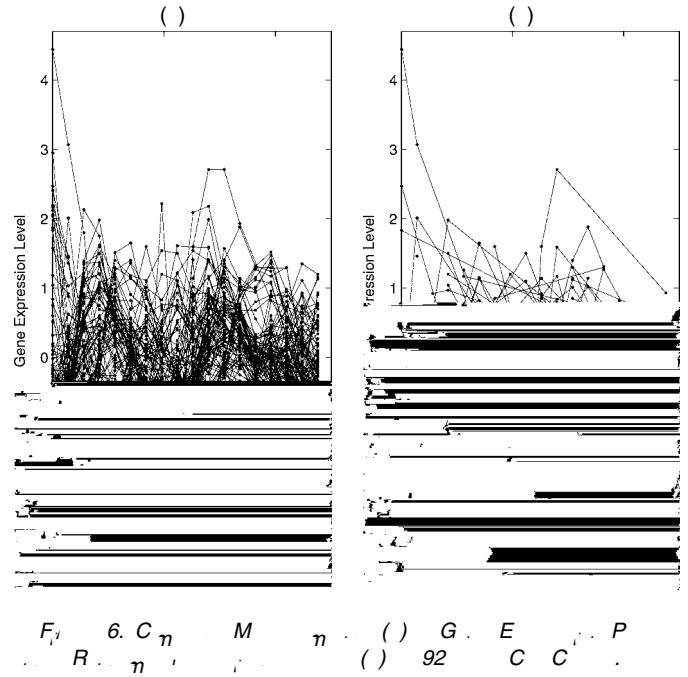
each subject are estimated using the PACE method. The predicted curves are seen to be reasonably close to the observations. Even for individuals with very sparse measurements, one is still able to effectively recover their random trajectories, combining the information from that individual and the entire collection. For example, the PACE principle of borrowing strength from the entire sample for predicting individual trajectories makes it feasible to predict trajectories and construct corresponding prediction bands for those cases where only one observation available per subject, as exemplified in the lower left panel of Figure 5. The predictions based on only one observation per subject work reasonably well, as is demonstrated in the second example described in Section 5.2 (see the lower right panel in Fig. 9, where only one single measurement enclosed in the circle is used for the prediction of the trajectory). Because we need to be able to consistently estimate the covariance structure, it is, however, not feasible to apply the method if there is only one observation available per subject for all subjects. Note that the 95% simultaneous bands show a widening near the endpoints due to end effects and increased variance near the ends, and that all observed data fall within these bands.

5.2 C C G . E . P

Time-course gene expression data (factor-synchronized) for the yeast cell cycle were obtained by Spellman et al. (1998). The experiment started with a collection of yeast cells, whose cycles were synchronized (factor-based) by a chemical process. There are 6,178 genes in total, and each gene expression profile consists of 18 data points, measured every 7 minutes between 0 and 119 minutes, covering two cell cycles. Of these genes, 92 had sufficient data and were identified by traditional methods, of which 43 are related to G1 phase regulation and 49 are related to non-G1 phase regulation (i.e., S, S/G2, G2/M, and M/G1 phases) of the yeast cell cycle; these genes serve as a training set. The gene expression level measurement at each time point is obtained as a logarithm of the expression-level ratio.

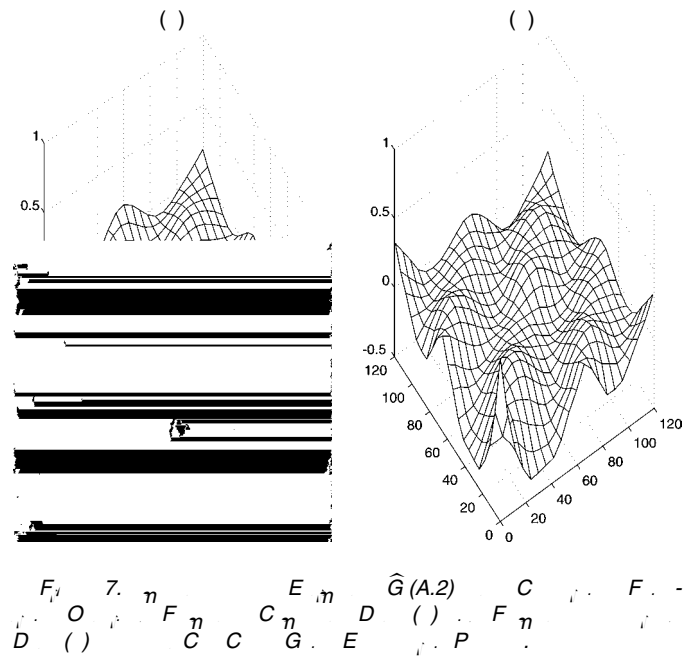
To demonstrate the usefulness of the PACE method for sparse functional data, we artificially “sparsify” the measurements made for the genes in the training data, then compare the results obtained from this “sparsified” data with those obtained from the complete data. To sparsify the expression measurements made for the n th gene expression profile, the number of measurements is randomly chosen between 1 and 6 with equal probability, and the measurement locations are then randomly selected from the 18 recorded gene expression measurements per profile. The median number of observations per gene expression profile for the resulting sparse data is just 3.

Analyses of both complete and sparsified yeast cell cycle profile data are illustrated in Figures 6–8. The two mean function estimates for the sparse and complete data, obtained by local linear smoothing of the pooled data, are close to each other and demonstrate periodicity [see Fig. 8(a), presenting two cell cycles]. The two smooth covariance surface estimates revealing the structure of the underlying process are displayed in Figure 7. Both surfaces are very similar and exhibit periodic features. We use the first two eigenfunctions to approximate the



expression profiles [Figs. 8(a) and 8(c)]. The estimates of the first two eigenfunctions obtained from both sparse and complete data are also close and reflect periodicity, explaining approximately 75% of the total variation.

We randomly select four genes, and present the predicted profiles obtained from both sparse and complete data and the confidence bands using only the sparse data in Figure 9. We note that the trajectories obtained for the complete data are enclosed in the simultaneous 95% confidence bands constructed from the sparse data. The predictions obtained from the sparse data are similar to those constructed from the complete data and are reasonable when compared with the complete measurements. This demonstrates that the PACE method allows us to effectively recover entire individual trajectories from fragmental data.



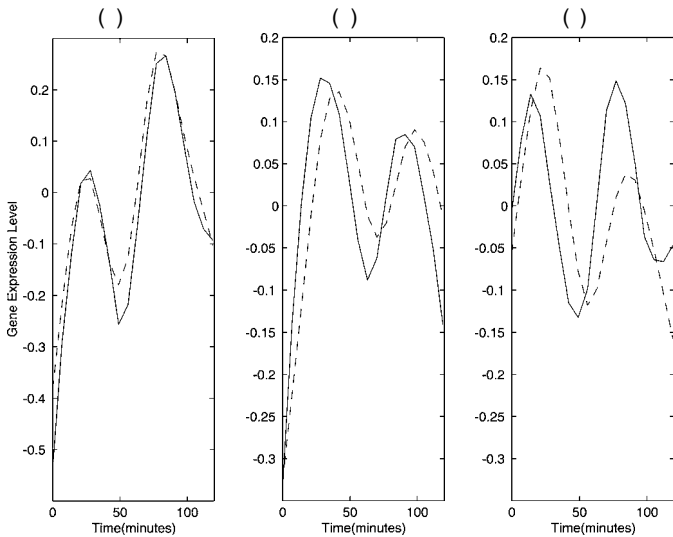


Figure 8. Imputed trajectories for subjects 1, 2, and 3. The solid line represents the observed data, and the dashed line represents the imputed data.

6. CONCLUDING REMARK

Besides the general application to FPC analysis for sparse and irregular data, an application of our proposed PACE method to impute missing data in longitudinal studies is also feasible. Consider a regular design where for some subjects many data are missing. The PACE method can then be used to impute the missing data from predicted trajectories.

An interesting finding from the simulation study is that the PACE method improves on traditional FPC analysis even under dense and regular designs. This improvement is due to replacing integrals by conditional expectations when determining FPC scores. The conditioning step can be interpreted as shrinkage

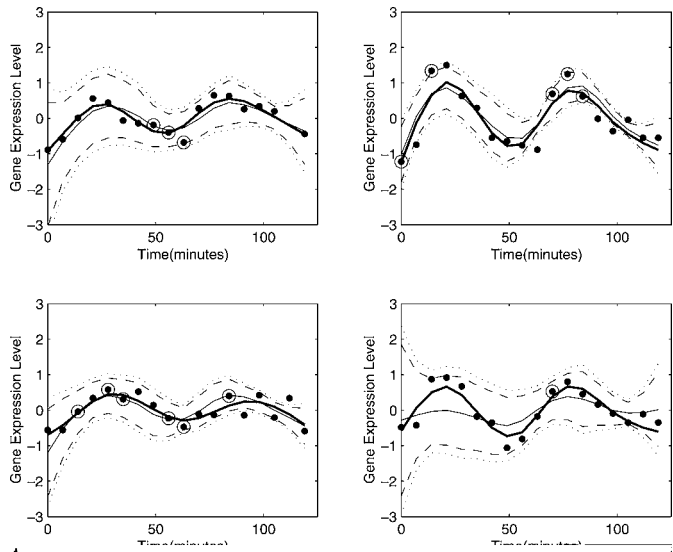


Figure 9. Fitted trajectories for subjects 1, 2, 3, and 4. The solid line with dots represents observed data, the dashed line represents imputed data, and the solid line with circles represents the fitted curve.

of these random effects toward 0. The observed improvement indicates that PACE can also be used to advantage for regularly spaced data, which enhances this method's appeal. We conclude that the underlying principle of borrowing strength from an entire sample of curves to predict individual strength trajectories shows promise in applications.

APPENDIX: PROOF AND AUXILIARY RESULTS

We assume regularity conditions for the marginal and joint densities $f_1(\cdot)$, $f_2(\cdot)$, and $f_2(t_1, t_2, \cdot, \cdot)$. Let ν_1, ν_2 , and ν be given integers, with $0 \leq \nu_1 + \nu_2 < \nu$. Then the following conditions apply:

- (B1.1) $f_1(\cdot)$ exists and is continuous on \mathcal{T} with $f_1(\cdot) > 0$ on \mathcal{T} .
- (B1.2) $f_2(\cdot, \cdot)$ exists and is uniformly continuous on $\mathcal{T} \times \mathcal{R}$.
- (B1.3) $f_2(\cdot, \cdot) / (f_1^{\nu_1} f_2^{\nu_2})$ exists and is uniformly continuous on $\mathcal{T}^2 \times \mathcal{R}^2$, for $\nu_1 + \nu_2 = \nu, 0 \leq \nu_1, \nu_2 \leq \nu$.

The assumptions for kernel functions $k_1: \mathcal{R} \rightarrow \mathcal{R}$ and $k_2: \mathcal{R}^2 \rightarrow \mathcal{R}$ are as follows. We say that a bivariate kernel function k_2 is of order (ν_1, ν_2) , where ν is a multi-index $\nu = (\nu_1, \nu_2)$, if

$$\iint k_2(t_1, t_2) dt_1 dt_2 = 0, \quad 0 \leq \nu_1 + \nu_2 < \nu, \nu_1 \neq \nu_1, \nu_2 \neq \nu_2$$

$$= \begin{cases} 0, & 0 \leq \nu_1 + \nu_2 < \nu, \nu_1 \neq \nu_1, \nu_2 \neq \nu_2 \\ (-1)^{|\nu|} |\nu|!, & \nu_1 = \nu_1, \nu_2 = \nu_2 \\ \neq 0, & \nu_1 + \nu_2 = \nu, \end{cases} \quad (A.1)$$

where $|\nu| = \nu_1 + \nu_2$. A univariate kernel k_1 is of order (ν) for a univariate $\nu = \nu_1$, if (A.1) holds with $\nu_2 = 0$ on the right side, integrating only over the argument t_2 on the left side.

- (B2.1a) k_1 is compactly supported, $\|k_1\|^2 = \int k_1^2(t) dt < \infty$.
- (B2.2a) k_1 is a kernel function of order (ν) .
- (B2.1b) k_2 is compactly supported, $\|k_2\|^2 = \iint k_2^2(t_1, t_2) dt_1 dt_2 < \infty$.
- (B2.2b) k_2 is a kernel function of order (ν) .

We define the local linear scatterplot smoother for $\mu(\cdot)$ by minimizing

$$\sum_{i=1}^n \sum_{j=1}^n \frac{1}{n} \left(\frac{y_i - y_j}{\mu(t_i) - \mu(t_j)} \right) \{ y_i - y_j - \mu(t_i) - \mu(t_j) \}^2 \quad (A.2)$$

with respect to μ_0 and μ_1 . The estimate of $\mu(\cdot)$ is then $\hat{\mu}(\cdot) = \hat{\mu}_0(\cdot)$. The local linear surface smoother for (\cdot, \cdot) is defined by minimizing

$$\sum_{i=1}^n \sum_{1 \leq j \neq i \leq n} \frac{1}{n} \left(\frac{y_i - y_j}{\mu(t_i) - \mu(t_j)} \right) \times \{ (y_i, t_i) - (y_j, t_j) - (\mu_0(t_i), \mu_1(t_i)) \}^2, \quad (A.3)$$

where $(y_i, t_i) = y_i + \mu_1(t_i) + \mu_2(t_i)$. Minimization is with regard to $\mu = (\mu_0, \mu_1, \mu_2)$, yielding the estimate $\hat{\mu}(\cdot, \cdot) = \hat{\mu}_0(\cdot, \cdot)$. To obtain the adjusted estimate of (\cdot, \cdot) on the diagonal [i.e., $\hat{\mu}(\cdot, \cdot)$], we first rotate both the t -axis and y -axis by 45 degrees clockwise and obtain the coordinates of (y_i, t_i) in the rotated axes, denoted by (y_i^*, t_i^*) , that is, $(y_i^*, t_i^*) = \begin{pmatrix} \sqrt{2}/2 & \sqrt{2}/2 \\ -\sqrt{2}/2 & \sqrt{2}/2 \end{pmatrix} \begin{pmatrix} y_i \\ t_i \end{pmatrix}$.

We then define the surface estimate $\hat{\mu}(\cdot, \cdot)$ by minimizing the weighted least squares,

$$\sum_{i=1}^n \sum_{1 \leq j \neq i \leq n} \frac{1}{n} \left(\frac{y_i^* - y_j^*}{\mu(t_i^*) - \mu(t_j^*)} \right) \times \{ (y_i^*, t_i^*) - (\mu_0(t_i^*), \mu_1(t_i^*)) \}^2, \quad (A.4)$$

where $\hat{\mu}(\cdot) = \hat{\mu}_0(\cdot) + \hat{\mu}_1(\cdot) + \hat{\mu}_2(\cdot)$. Minimization is with respect to $\hat{\mu} = (\hat{\mu}_1, \hat{\mu}_2, \hat{\mu}_3)$, leading to $\hat{\mu}(\cdot) = \hat{\mu}_0(\cdot)$. Because of the rotation, the estimate of the covariance surface on the diagonal, $\hat{\mu}(\cdot)$, is now indeed $\hat{\mu}(0, \sqrt{2})$, obtained with the rotated coordinates.

The following auxiliary results provide the weak uniform convergence rate for univariate weighted averages defined later (cf. Bhattacharya and Müller 1993). For a positive integer $\mu \geq 1$, let $(f_j)_{j=1, \dots, \mu}$ be a collection of real functions $f_j: \mathbb{R}^2 \rightarrow \mathbb{R}$, that satisfy the following conditions:

- (C1.1a) f_j are uniformly continuous on $\mathcal{T} \times \mathfrak{R}$.
- (C1.2a) The functions $(f_j / \mu)(\cdot, \cdot)$ exist for all arguments (\cdot, \cdot) and are uniformly continuous on $\mathcal{T} \times \mathfrak{R}$.
- (C1.3a) $\int f_j^2(\cdot, \cdot) < \infty$.

Bandwidths $\mu = \mu(\cdot)$ used for one-dimensional smoothers are assumed to satisfy the following:

(C2.1a) $\mu \rightarrow 0, \mu^{+1} \rightarrow \infty, \mu^{+2} < \infty, \text{ as } \mu \rightarrow \infty$.

Define the weighted averages

$$\begin{aligned} \hat{\mu}_j &= \hat{\mu}_j(\cdot) \\ &= \frac{1}{\mu^{+1}} \sum_{i=1}^{\mu} \frac{1}{\mu} \sum_{l=1}^{\mu} (f_j(\cdot, \cdot))_l \left(\frac{\cdot}{\mu} \right), \quad j = 1, \dots, \mu \end{aligned}$$

and the quantity

$$\begin{aligned} \mu &= \mu(\cdot) \\ &= \int (f_j(\cdot, \cdot))_l, \quad j = 1, \dots, \mu \end{aligned}$$

Under (A1.1), (A1.2), (A3.1), (B1.1), (B1.2), (B2.1a), (B2.2a), (C1.1a)–(C1.3a), and C(2.1a), $\sup_{\mathcal{T}} |\hat{\mu}_j - \mu_j| = O_p(1/\sqrt{\mu^{+1}})$.

Note that $|\hat{\mu}_j - \mu_j| \leq \sup_{\mathcal{T}} |\hat{\mu}_j - \mu_j| + \{ \sup_{\mathcal{T}} |\hat{\mu}_j - \mu_j| \}$, where \mathcal{T} takes values in \mathcal{T} and $|\hat{\mu}_j - \mu_j| = O_p(1/\sqrt{\mu^{+1}})$ implies that $\hat{\mu}_j = O_p(1/\sqrt{\mu^{+1}})$.

Using a Taylor expansion to order μ^{-1} , it is easy to show that $\hat{\mu}_j = \mu_j + O_p(\mu^{-1})$, where the remainder term is uniform in \mathcal{T} , observing that $(f_j / \mu)(\cdot, \cdot)$ and $(f_j / \mu)(\cdot, \cdot)$ are uniformly continuous. It remains to show that $\{ \sup_{\mathcal{T}} |\hat{\mu}_j - \mu_j| \} = O_p(1/\sqrt{\mu^{+1}})$. Recall that the inverse Fourier transform is $f_1(\cdot) = \int f_1(\cdot) d\mu$. We may insert $f_1((\cdot - \cdot) / \mu) = \int f_1(\cdot) / \mu d\mu$ into $\hat{\mu}_j$. Letting

$$\hat{\mu}_j = \frac{1}{\mu^{+1}} \sum_{i=1}^{\mu} \frac{1}{\mu} \sum_{l=1}^{\mu} (f_j(\cdot, \cdot))_l,$$

we obtain

$$\begin{aligned} &= \frac{1}{\mu^{+1}} \sum_{i=1}^{\mu} \frac{1}{\mu} \sum_{l=1}^{\mu} \left(\frac{\cdot}{\mu} \right)_l (f_j(\cdot, \cdot))_l \\ &= \frac{1}{2} \frac{1}{\mu} \int (f_j(\cdot, \cdot))_l^{-1}(\mu) d\mu, \end{aligned}$$

and thus

$$\sup_{\mathcal{T}} |\hat{\mu}_j - \mu_j| \leq \frac{1}{2} \frac{1}{\mu} \int |f_j(\cdot, \cdot) - \mu_j| \cdot |f_j(\cdot, \cdot)| d\mu.$$

Note that $|f_j(\cdot, \cdot) - \mu_j| \leq \sqrt{[f_j(\cdot, \cdot) - \mu_j]^2}$, and be-

cause $\{\tilde{\mathbf{T}}, \tilde{\mathbf{Y}}, \cdot\}$ are iid, using the Cauchy–Schwarz inequality,

$$\begin{aligned} \text{var}(\hat{\mu}_j(\cdot)) &= \frac{1}{\mu^{+1}} \text{var} \left\{ \frac{1}{\mu} \sum_{i=1}^{\mu} (f_j(\cdot, \cdot))_i \right\} \\ &\leq \frac{1}{\mu^{+1}} \left\{ \left(\frac{1}{\mu} \sum_{i=1}^{\mu} (f_j(\cdot, \cdot))_i \right)^2 \right\} \\ &\leq \frac{1}{\mu^{+1}} \left\{ \frac{1}{\mu^2} \left(\sum_{i=1}^{\mu} f_j^2(\cdot, \cdot) \right) \left(\sum_{i=1}^{\mu} 1 \right) \right\} \\ &\leq \frac{1}{\mu^{+1}} \left\{ \frac{1}{\mu^2} \sum_{i=1}^{\mu} (f_j^2(\cdot, \cdot))_i \right\} = \frac{1}{\mu^{+1}} f_j^2(\cdot, \cdot), \end{aligned}$$

implying that

$$\begin{aligned} &\left\{ \sup_{\mathcal{T}} |\hat{\mu}_j - \mu_j| \right\} \\ &\leq \frac{1}{2} \frac{1}{\mu} \int |f_j(\cdot, \cdot)| d\mu \end{aligned}$$

P. n 1

From (A.2), the local linear estimator $\hat{\mu}(\cdot)$ of the mean function $\mu(\cdot)$ can be written explicitly as

$$\hat{\mu}(\cdot) = \hat{\mu}_0(\cdot) = \frac{\sum \frac{1}{\sum} \sum \dots - \frac{\sum \frac{1}{\sum} \sum \dots (\dots)}{\sum \frac{1}{\sum} \sum} \hat{\mu}_1(\cdot), \tag{A.5}$$

where

$$\begin{aligned} \hat{\mu}_1(\cdot) &= \left(\sum \frac{1}{\sum} \sum \dots \right) \\ &\quad - \left(\sum \frac{1}{\sum} \sum \dots \right) \sum \frac{1}{\sum} \sum \dots \\ &\quad / \left(\sum \frac{1}{\sum} \sum \dots \right) \\ &\quad \times \left(\sum \frac{1}{\sum} \sum \dots \right)^2 \\ &\quad - \left(\sum \frac{1}{\sum} \sum \dots \right)^2 / \left(\sum \frac{1}{\sum} \sum \dots \right) \end{aligned} \tag{A.6}$$

Here $\dots = \dots(\dots) / \dots(\dots)$, where \dots_1 is a kernel function of order $(0, 2)$ satisfying (B2.1a) and (B2.2a), and $\hat{\mu}_1(\cdot)$ is an estimator for the first derivative $\mu'(\cdot)$ of μ at \dots .

Considering the Nadaraya–Watson estimator of μ , $\hat{\mu}_0(\cdot) = (\sum \dots / \sum \dots)$ and $\hat{\mu}_1(\cdot) = \sum \dots / \dots$, we choose $\dots = 0$, $\dots = 2$, $\dots = 2$, $\dots_1(\dots) = \dots$, and $\dots_2(\dots) \equiv 1$ in Lemma A.1. Obviously, $\hat{\mu}_0(\cdot) = \dots(\dots)$, with $\dots_1(\dots) = \dots/2$ and $\hat{\mu}_1(\cdot) = \dots$. Using Slutsky’s theorem and Lemma A.1, it follows that $\sup_{\mathcal{T}} |\hat{\mu}_0(\cdot) - \mu(\cdot)| = O_p(1/\sqrt{n})$ and $\sup_{\mathcal{T}} |\hat{\mu}_1(\cdot) - \mu'(\cdot)| = O_p(1/\sqrt{n})$.

For the uniform consistency of $\hat{\mu}_1$ as an estimator of the derivative μ' , define $\dots, 1 \leq \dots \leq 3$, $\dots_1 = \int \dots_1(\cdot) \dots$, and the kernel function $\tilde{\mu}_1(\cdot) = \dots_1(\cdot) / \dots_1$; furthermore, $\dots_1(\dots) = \dots$, $\dots_2(\dots) \equiv 1$ and $\dots_3(\dots) = \dots$. Observe that $\tilde{\mu}_1$ is of order $(1, 3)$, $\sup_{\mathcal{T}} |\hat{\mu}_1(\cdot) - \mu'(\cdot)| = O_p(1/\sqrt{n})$, and define

$$\begin{aligned} \tilde{\mu}_1(\dots, \dots, \dots) &= \frac{1 - 2\hat{\mu}_0(\dots)}{3 - 2\hat{\mu}_0^2(\dots)} \quad \text{and} \\ \dots_1(\dots, \dots, \dots) &= \frac{1 - 2\mu(\dots)}{3}. \end{aligned}$$

Then

$$\begin{aligned} \hat{\mu}_1(\cdot) &= \tilde{\mu}_1(\dots, \dots, \dots) \\ &= \left[\dots_1(\dots, \dots, \dots) + \frac{2(\mu(\dots) - \hat{\mu}_0(\dots))}{3} \right] \\ &\quad \times \frac{3}{3 + 2\hat{\mu}_0^2(\dots)}. \end{aligned}$$

Note that $\mu_1 = (\mu' + \dots)(\dots)$, $\mu_2 = \dots(\dots)$, and $\mu_3 = \dots(\dots)$, implying that $\sup_{\mathcal{T}} |\dots - \mu| = O_p(1/\sqrt{n})$, for $\dots = 1, 2, 3$, by Lemma A.1. Using the uniform version of Slutsky’s theorem, $\sup_{\mathcal{T}} |\dots_1(\dots, \dots, \dots) - \mu'(\cdot)| = O_p(1/\sqrt{n})$ follows.

Considering the uniform convergence of $\hat{\mu}_0$ for μ , note that

$$\hat{\mu}_0(\cdot) = \hat{\mu}_0(\cdot) + \frac{2\hat{\mu}_1(\cdot)}{\hat{\mu}_0(\cdot)} \dots$$

Because $\sup_{\mathcal{T}} |\dots - \mu| = O_p(1/\sqrt{n})$, $\sup_{\mathcal{T}} |\hat{\mu}_0(\cdot) - \mu(\cdot)| = O_p(1/\sqrt{n})$, and $\sup_{\mathcal{T}} |\hat{\mu}_1(\cdot) - \mu'(\cdot)| = O_p(1/\sqrt{n})$, we have $\sup_{\mathcal{T}} |\dots - \hat{\mu}_0(\cdot)| = O_p(1/\sqrt{n})$, as $\dots < \infty$. As $\sup_{\mathcal{T}} |\hat{\mu}_0(\cdot) - \mu(\cdot)| = O_p(1/\sqrt{n})$, the result (12) follows.

We proceed to show (13). In the local linear estimator for the covariance $\dots(\dots)$, we used the raw observations, $\dots(\dots) = (\dots - \hat{\mu}_0(\dots))(\dots - \hat{\mu}_0(\dots))$, instead of $\tilde{\mu}_1(\dots, \dots) = (\dots - \mu(\dots))(\dots - \mu(\dots))$. Note that

$$\begin{aligned} \dots(\dots) &= \tilde{\mu}_1(\dots, \dots) + (\dots - \mu(\dots))(\mu(\dots) - \hat{\mu}_0(\dots)) \\ &\quad + (\dots - \mu(\dots))(\mu(\dots) - \hat{\mu}_0(\dots)) \\ &\quad + (\mu(\dots) - \hat{\mu}_0(\dots))(\mu(\dots) - \hat{\mu}_0(\dots)). \end{aligned}$$

Because $\sup_{\mathcal{T}} |\hat{\mu}_0(\cdot) - \mu(\cdot)| = O_p(1/\sqrt{n})$ by (12), letting $\dots_1(1, 2, \dots) = (\dots - \mu(1))(\dots - \mu(2))$, $\dots_2(1, 2, \dots) = 1 - \mu(1)$, and $\dots_3(1, 2, \dots) \equiv 1$, then $\sup_{\mathcal{T}} |\dots| = O_p(1)$, for $\dots = 1, 2, 3$, by Lemma A.2. This implies that $\sup_{\mathcal{T}} |\dots| = O_p(1/\sqrt{n}) = O_p(1/\sqrt{n})$ and $\sup_{\mathcal{T}} |\dots| = O_p(1/\sqrt{n}) = O_p(1/\sqrt{n})$. Because $\sup_{\mathcal{T}} |\hat{\mu}_0(\cdot) - \mu(\cdot)|^2 = O_p(1/n)$ are negligible compared with \dots_1 , the local linear estimator, $\hat{\mu}_1(\dots)$, of $\dots(\dots)$ obtained from $\tilde{\mu}_1(\dots, \dots)$ is asymptotically equivalent to that obtained from $\dots_1(\dots, \dots)$, denoted by $\tilde{\mu}_1(\dots, \dots)$. Analogously to the proof of (12), using Lemma A.2 and the uniform version of Slutsky’s theorem, we obtain the uniform consistency of the local linear estimator $\hat{\mu}_1(\dots)$.

P. C. 1

Because $\hat{\mu}_0(\cdot)$ is a uniformly consistent estimator of $\{\dots(\dots) + \dots\}$, analogously to (12), (14) follows by applying (13).

P. n 2

Define the rank-one operator $\otimes = \langle \dots, \dots \rangle$ for $\dots \in \dots$, and denote the separable Hilbert space of Hilbert–Schmidt operators on \dots by $\dots = \dots(\dots)$, endowed by $\langle \dots, \dots \rangle = \text{tr}(\dots_1^* \dots_2) = \sum \langle \dots_1, \dots_2 \rangle$ and $\|\dots\|^2 = \langle \dots, \dots \rangle$, where $\dots_1, \dots_2 \in \dots$, \dots_1^* is the adjoint of \dots_1 and $\{\dots : \dots \geq 1\}$ is any complete orthonormal system in \dots . The covariance operator \mathbf{G} (resp. $\hat{\mathbf{G}}$) is generated by the kernel $\dots(\dots)$ (resp. $\hat{\mu}_0(\dots)$), that is, $\mathbf{G}(\dots) = \int_{\mathcal{T}} \dots(\dots) \dots$ [resp. $\hat{\mathbf{G}}(\dots) = \int_{\mathcal{T}} \hat{\mu}_0(\dots) \dots$]. It is obvious that \mathbf{G} and $\hat{\mathbf{G}}$ are Hilbert–Schmidt operators, and (13) implies that $\|\hat{\mathbf{G}} - \mathbf{G}\| = O_p(1/\sqrt{n})$.

Let $\mathcal{I} = \{\dots : \dots = 1\}$, $\mathcal{I}' = \{\dots : |\mathcal{I}| = 1\}$, where $|\mathcal{I}|$ denotes the number of elements in \mathcal{I} . To obtain (16), let $\mathbf{P} = \sum_{\mathcal{I}} \dots \otimes \dots$ and $\hat{\mathbf{P}} = \sum_{\mathcal{I}} \dots \otimes \hat{\mu}_0(\dots)$ denote the true and estimated orthogonal projection operators from \dots to the subspace spanned by $\{\dots : \dots \in \mathcal{I}\}$. For fixed $0 < \dots < \min\{|\dots - \dots| : \dots \notin \mathcal{I}\}$, let $\mathbf{A} = \{\dots \in \mathcal{C} : |\dots - \dots| = \dots\}$, where \mathcal{C} stands for the complex numbers. The resolvent of \mathbf{G} (resp. $\hat{\mathbf{G}}$) is denoted by \mathbf{R} (resp. $\hat{\mathbf{R}}$), that is, $\mathbf{R}(\dots) = (\mathbf{G} - \dots)^{-1}$ [resp. $\hat{\mathbf{R}}(\dots) = (\hat{\mathbf{G}} - \dots)^{-1}$]. As $\hat{\mathbf{R}}(\dots) = \mathbf{R}(\dots) [1 + (\hat{\mathbf{G}} - \mathbf{G})\mathbf{R}(\dots)]^{-1} = \mathbf{R}(\dots) \sum_{j=0}^{\infty} [(\hat{\mathbf{G}} - \mathbf{G})\mathbf{R}(\dots)]^j$, $\|\hat{\mathbf{R}}(\dots) - \mathbf{R}(\dots)\| \leq (\|\hat{\mathbf{G}} - \mathbf{G}\| \|\mathbf{R}(\dots)\|) / (1 - \|\hat{\mathbf{G}} - \mathbf{G}\| \|\mathbf{R}(\dots)\|)$. Note that $\mathbf{P} = (\dots)^{-1} \int_{\mathbf{A}} \mathbf{R}(\dots)$, $\hat{\mathbf{P}} = (\dots)^{-1} \int_{\mathbf{A}} \hat{\mathbf{R}}(\dots)$. Let $\dots = \sup\{\|\mathbf{R}(\dots)\| : \dots \in \mathbf{A}\} < \infty$, and let \dots be such that $0 < \dots < 1/(2 \dots)$; then

$$\begin{aligned} \|\hat{\mathbf{P}} - \mathbf{P}\| &\leq \int_{\mathbf{A}} \|\hat{\mathbf{R}}(\dots) - \mathbf{R}(\dots)\| \dots \\ &\leq \frac{\|\hat{\mathbf{G}} - \mathbf{G}\|}{1 - \|\hat{\mathbf{G}} - \mathbf{G}\|} \dots \leq 2 \dots \end{aligned}$$

Considering corresponding to $\in \mathcal{T}'$, choose $\hat{\cdot}$ such that $(\hat{\cdot}, \cdot) > 0$. Then

$$\begin{aligned} \|\hat{\mathbf{P}} - \mathbf{P}\|^2 &= 2(1 - \langle \hat{\cdot} \otimes \hat{\cdot}, \otimes \cdot \rangle) \\ &= 2(1 - \langle \hat{\cdot}, \cdot \rangle^2) \geq \|\hat{\cdot} - \cdot\|^2, \end{aligned}$$

and (16) follows. Note that $\hat{\cdot} = \langle \cdot, \mathbf{G}(\cdot) \rangle$ and $\hat{\cdot} = \langle \hat{\cdot}, \hat{\mathbf{G}}(\hat{\cdot}) \rangle$; then (15) follows by applying Slutsky's theorem. To obtain (17), for fixed $\in \mathcal{T}'$,

$$\begin{aligned} &|\hat{\cdot}(\cdot) - \cdot(\cdot)| \\ &= \left| \int_0^{\mathcal{T}} \hat{\cdot}(\cdot, \cdot) \hat{\cdot}(\cdot) - \int_0^{\mathcal{T}} \cdot(\cdot, \cdot) \cdot(\cdot) \right| \\ &\leq \int_0^{\mathcal{T}} |\hat{\cdot}(\cdot, \cdot) - \cdot(\cdot, \cdot)| \cdot |\hat{\cdot}(\cdot)| \\ &\quad + \int_0^{\mathcal{T}} |\cdot(\cdot, \cdot)| \cdot |\hat{\cdot}(\cdot) - \cdot(\cdot)| \\ &\leq \sqrt{\int_0^{\mathcal{T}} (\hat{\cdot}(\cdot, \cdot) - \cdot(\cdot, \cdot))^2} + \sqrt{\int_0^{\mathcal{T}} \cdot^2(\cdot, \cdot)} \|\hat{\cdot} - \cdot\|. \end{aligned}$$

Due to (13) and (16), assuming > 0 without loss of generality, we have $|\hat{\cdot}(\cdot)/\hat{\cdot} - \cdot(\cdot)| = (1/\sqrt{\cdot^2})$, uniformly in $\in \mathcal{T}$. Then (17) follows by applying (15).

The next result ensures that the target trajectory $\tilde{\cdot}$ is well defined.

For the positive definite covariance operator \mathbf{G} generated by the continuous symmetric function (\cdot, \cdot) on \mathcal{T}^2 , as $\rightarrow \infty$,

$$\sup_{\in \mathcal{T}} [\tilde{\cdot}(\cdot) - \cdot(\cdot)]^2 \rightarrow 0. \tag{A.7}$$

Because the covariance operator \mathbf{G} generated by the continuous symmetric function (\cdot, \cdot) is positive definite, by Mercer's theorem, $\sum_{=+1}^{\infty} \cdot(\cdot)$ converges to 0 uniformly in $(\cdot, \cdot) \in \mathcal{T}^2$. Note that $\tilde{\cdot}(\cdot) - \cdot(\cdot) = [\sum_{=+1}^{\infty} \cdot(\cdot) \tilde{\mathbf{Y}}]$. From

$$\begin{aligned} &\sup_{\in \mathcal{T}} \text{var} \left(\sum_{=+1}^{\infty} \cdot(\cdot) \right) \\ &= \sup_{\in \mathcal{T}} \left\{ \left[\left[\sum_{=+1}^{\infty} \cdot(\cdot) \tilde{\mathbf{Y}} \right]^2 \right] \right. \\ &\quad \left. + \left[\text{var} \left(\sum_{=+1}^{\infty} \cdot(\cdot) \tilde{\mathbf{Y}} \right) \right] \right\} \\ &= \sup_{\in \mathcal{T}} \sum_{=+1}^{\infty} \cdot^2(\cdot) \rightarrow 0, \end{aligned}$$

and $[\text{var}(\sum_{=+1}^{\infty} \cdot(\cdot) \tilde{\mathbf{Y}})] \geq 0$, (A.7) follows.

P n 3

Recall that $\hat{\cdot} = \hat{\cdot}, \hat{\cdot} \hat{\Sigma}_{\mathbf{Y}}^{-1} (\tilde{\mathbf{Y}} - \hat{\cdot})$, where the (\cdot, \cdot) th entry of the \times matrix $\hat{\Sigma}_{\mathbf{Y}}$ is $(\hat{\Sigma}_{\mathbf{Y}})_{\cdot, \cdot} = \hat{\cdot}(\cdot, \cdot) + \hat{\cdot}^2$ with $= 1$ if $=$ and 0 if \neq . Applying Theorems 1 and 2, Corollary 1, and Slutsky's theorem, (20) follows. We next prove (21) for each fixed $\in \mathcal{T}$. Let $\tilde{\cdot}(\cdot) = \mu(\cdot) + \sum_{=1} \tilde{\cdot}(\cdot)$, where $\tilde{\cdot}$ is as defined in (4). Note that

$$|\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)| \leq |\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)| + |\tilde{\cdot}(\cdot) - \cdot(\cdot)|.$$

Lemma A.3 implies that $\tilde{\cdot}(\cdot) \rightarrow \cdot(\cdot)$ as $\rightarrow \infty$. For fixed \cdot , observing that $\hat{\cdot} \rightarrow \tilde{\cdot}$ as $\rightarrow \infty$, $\sup_{\in \mathcal{T}} |\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)| \rightarrow 0$ as $\rightarrow \infty$ by (12), (17), and Slutsky's theorem. This implies that for

given $\cdot, > 0$, there exists \cdot_0 such that for $\geq \cdot_0$, $\{|\tilde{\cdot}(\cdot) - \cdot(\cdot)| > \cdot/2\} \leq \cdot/2$. For each \cdot , there exists $\cdot_0(\cdot)$ such that for $\geq \cdot_0(\cdot)$, $\{|\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)| \geq \cdot/2\} \leq \cdot/2$. Thus for $\geq \cdot_0$ and $\geq \cdot_0(\cdot)$, $\{|\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)| \geq \cdot\} \leq \{|\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)| > \cdot/2\} + \{|\tilde{\cdot}(\cdot) - \cdot(\cdot)| \geq \cdot/2\} \leq \cdot$, which leads to (21).

P n 4

Under the Gaussian assumption, for any fixed ≥ 1 , from Section 2.4, we have $(\tilde{\cdot}, \cdot, \cdot) \sim \mathcal{N}(\mathbf{0}, \cdot)$. Observing (12), (17), and (20), $\lim_{\rightarrow \infty} \sup_{\in \mathcal{T}} |\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)| \rightarrow 0$. Because $\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot) = \hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot) + \tilde{\cdot}(\cdot) - \cdot(\cdot)$ for fixed \cdot , it follows that $\{\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)\} \xrightarrow{\mathcal{D}} \sim \mathcal{N}(\mathbf{0}, \cdot(\cdot, \cdot))$. Under condition (A7), letting $\rightarrow \infty$ leads to $\xrightarrow{\mathcal{D}} \sim \mathcal{N}(\mathbf{0}, \cdot(\cdot, \cdot))$. From the Karhunen-Loève theorem, $|\cdot(\cdot) - \cdot(\cdot)| \rightarrow 0$, as $\rightarrow \infty$. Therefore, $\lim_{\rightarrow \infty} \lim_{\rightarrow \infty} \{\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)\} \xrightarrow{\mathcal{D}} \cdot$. From Theorems 1 and 2, it can be shown that $\hat{\cdot}(\cdot, \cdot) \rightarrow \cdot(\cdot, \cdot)$ as $\rightarrow \infty$. Under condition (A7), it follows that $\lim_{\rightarrow \infty} \lim_{\rightarrow \infty} \hat{\cdot}(\cdot, \cdot) = \cdot(\cdot, \cdot)$ in probability. Applying Slutsky's theorem, (22) follows.

P n 5

We first prove

$$\left\{ \sup_{\in \mathcal{T}} \frac{|\tilde{\cdot}(\cdot) - \cdot(\cdot)|}{\sqrt{\cdot(\cdot, \cdot)}} \leq \sqrt{\cdot^2, \cdot, 1-} \right\} \geq 1 - \cdot. \tag{A.8}$$

It is obvious that $\tilde{\cdot}(\cdot) - \cdot(\cdot) = \mathbf{v} \cdot (\tilde{\cdot}, \cdot, \cdot)$. Due to orthogonality, $\mathcal{F} = \{\cdot, \cdot, \cdot : \in \mathcal{T}\}$ is a \cdot -dimensional compact set. Because \cdot is positive definite, there exists a \times nonsingular matrix \mathbf{U} such that $\mathbf{U} \cdot \mathbf{U} = \mathbf{I}$. Let $\cdot_1 = \mathbf{U} \cdot$, and $\tilde{\cdot}_1 = \mathbf{U} \tilde{\cdot}$; then $(\tilde{\cdot}_1, \cdot_1) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. This leads to $(\tilde{\cdot}_1, \cdot_1) \sim \cdot^2$ and $\{(\tilde{\cdot}_1, \cdot_1) | (\tilde{\cdot}_1, \cdot_1)\} = 1 - \cdot$. We use the following result, known from linear algebra.

For a fixed \cdot -vector \mathbf{x} and a constant > 0 , $\mathbf{x} \cdot \mathbf{x} \leq 2$ if and only if $|\mathbf{a} \cdot \mathbf{x}| \leq \sqrt{\mathbf{a} \cdot \mathbf{a}}$, for all $\mathbf{a} \in \mathfrak{R}$.

Hence $\{|\mathbf{a} \cdot (\tilde{\cdot}_1, \cdot_1)| \leq \sqrt{\cdot^2, \cdot, 1-} \mathbf{a} \cdot \mathbf{a} : \text{for all } \mathbf{a} \in \mathfrak{R}\} = 1 - \cdot$. Let $\mathcal{E} = \{\mathbf{a} \in \mathfrak{R} : \cdot, \cdot, \cdot = \mathbf{U} \cdot, \cdot \in \mathcal{T}\}$, which is a compact subset of \mathfrak{R} .

Then $\{|\mathbf{a} \cdot (\tilde{\cdot}_1, \cdot_1)| \leq \sqrt{\cdot^2, \cdot, 1-} \mathbf{a} \cdot \mathbf{a} : \text{for all } \mathbf{a} \in \mathcal{E}\} \geq 1 - \cdot$, that is,

$$\begin{aligned} &\left\{ \left| \cdot, \tilde{\cdot}, \cdot, \cdot \right| \right. \\ &\quad \left. \leq \sqrt{\cdot^2, \cdot, 1-} \cdot, \mathbf{U}^{-1}(\mathbf{U} \cdot)^{-1} \cdot, : \text{for all } \cdot \in \mathcal{T} \right\} \\ &\geq 1 - \cdot. \end{aligned}$$

Observing that $\mathbf{U} \cdot \mathbf{U} = \mathbf{I}$, (A.8) follows.

To prove (23), note that

$$\begin{aligned} &\sup_{\in \mathcal{T}} \frac{|\hat{\cdot}(\cdot) - \cdot(\cdot)|}{\sqrt{\cdot(\cdot, \cdot)}} \\ &\leq \left(\sup_{\in \mathcal{T}} \frac{|\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)|}{\sqrt{\cdot(\cdot, \cdot)}} + \sup_{\in \mathcal{T}} \frac{|\tilde{\cdot}(\cdot) - \cdot(\cdot)|}{\sqrt{\cdot(\cdot, \cdot)}} \right) \sup_{\in \mathcal{T}} \sqrt{\frac{\cdot(\cdot, \cdot)}{\hat{\cdot}(\cdot, \cdot)}}. \end{aligned}$$

Let $\cdot = \sup_{\in \mathcal{T}} |\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)|/\sqrt{\cdot(\cdot, \cdot)}$, $\cdot = \sup_{\in \mathcal{T}} |\tilde{\cdot}(\cdot) - \cdot(\cdot)|/\sqrt{\cdot(\cdot, \cdot)}$, and $\cdot = \sup_{\in \mathcal{T}} \sqrt{\cdot(\cdot, \cdot)}/\hat{\cdot}(\cdot, \cdot)$. Because (\cdot, \cdot) is a continuous positive definite function on the bounded interval \mathcal{T} , it is bounded from above and below, say $0 < \cdot \leq (\cdot, \cdot) \leq < \infty$. Because $\sup_{\in \mathcal{T}} |\hat{\cdot}(\cdot) - \tilde{\cdot}(\cdot)| \rightarrow 0$ as $\rightarrow \infty$, we have $\cdot \rightarrow 0$ as $\rightarrow \infty$. In the proof of (22), we established that $\hat{\cdot}(\cdot, \cdot) \rightarrow (\cdot, \cdot)$, as $\rightarrow \infty$, implying that $\cdot \rightarrow 1$. We now show that

$$\lim_{\rightarrow \infty} \left\{ (\cdot + \cdot) \geq (\cdot + \sqrt{\cdot^2, \cdot, 1-}) (1 + \cdot) \right\} \leq \cdot. \tag{A.9}$$

Note that

$$\begin{aligned} & \{(\hat{c}_n + \hat{c}_n) \geq (\sqrt{2} + \sqrt{2})\} \\ & \subseteq \{(\hat{c}_n + \hat{c}_n) \geq (\sqrt{2} + \sqrt{2})\} \cup \{\hat{c}_n \geq (\sqrt{2} + \sqrt{2})\} \\ & \subseteq \{\hat{c}_n \geq \sqrt{2}\} \cup \{\hat{c}_n \geq \sqrt{2}\} \cup \{\hat{c}_n \geq (\sqrt{2} + \sqrt{2})\}. \end{aligned}$$

Because $\hat{c}_n \rightarrow 0$ and $\hat{c}_n \rightarrow 1$ as $n \rightarrow \infty$, for sufficiently large n , $(\hat{c}_n \geq \sqrt{2}) \leq \sqrt{3}$ and $(\hat{c}_n - 1 \geq \sqrt{2}) \leq \sqrt{3}$. We have shown that $(\hat{c}_n \geq \sqrt{2} + \sqrt{2}) \leq \sqrt{3}$ in (A.8). This implies (A.9), and then (23), by letting $n \rightarrow \infty$.

P C 2

There exists a $n \times n$ matrix \mathbf{Q} with rank $\leq n$ such that \mathcal{F} is spanned by the column vectors of \mathbf{Q} . Letting $\tilde{\mathbf{c}}_n = \mathbf{Q}^{-1} \mathbf{c}_n$ and $\tilde{\mathbf{c}}_n = \mathbf{Q}^{-1} \mathbf{c}_n$, for any $\mathbf{l} \in \mathcal{A}$, where $\mathcal{A} \subseteq \mathfrak{R}$ is a linear space with dimension n , there exists a vector $\tilde{\mathbf{l}} \in \mathfrak{R}$ such that $\mathbf{l} = \mathbf{Q} \tilde{\mathbf{l}}$. Then

$$\tilde{\mathbf{l}}_1, \dots, \tilde{\mathbf{l}}_n = \tilde{\mathbf{c}}_n - \tilde{\mathbf{c}}_n, \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{-1} \mathbf{Q}^{-1}).$$

Because \mathbf{Q} is of rank n and $\tilde{\mathbf{c}}_n$ is positive definite, which implies that $\mathbf{Q}^{-1} \mathbf{Q}^{-1}$ is also positive definite, there exists a nonsingular $n \times n$ matrix \mathbf{P} such that $\mathbf{PQ}^{-1} \mathbf{Q}^{-1} \mathbf{P} = \mathbf{I}$, where \mathbf{I} is the $n \times n$ identity matrix. Letting $\tilde{\mathbf{b}}_n = \mathbf{P} \tilde{\mathbf{c}}_n$ and $\tilde{\mathbf{b}}_n = \mathbf{P} \tilde{\mathbf{c}}_n$, we have $(\tilde{\mathbf{b}}_n - \tilde{\mathbf{b}}_n) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, that is, $(\tilde{\mathbf{b}}_n - \tilde{\mathbf{b}}_n) | (\tilde{\mathbf{b}}_n - \tilde{\mathbf{b}}_n) \sim \chi^2_n$. Therefore, $\{(\tilde{\mathbf{b}}_n - \tilde{\mathbf{b}}_n) | (\tilde{\mathbf{b}}_n - \tilde{\mathbf{b}}_n) \leq \sqrt{2}\} = 1 - \alpha$. Applying Lemma A.4, we obtain $\{|\mathbf{a} | (\tilde{\mathbf{b}}_n - \tilde{\mathbf{b}}_n) \leq \sqrt{2}\} = 1 - \alpha$ for all $\mathbf{a} \in \mathfrak{R}$. Because \mathbf{P} is nonsingular and \mathbf{Q} is of rank n , there exists $\tilde{\mathbf{l}} \in \mathfrak{R}$ and $\mathbf{l} \in \mathcal{A}$, such that $\mathbf{l} = \mathbf{P} \tilde{\mathbf{l}}$ and $\mathbf{l} = \mathbf{Q} \tilde{\mathbf{l}}$. If \mathbf{a} takes all values in \mathfrak{R} , then \mathbf{l} will also take all values in \mathcal{A} , that is,

$$\begin{aligned} & \{|\mathbf{l} | (\tilde{\mathbf{b}}_n - \tilde{\mathbf{b}}_n) \leq \sqrt{2}\} \\ & \leq \sqrt{2} | \mathbf{l} | (\mathbf{PQ}^{-1} \mathbf{Q}^{-1} \mathbf{P})^{-1} \mathbf{l} = 1 - \alpha. \end{aligned}$$

Because $\mathbf{PQ}^{-1} \mathbf{Q}^{-1} \mathbf{P} = \mathbf{I}$, the result (24) follows.

REFERENCE

Berkey, C. S., Laird, N. M., Valadian, I., and Gardner, J. (1991), "Modeling Adolescent Blood Pressure Patterns and Their Prediction of Adult Pressures," *Journal of the American Statistical Association*, 86, 1005–1018.

Besse, P., Cardot, H., and Ferraty, F. (1997), "Simultaneous Nonparametric Regression of Unbalanced Longitudinal Data," *Journal of Nonparametric Statistics*, 11, 255–270.

Besse, P., and Ramsay, J. O. (1986), "Principal Components Analysis of Sampled Functions," *Journal of the Royal Statistical Society B*, 48, 285–311.

Bhattacharya, P. K., and Müller, H. G. (1993), "Asymptotics for Nonparametric Regression," *Journal of Nonparametric Statistics*, 5, 420–441.

Boente, G., and Fraiman, R. (2000), "Kernel-Based Functional Principal Components," *Journal of Nonparametric Statistics*, 14, 335–345.

Boulanar, J., Ferré, L., and Vieu, P. (1993), "Growth Curves: A Two-Stage Nonparametric Approach," *Journal of Nonparametric Statistics*, 7, 327–350.

Bosq, D. (1991), "Modelization, Nonparametric Estimation and Prediction for Continuous Time Processes," in *Nonparametric Statistics*, ed. G. Roussas, Dordrecht, The Netherlands: Kluwer Academic, pp. 509–529.

Capra, W. B., and Müller, H. G. (1997), "An Accelerated-Time Model for Response Curves," *Journal of Nonparametric Statistics*, 11, 72–83.

Cardot, H., Ferraty, F., and Sarda, P. (1999), "Functional Linear Model," *Journal of Nonparametric Statistics*, 13, 11–22.

Castro, P. E., Lawton, W. H., and Sylvestre, E. A. (1986), "Principal Modes of Variation for Processes With Continuous Sample Curves," *Journal of the American Statistical Association*, 81, 329–337.

Courant, R., and Hilbert, D. (1953), *Methods of Mathematical Physics*, New York: Wiley.

Dauxois, J., Pousse, A., and Romain, Y. (1982), "Asymptotic Theory for the Principal Component Analysis of a Vector Random Function: Some Applications to Statistical Inference," *Journal of Nonparametric Statistics*, 12, 136–154.

Diggle, P. J., Liang, K. Y., and Zeger, S. L. (1994), *Generalized Linear Models*, Oxford, U.K.: Oxford University Press.

Fan, J., and Gijbels, I. (1996), *Functional Nonparametric Statistics*, London: Chapman and Hall.

Fan, J., and Zhang, J. T. (2000), "Two-Step Estimation of Functional Linear Models With Applications to Longitudinal Data," *Journal of Nonparametric Statistics*, Ser. B, 62, 303–322.

Ferré, L. (1995), "Improvement of Some Multivariate Estimates by Reduction of Dimensionality," *Journal of Nonparametric Statistics*, 9, 147–162.

James, G., Hastie, T. G., and Sugar, C. A. (2001), "Principal Component Models for Sparse Functional Data," *Journal of Nonparametric Statistics*, 15, 587–602.

James, G., and Sugar, C. A. (2003), "Clustering for Sparsely Sampled Functional Data," *Journal of Nonparametric Statistics*, 17, 397–408.

Jones, M. C., and Rice, J. (1992), "Displaying the Important Features of Large Collections of Similar Curves," *Journal of Nonparametric Statistics*, 6, 140–145.

Kaslow, R. A., Ostrow, D. G., Detels, R., Phair, J. P., Polk, B. F., and Rinaldo, C. R. (1987), "The Multicenter AIDS Cohort Study: Rationale, Organization and Selected Characteristics of the Participants," *American Journal of Epidemiology*, 126, 310–318.

Kneip, A. (1994), "Nonparametric Estimation of Common Regressors for Similar Curve Data," *Journal of Nonparametric Statistics*, 22, 1386–1472.

Kneip, A., and Utikal, K. (2001), "Inference for Density Families Using Functional Principal Component Analysis," *Journal of Nonparametric Statistics*, 15, 519–532.

Lin, X., and Carroll, R. J. (2000), "Nonparametric Function Estimation for Clustered Data When the Predictor Is Measured Without/With Error," *Journal of Nonparametric Statistics*, 14, 520–534.

Mardia, K. V., Kent, J. T., and Bibby, J. M. (1979), *Biometrika*, London: Academic Press.

Müller, H. G., and Prewitt, K. (1993), "Multiparameter Bandwidth Processes and Adaptive Surface Smoothing," *Journal of Nonparametric Statistics*, 7, 1–21.

Ramsay, J., and Silverman, B. (1997), *Functional Data Analysis*, New York: Springer-Verlag.

Rao, C. R. (1958), "Some Statistical Methods for Comparison of Growth Curves," *Journal of the American Statistical Association*, 53, 1–17.

Rice, J., and Silverman, B. (1991), "Estimating the Mean and Covariance Structure Nonparametrically When the Data Are Curves," *Journal of Nonparametric Statistics*, Ser. B, 53, 233–243.

Rice, J., and Wu, C. (2000), "Nonparametric Mixed Effects Models for Unequally Sampled Noisy Curves," *Journal of Nonparametric Statistics*, 14, 253–259.

Shi, M., Weiss, R. E., and Taylor, J. M. G. (1996), "An Analysis of Paediatric CD4 Counts for Acquired Immune Deficiency Syndrome Using Flexible Random Curves," *Journal of Nonparametric Statistics*, 10, 151–163.

Shibata, R. (1981), "An Optimal Selection of Regression Variables," *Journal of Nonparametric Statistics*, 15, 45–54.

Silverman, B. (1996), "Smoothed Functional Principal Components Analysis by Choice of Norm," *Journal of Nonparametric Statistics*, 10, 45–54.

Spellman, P. T., Sherlock, G., Zhang, M. Q., Tyer, V. R., Eisen, M. B., Brown, P. O., Botstein, D., and Futcher, B. (1998), "Comprehensive Identification of Cell Cycle-Regulated Genes of the Yeast by Microarray Hybridization," *Science*, 280, 3273–3297.

Staniswalis, J. G., and Lee, J. J. (1998), "Nonparametric Regression Analysis of Longitudinal Data," *Journal of Nonparametric Statistics*, 12, 1403–1418.

Wu, C., and Chiang, C. (2000), "Kernel Smoothing on Varying Coefficient Models With Longitudinal Dependent Variable," *Journal of Nonparametric Statistics*, 14, 433–456.

Yao, F., Müller, H. G., Clifford, A. J., Dueker, S. R., Follett, J., Lin, Y., Buchholz, B. A., and Vogel, J. S. (2003), "Shrinkage Estimation for Functional Principal Component Scores With Application to the Population Kinetics of Plasma Folate," *Journal of Nonparametric Statistics*, 17, 676–685.