

# Functional Data Analysis for Sparse Longitudinal Data

Fang YAO, Hans-Georg MÜLLER, and Jane-Ling WANG  
Journal of the American Statistical Association

Wangfei Wang

April 27, 2020

# What is the main problem the paper trying to address?

## ► Introduction

Functional principal components (FPC) analysis characterizes the dominant mode of variation around an overall mean trend function, and therefore is popular in longitudinal data analysis.

## ► Limitations of available models:

- Cannot deal with infrequent, irregularly-spaced repeated measures.
- Some kernel-based FPC analysis cannot be approximated by the usual integration method.
- Linear mixed models or reduced-rank mixed effects models using B-splines to model the individual curves with random coefficients [refs] are too complex, and the asymptotic properties of the estimated components were not investigated.

# What is the proposed solution?

- ▶ They proposed a version of FPC analysis, in which they framed the FPC scores as conditional expectations. And thus they coined this method “principal components analysis through conditional expectation (**PACE**)”.
- ▶ **Contributions of the paper**
  - ▶ In the model, they took into account the additional measurement errors.
  - ▶ They derived the asymptotic consistency properties.
  - ▶ They derived the asymptotic distribution needed for obtaining point-wise confidence intervals for individual trajectories.

# Innovation

- ▶ The proposed conditional model is designed for sparse and irregular longitudinal data.
- ▶ Under Gaussian assumptions, the authors showed that estimation of individual FPC scores are the best prediction; and under non-Gaussian assumption, they provide estimates for best linear prediction.
- ▶ One-curve-leave-out cross-validation was proposed to choose auxiliary parameters.
- ▶ Akaike information criterion (AIC) was used for faster computation to select eigenfunctions.

# METHOD: PACE

- ▶ Model with Measurement Errors
- ▶ Estimation of the Model Components
- ▶ Functional Principal Components Analysis Through Conditional Expectation
- ▶ Asymptotic Confidence Bands for Individual Trajectories
- ▶ Selection of the Number of Eigenfunctions

## Methods: Model with Measurement Errors

Assume: 1) Trajectories are independent realizations of a smooth random function with unknown mean  $EX(t) = \mu(t)$  and covariance  $cov(X(s), X(t)) = G(s, t)$ , where domain of  $X(\cdot)$  is  $\mathcal{T}$ .

2)  $G$  has an orthogonal expansion in terms of eigenfunction  $\phi_k$  and eigenvalues  $\lambda_k$ :  $G(s, t) = \sum_k \lambda_k \phi_k(s) \phi_k(t)$ ,  $t, s \in \mathcal{T}$ , where  $\lambda_1 \geq \lambda_2 \geq \dots$ .

Model:

$$Y_{ij} = X_i(T_{ij}) + \epsilon_{ij} \quad (1)$$

$$= \mu(T_{ij}) + \sum_{k=1}^{\infty} \xi_{ik} \phi_k(T_{ij}) + \epsilon_{ij}, \quad T_{ij} \in \mathcal{T} \quad (2)$$

where  $E\epsilon_{ij} = 0$ ,  $\text{var}(\epsilon_{ij}) = \sigma^2$ .

$Y_{ij}$  is the  $j$ th observation of the random function  $X(\cdot)$ , and  $\epsilon_{ij}$  is the measurement errors that are iid and are independent of random coefficients  $\xi_{ik}$ , where  $i = 1, \dots, n; j = 1, \dots, N_i; k = 1, 2, \dots$

# Methods: Estimation of the Model Components

- Estimation of mean function  $\mu$

Minimizing the following equation (3) respect to  $\beta_0$  and  $\beta_1$

$$\sum_{i=1}^n \sum_{j=1}^{N_i} \kappa_1\left(\frac{T_{ij} - t}{h_\mu}\right) \{Y_{ij} - \beta_0 - \beta_1(t - T_{ij})\}^2 \quad (3)$$

where  $\kappa_1$  is a kernel function:  $\mathbb{R} \rightarrow \mathbb{R}$ .

Then estimation of mean function  $\mu$  can be obtained:

$$\hat{\mu}(t) = \hat{\beta}_0(t)$$

# Methods: Estimation of the Model Components

- Estimation of measurement errors  $\sigma^2$

$$\hat{\sigma}^2 = \frac{2}{|\mathcal{T}|} \int_{\mathcal{T}_1} \{\hat{V}(t) - \tilde{G}(t)\} dt \quad (4)$$

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

where  $|\mathcal{T}|$  is the length of  $\mathcal{T}$ ,  $\mathcal{T}_\infty = [\inf\{x : x \in \mathcal{T}\} + |\mathcal{T}|/4]$ ,

$\tilde{G}$  is the diagonal of the surface estimate

$\hat{V}(t)$  is a local linear smoother focusing on diagonal values  $\{G(t, t) + \sigma^2\}$ .

Estimation procedures for  $\tilde{G}$ :

$\hat{G}(s, t) \rightarrow$  surface estimate  $\bar{G}(s, t) \rightarrow \tilde{G}(t) = \bar{G}(0, t/\sqrt{(2)})$ ,  
where  $G(s, t)$  is the “raw covariance”  $\text{cov}(X(s), X(t))$ .



# Methods: Estimation of the Model Components

- Estimation of eigenfunctions and eigenvalues  $\phi_k$  and  $\lambda_k$   
Solutions  $\phi_k$  and  $\lambda_k$  of the following eigenequation:

$$\int_{\mathcal{T}} \hat{G}(s, t) \hat{\phi}_k(s) ds = \hat{\lambda}_k \hat{\phi}_k(t) \quad (5)$$

where the  $\hat{\phi}_k$  are subject to  $\int_{\mathcal{T}} \hat{\phi}_k(t)^2 dt = 1$  and  $\int_{\mathcal{T}} \hat{\phi}_k(t) \times \hat{\phi}_m(t) dt = 0$  for  $m < k$ .

# Methods: Functional Principal Components Analysis Through Conditional Expectation

- Under the assumption that  $\xi_{ik}$  and  $\epsilon_{ij}$  are jointly Gaussian:

$$\hat{\xi}_{ik} = \hat{E}[\xi_{ik} | \tilde{\mathbf{Y}}_i] = \hat{\lambda}_k \hat{\phi}_{ik}^T \hat{\Sigma}_{\mathbf{Y}_i}^{-1} (\tilde{\mathbf{Y}}_i - \hat{\boldsymbol{\mu}}_i) \quad (6)$$

where  $(\hat{\Sigma}_{\mathbf{Y}_i})_{j,l} = \hat{G}(T_{ij}, T_{il}) + \sigma^2 \delta_{jl}$  is the  $(j, l)$ th element of  $\hat{\Sigma}_{\mathbf{Y}_i}$ . Under the Gaussian assumption, the  $\hat{\xi}_{ik} = E[\xi_{ik} | \tilde{\mathbf{Y}}_i]$  is the best prediction of the FPC score.

- The prediction for the trajectory  $X_i(t)$  for the  $i$ th subject using the first  $K$  eigenfunctions is then:

$$\hat{X}_i^K(t) = \hat{\mu}(t) + \sum_{k=1}^K \hat{\xi}_{ik} \hat{\phi}_k(t) \quad (7)$$

- In simulation result, the authors showed that this proposed model is also robust when the Gaussian assumption does not hold.

## Methods: Asymptotic Confidence Bands for Individual Trajectories

- ▶ The  $(1 - \alpha)$  asymptotic simultaneous confidence bands for  $X_i(t)$  can be obtained:

$$\hat{X}_i^K(t) \pm \sqrt{\chi_{K,1-\alpha}^2 \hat{\phi}_{K,t}^T \hat{\Omega}_K \hat{\phi}_{K,t}} \quad (8)$$

where  $\chi_{K,1-\alpha}^2$  is the  $100(1 - \alpha)$ th percentile of the chi-squared distribution with  $K$  degrees of freedom.

- ▶ For all linear combinations of the FPC scores, the authors proved that they could be obtained by:

$$\mathbf{l}^T \boldsymbol{\xi}_{K,i} \in \mathbf{l}^T \hat{\boldsymbol{\xi}}_{K,i} \pm \sqrt{\chi_{d,1-\alpha}^2 \mathbf{l}^T \hat{\boldsymbol{\Omega}} \mathbf{l}} \quad (9)$$

with approximate probability  $(1 - \alpha)$ , where  $\mathbf{l} \in \mathcal{A}$ ,  $\mathcal{A} \subseteq \mathbb{R}^K$  is a linear space with dimension  $d \leq K$ .

## Methods: Selection of the Number of Eigenfunctions

- ▶ Choose the number of eigenfunctions  $K$  that minimizes the cross-validation score:

$$CV(K) = \sum_{i=1}^n \sum_{j=1}^{N_i} \{Y_{ij} - \hat{Y}_i^{(-i)}(T_{ij})\}^2 \quad (10)$$

where  $\hat{Y}_i^{(-i)}(t) = \hat{\mu}^{(-i)}(t) + \sum_{k=1}^K \hat{\xi}_{ik}^{(-i)}(t) \hat{\phi}_k^{(-i)}(t)$ .

- ▶ AIC-type criteria was found to be more computationally efficient. A pseudo-Gaussian log-likelihood was generated:

$$\begin{aligned} \hat{L} = \sum_{i=1}^n \left\{ -\frac{N_i}{2} \log(2\pi) - \frac{N_i}{2} \log \hat{\sigma}^2 - \right. \\ \left. \frac{1}{2\hat{\sigma}^2} (\tilde{\mathbf{Y}}_i - \hat{\boldsymbol{\mu}}_i - \sum_{k=1}^K \hat{\xi}_{ik} \hat{\boldsymbol{\phi}}_{ik})^T \times (\tilde{\mathbf{Y}}_i - \hat{\boldsymbol{\mu}}_i - \sum_{k=1}^K \hat{\xi}_{ik} \hat{\boldsymbol{\phi}}_{ik}) \right\} \end{aligned} \quad (11)$$

where  $AIC = -\hat{L} + K$ .

# Asymptotic Properties

- ▶  $\sup_{t \in \mathcal{T}} |\hat{\mu}(t) - \mu(t)| = O_p \left( \frac{1}{\sqrt{nh_{\mu}}} \right)$
- ▶  $\sup_{t, s \in \mathcal{T}} |\hat{G}(s, t) - G(s, t)| = O_p \left( \frac{1}{\sqrt{nh_G^2}} \right), |\hat{\lambda}_k - \lambda_k| = O_p \left( \frac{1}{\sqrt{nh_G^2}} \right),$   
 $\|\hat{\phi}_k - \phi_k\|_H = O_p \left( \frac{1}{\sqrt{nh_G^2}} \right), k \in \mathcal{T}'$
- ▶  $\sup_{t \in \mathcal{T}} |\hat{\phi}_k(t) - \phi_k(t)| = O_p \left( \frac{1}{\sqrt{nh_G^2}} \right), k \in \mathcal{T}'$
- ▶  $\lim_{n \rightarrow \infty} \hat{\xi}_{ik} = \tilde{\xi}_{ik}, \lim_{K \rightarrow \infty} \lim_{n \rightarrow \infty} \hat{X}_i^K(t) = \tilde{X}_i(t) \quad \forall t \in \mathcal{T} \text{ in probability.}$
- ▶  $\lim_{K \rightarrow \infty} \lim_{n \rightarrow \infty} P \left\{ \frac{\hat{X}_i^K(t) - X_i(t)}{\sqrt{\omega_K(t, t)}} \leq x \right\} = \Phi(x),$  where  $\Phi(x)$  is the standard Gaussian cdf.
- ▶  $\lim_{n \rightarrow \infty} P \left\{ \sup_{t \in \mathcal{T}} \frac{|\hat{X}_i^K(t) - X_i^K(t)|}{\sqrt{\omega_K(t, t)}} \leq \sqrt{\chi_{K, 1-\alpha}^2} \right\} \geq 1 - \alpha,$  where  $\chi_{K, 1-\alpha}^2$  is the  $1 - \alpha$ th percentile of the chi-squared distribution with  $K$  degrees of freedom.

# Simulation Studies

*Table 1. Results for FPC Analysis Using Conditional Expectation (CE, corresponding to PACE) and Integration (IN) Methods for 100 Monte Carlo Runs With  $N = 100$  Random Trajectories per Sample, Generated With Two Random Components*

$N = 100$ FPC		Normal			Mixture		
		MSE	ASE( $\xi_1$ )	ASE( $\xi_2$ )	MSE	ASE( $\xi_1$ )	ASE( $\xi_2$ )
Sparse	CE	1.33	.762	.453	1.30	.737	.453
	IN	2.32	1.58	.622	2.25	1.53	.631
Nonsparse	CE	.259	.127	.110	.256	.132	.105
	IN	.286	.159	.115	.286	.168	.114

NOTE: Shown are the averages of estimated mean squared prediction error, MSE, and average squared error, ASE( $\xi_k$ ),  $k = 1, 2$ , as described in Section 4. The number of components for each Monte Carlo run is chosen by the AIC criterion (11).

$$MSE = \sum_{i=1}^n \int_0^{10} \left\{ X_i(t) - \hat{X}_i^K(t) \right\}^2 dt / n$$

$$ASE(\xi_k) = \sum_{i=1}^n (\hat{\xi}_{ik} - \xi_{ik})^2 / n \quad k = 1, 2.$$

# Applications

## Objectives:

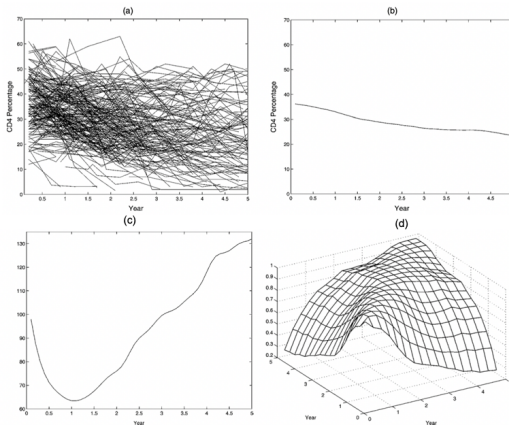
- ▶ To estimate the overall trend over time
- ▶ To study subject-specific variation patterns
- ▶ To uncover the dominant modes of variation
- ▶ To recover individual trajectories from sparse measurements.

## Datasets:

- ▶ Longitudinal CD4 Counts
- ▶ Yeast Cell Cycle Gene Expression Profiles

# Longitudinal CD4 Counts

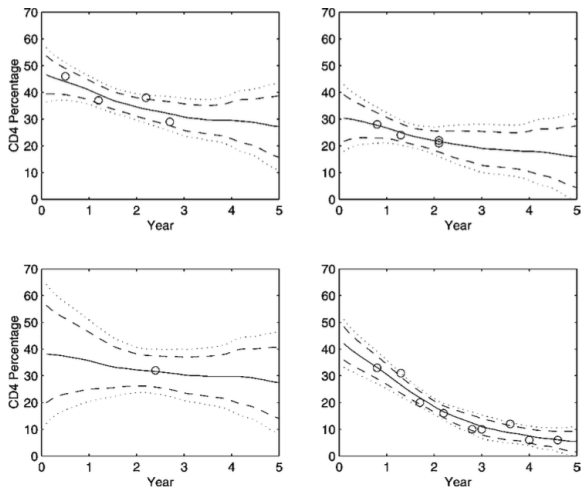
A cohort of 283 homosexual men who became HIV-positive between 1984 and 1991. CD4 counts and CD4 percentage were recorded over 5 years at their semiannual visits.



**Figure 1.** (a) Individual trajectories of CD4 percentage in 283 individuals. (b) Smooth estimate of the mean function. (c) Smooth estimate of the variance function for CD4 counts. (d) Smooth estimate of correlation function.



# Longitudinal CD4 Counts

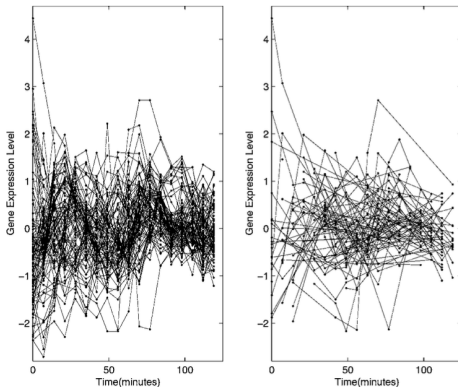


**Figure 2.** Observations (circles), predicted (solid lines) trajectories, and 95% pointwise (dashed lines) and simultaneous (dotted lines) bands for four randomly chosen individuals, for the CD4 count data.

# Yeast Cell Cycle Gene Expression Profiles

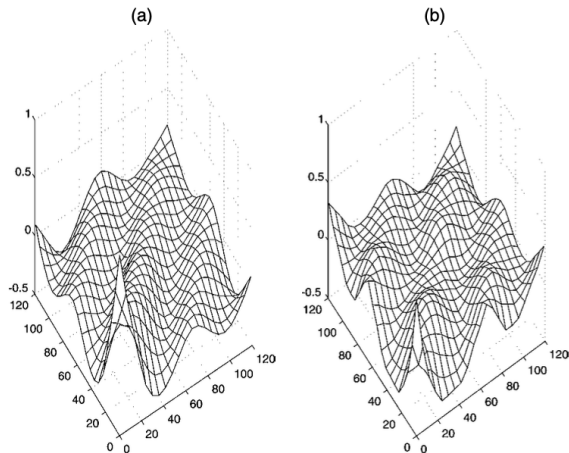
**Whole dataset:** The training set include 6178 genes with each gene expression profile consists of 18 data points, measured every 7 minutes in a span of 0 to 119 minutes.

**“Sparsified” dataset:** The authors artificially induced sparsity to the data by randomly selecting  $N_i \in (1 - 6)$  with equal probability, and then randomly select from the 18 recorded gene expression measurements (the median of the number of observations per gene expression profile is 3).



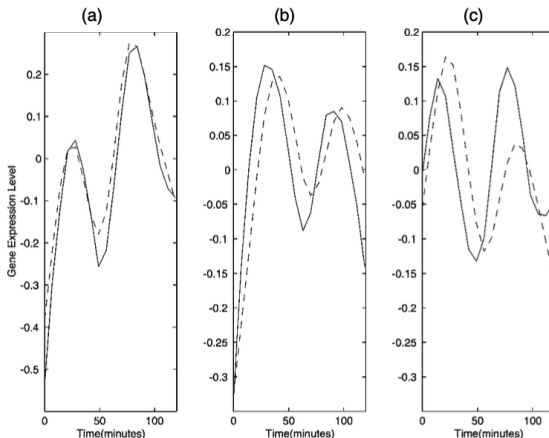
**Figure 3.** Left: the whole dataset. Right: the “sparsified” dataset.

# Yeast Cell Cycle Gene Expression Profiles



**Figure 4.** Smooth surface estimates  $\hat{G}$  of the covariance functions from the complete data (a) and from the sparsified data (b).

# Yeast Cell Cycle Gene Expression Profiles



**Figure 5.** Smooth estimates of the mean function (a), the first (b) and second (c) eigenfunctions, obtained from sparse (solid lines) and complete (dashed lines) data.

# Potential applications of the proposed method

- ▶ Longitudinal data with sparse measurements per subject.
- ▶ Longitudinal data with irregularly spaced measurements.
- ▶ To help impute missing data in longitudinal studies.
- ▶ PACE can also be used for regularly spaced data, which was also shown to be advantageous in this study.

Propose one or two possible topics/questions for future research in this area.

- ▶ Additive models are more flexible in practice. The authors could extend “PACE” to additive models. The challenge of additive model in functional study could be that the predictor set is not countable and the additive models may require many additive components.

$$E(Y|X) = \mu_Y + \sum_{k=1}^{\infty} f_k(\xi_k)$$

and

$$E(Y(t)|X) = \mu_Y(t) + \sum_{k=1}^{\infty} \sum_{m=1}^{\infty} f_{km}(\xi_k) \psi_m(t),$$

- ▶ Another area the authors can investigate is how to classify the functional data. In practice, it's very often that we need to compare two or more groups of longitudinal data.