

# Some selected topics in Functional Data Analysis

Juhyun Park

ENSIIE, Evry, France

Motivation

What is functional data

Exploring functional data

- Principal Component Analysis (PCA)

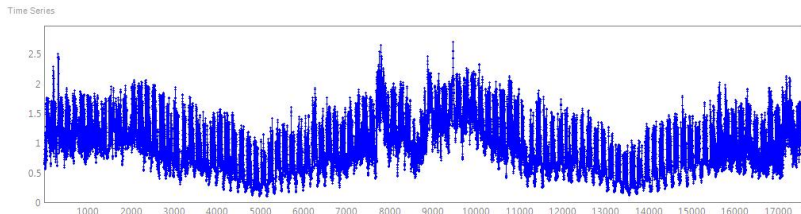
- Functional Principal Component Analysis

Smoothing methods

## Section 1

### Motivation

# Forecasting challenges

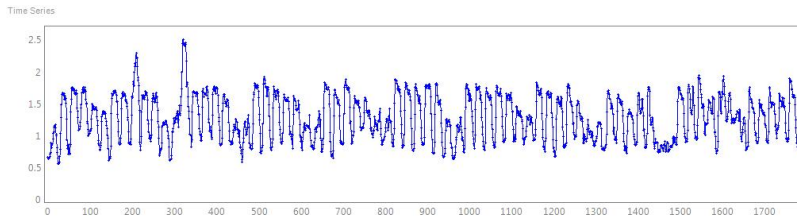


- ▶ Novel time series data in high frequency
  - ▶ Advances in data processing, recording and storage of vast amounts of data gives observations on previously unobserved periods, recorded at higher, near continuous sampling frequencies
- ▶ Time series exhibit novel patterns and pose new challenges in the analysis

## Examples of high frequency data

- ▶ Blood pressure of patients is constantly monitored and recorded over 24 hours
- ▶ EEG brain signals recorded continuously over a period of time
- ▶ Meteorological measurements (temperature, humidity..) monitored continuously
- ▶ Retail sales volumes store daily and intra-day sales observations or even transaction records (ticks)
- ▶ Bidding history at online auction sites such as eBay can be tracked continuously
- ▶ Stock market transactions are recorded (and executed) at millisecond intervals of tick-by-tick data

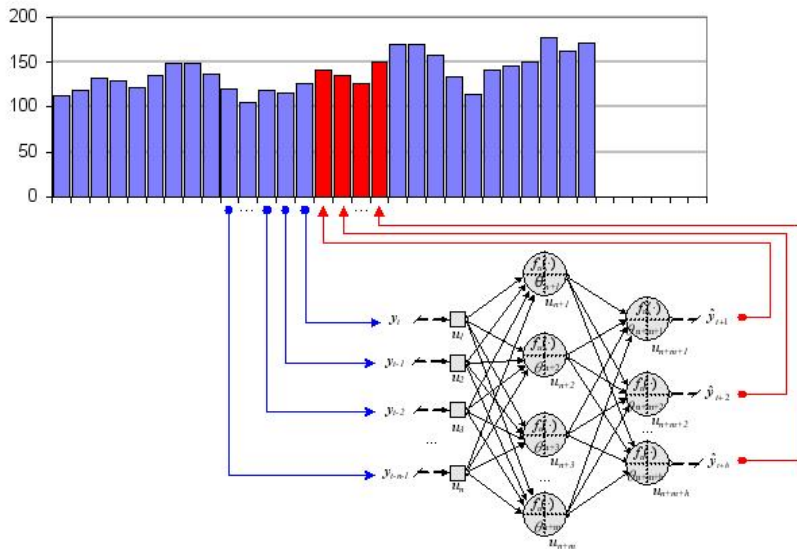
# Forecasting gas consumption - zoomed in



Model multiple seasonal patterns... for ARIMA /Neural Networks /Smoothing?

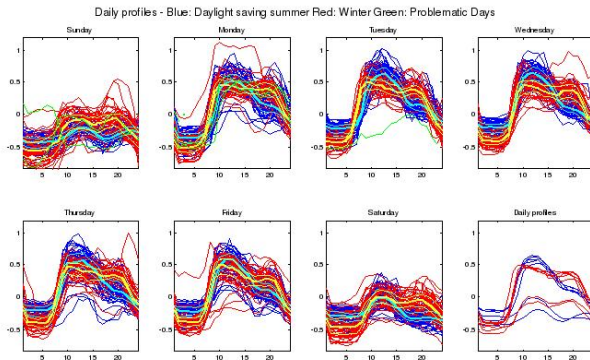
- ▶ Daylight savings in summer  $\neq$  winter?
- ▶ Leap years in AR-processes?
- ▶ Variable selection for high-frequency data?

## A standard approach - discrete time series



# Functional data view

Model multiple seasonal patterns... as functions!



- ▶ Model (smooth) transition between functions?
- ▶ Consider covariates (temperature, windspeed etc.) in transition



# Advantages of functional data view

- ▶ Enrich information from discrete data to replicates of functions
  - ▶ Periodic, non-stationary behaviour
  - ▶ Smooth change between measurements
- ▶ Guided analysis by focusing on recovering the common functional structure while taking into account the variations among the replicates.
- ▶ Time series forecasting is built on limited number of past observations whereas functional data view encourages us to look at all the past observations.
- ▶ Functional data facilitates an analysis of the smooth transitions between these different groups, possibly formulated as having hierarchical (multi-level) structure

## Section 2

### What is functional data

# Big picture: Statistical Learning

## Statistical framework

- ▶ Population: the set of all possible *units* of interest
- ▶ Sample: a (random) subset of population

From finite sample, we would like to draw conclusions about the population.

## Statistical analysis

- ▶ Data: multiple observation  $\{x_i, i = 1, \dots, n\}$
- ▶ Aim: understanding and characterizing variability in the sample of observations

# What is the unit of the analysis

Multiple observation  $\{X_1, \dots, X_n\}$ ,  $X_i$  characteristics of subject  $i$

- ▶ Number: univariate analysis

$X_i$  : commuting distance in km,  $x_i = 5.2$

- ▶ Vector: multivariate analysis

$X_i = (\text{distance(km)}, \text{age (y)}), \quad x_i = (5.2, 32)$

- ▶ Function: functional data analysis

$X_i = \text{commuting trajectory}, \quad x_i = ?$

- ▶ Shape analysis, image analysis, topological data analysis, manifold analysis ...

- ▶ *Object Oriented Data Analysis*

*The unit of analysis is the whole trajectory, not the finite number of available observations!*

# Functional data

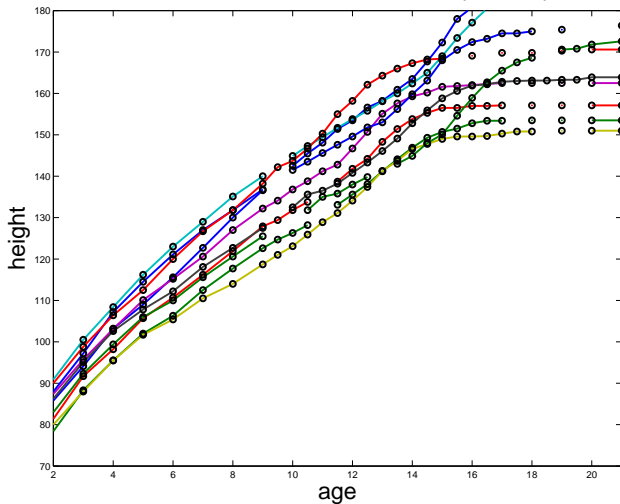
Most often, functional data refers to measurements on a curve but in a broader sense, it is also used to encompass images, tree-like objects and many other non-Euclidian objects arising in modern applications → *Object Oriented Data Analysis*

Observations from functional data:

- ▶ repeated measurements available from multiple subjects
- ▶ often densely observed, though sparse observations are also dealt with.
- ▶ often represent the underlying *continuous*, possibly smooth, (physical or biological) process

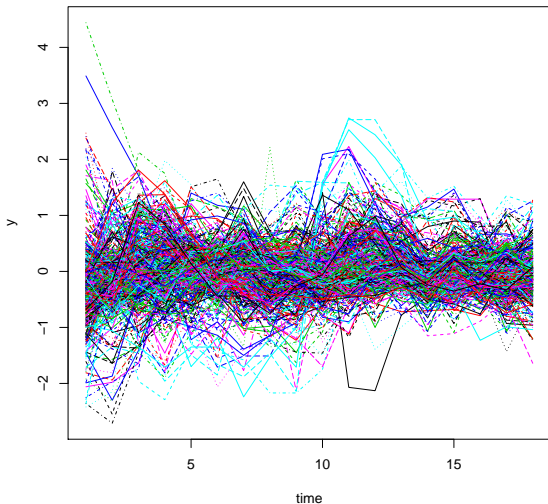
## Example: Growth curves

A sample of child's growth measurements (height) until 21 years



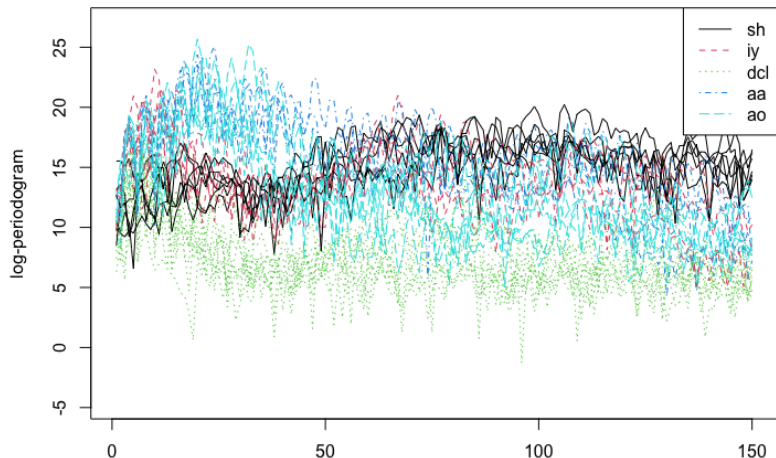
## Example: Genome's cell cycle regulation

A sample of gene expression (mRNA levels) data measured every 7 minutes during 119 minutes



## Example: Speech recognition

A sample of log-periodograms of speech recording of 32 ms for 5 classes of phonemes in the first 150 frequencies



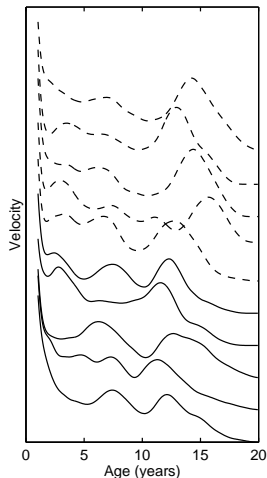
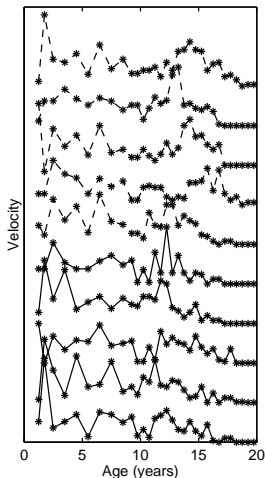


# Common characteristics of measurements

- ▶ only discrete measurements are available:  $X_j, j = 1, \dots, n$
- ▶ often replicates of functions are available:  
 $X_{ij}, i = 1, \dots, n; j = 1, \dots, m$
- ▶ sampling points may vary from one record to another;  
 $X_{ij}, i = 1, \dots, n; j = 1, \dots, n_i$
- ▶ measurement error may be present:  $Y_{ij} = X_{ij} + \varepsilon_{ij}$
- ▶ values reflect a smooth variation that could be assessed at *any* time or as often as desired:  $X_{ij} = f_i(t_{ij})$
- ▶ values are continuous in nature thus should be viewed as a *function*:  $f_i(\cdot)$  is continuous
- ▶ functional features such as derivatives could be of main interest:  $f_i(\cdot)$  differentiable

## Estimating velocity curves for trunk length

Velocities of trunk length for 5 boys (above) and 5 girls (below).  
left = raw velocities, right = kernel estimated velocities.



# Velocity estimation

Raw velocity:

$$\tilde{y}_j = \frac{y(t_j) - y(t_{j-1})}{t_j - t_{j-1}} \quad s_j = \frac{t_j + t_{j-1}}{2}$$

Estimated velocity:

$$\tilde{y}(t) = \text{smooth}\{(s_1, \tilde{y}_1), (s_2, \tilde{y}_2), \dots, (s_p, \tilde{y}_p)\} \quad \text{and evaluate at } t$$

# Function representation: uni-dimensional case

- ▶ Functions are infinite-dimensional objects:  $f : \mathcal{I} \rightarrow \mathbb{R}$
- ▶ Finite observations are available:  $f(t_1), \dots, f(t_n), t_j \in \mathcal{I}$ .
- ▶ Observations without error made on fine grid points: *numerical interpolation* on a finite grid

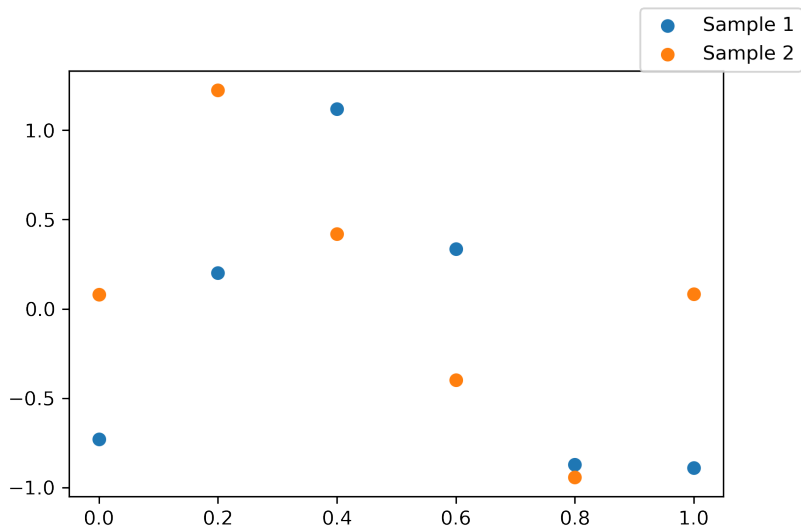
$$\mathbf{X} = (f(t_1), \dots, f(t_m))$$

- ▶ Measurement error in data:  $Y_1, \dots, Y_n$ , noisy observations of  $X_i$ :

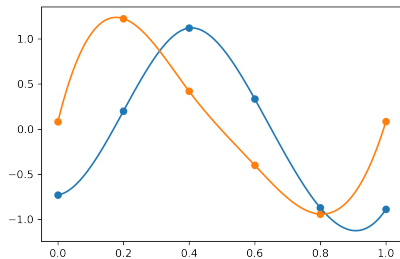
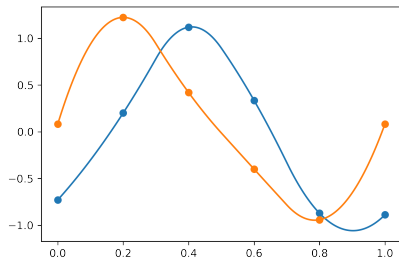
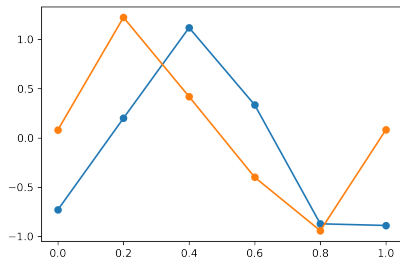
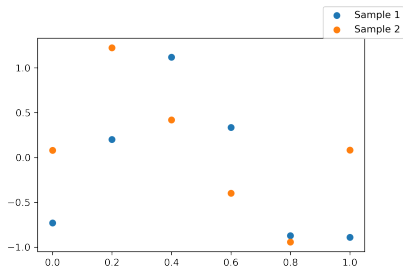
$$Y_i(t) = X_i(t) + \varepsilon_i(t)$$

where  $\varepsilon_i$ 's are independent of  $X_i$ 's and are independent and identically distributed zero mean stationary processes.  $\Rightarrow$  *statistical interpolation = smoothing*

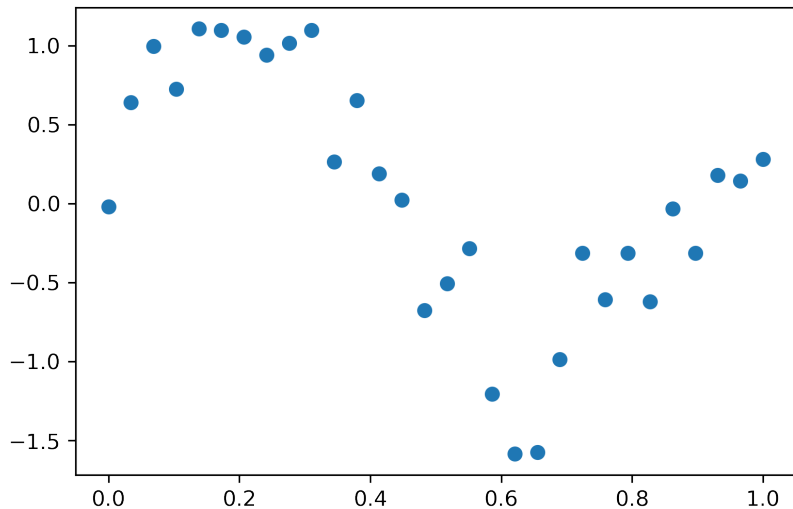
# Function representation: without noise



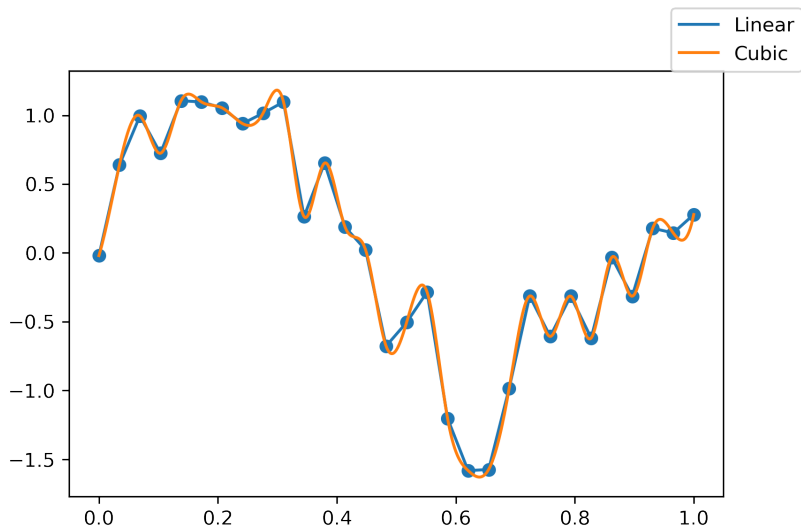
# Interpolation (linear, quadratic, cubic)



## Function representation: with noise

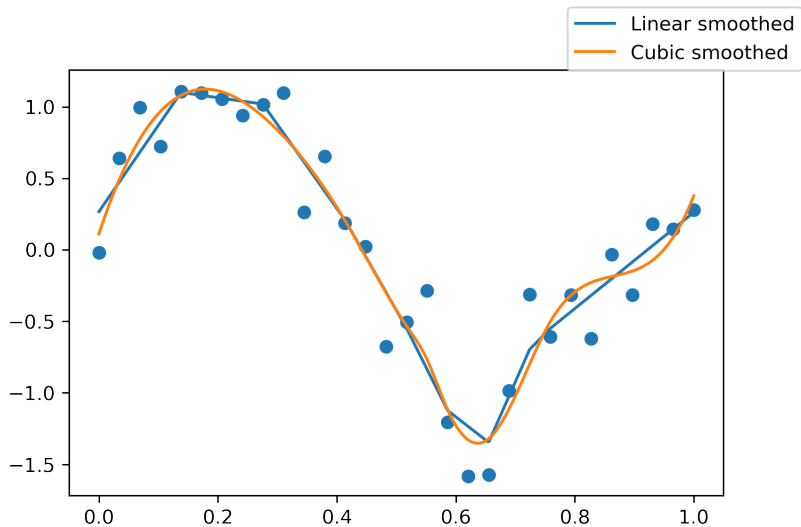


# Interpolation (linear, cubic)





# Statistical interpolation (linear, cubic)



# Vector vs Function

- ▶ Although a vector representation gives a convenient mean to represent/compute/manipulate function valued variables, they are not equivalent.
- ▶ Multivariate analysis deals with *fixed dimension* (number of available observations per curve), whereas functional data does not have a fixed dimension to begin with.
- ▶ Vector representation cannot account for the fact that future observations (test data) may not have the same representation as the available data (training data) for functional data
- ▶ For functional data, often the interest lies in derivatives of the functions!

# Overview of functional data analysis

## Features of data:

- ▶ similar pattern of variation among curves: share common structure or *common shape*
- ▶ complex relationship and complex variability
- ▶ interest in functional characteristics: nonlinear relationship, derivatives ...

## Aims of the analysis:

- ▶ understand and characterize the *common features* of the homogeneous population
- ▶ discriminate and classify *distinct* populations
- ▶ extract *maximal* information with an *efficient* representation

## Some references

- ▶ Ramsay, J. O. and Silverman, B. W. (1997, 2005) Functional Data Analysis
- ▶ Ramsay, J. O. and Silverman, B. W. (2002) Applied Functional Data Analysis: Methods and Case Studies
- ▶ Ramsay, J. O., Hooker, G., Gaves, S. (2005) Functional Data Analysis with R and Matlab
- ▶ Ferraty, F. and Vieu, P. (2006) Nonparametric Functional Data Analysis: Theory and Practice
- ▶ Hastie, T., Tibshirani, R. and Friedman, J. (2001, 2008) The Elements of Statistical Learning: Data Mining, Inference and Prediction

## Section 3

### Exploring functional data

# Multivariate data

Let  $X = (X_1, \dots, X_p)^\top \in \mathbb{R}^p$  be a  $p$ -dimensional random vector with mean  $\mu$  and variance  $\Sigma$ . Then  $\Sigma$  is a  $p \times p$  symmetric and positive-definite (excluding constant variable case) matrix.

Data:

$$X_i = (X_{i1}, \dots, X_{ip})^\top \quad i = 1, \dots, n$$

Mean  $\mu = (\mu_1, \dots, \mu_p)$  and  $\Sigma = (\sigma_{ij})$  where

$$E[X_{ij}] = \mu_j \quad \text{cov}(X_{ij}, X_{ik}) = \sigma_{ij}$$

$$\Sigma = \text{cov}(X) = \begin{pmatrix} \text{var}(X_1) & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_p) \\ \vdots & \vdots & \vdots & \vdots \\ \text{cov}(X_p, X_1) & \text{cov}(X_p, X_2) & \dots & \text{var}(X_p) \end{pmatrix}$$

# Sample mean and variance

- ▶ Sample mean for  $j$ th component:  $\bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}$ .
- ▶ Sample mean vector:

$$\bar{X} = (\bar{X}_1, \dots, \bar{X}_p)^\top$$

- ▶ Residual:  $\tilde{X}_i = X_i - \bar{X}$
- ▶ Sample covariance:

$$\begin{aligned}\Sigma_n &= \frac{1}{n} \sum_{i=1}^n \tilde{X}_i \tilde{X}_i^\top \\ &= \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n \tilde{X}_{i1}^2 & \frac{1}{n} \sum_{i=1}^n \tilde{X}_{i1} \tilde{X}_{i2} & \dots & \frac{1}{n} \sum_{i=1}^n \tilde{X}_{i1} \tilde{X}_{ip} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{1}{n} \sum_{i=1}^n \tilde{X}_{ip} \tilde{X}_{i1} & \frac{1}{n} \sum_{i=1}^n \tilde{X}_{ip} \tilde{X}_{i2} & \dots & \frac{1}{n} \sum_{i=1}^n \tilde{X}_{ip}^2 \end{pmatrix}\end{aligned}$$

# Functional random variable

- Population: let  $X \in L^2(\mathcal{I})$  be the functional random variable on  $\mathcal{I}$  with

$$\mathbb{E}[X(t)] = \mu(t) \quad \text{cov}(X(s), X(t)) = \Gamma(s, t)$$

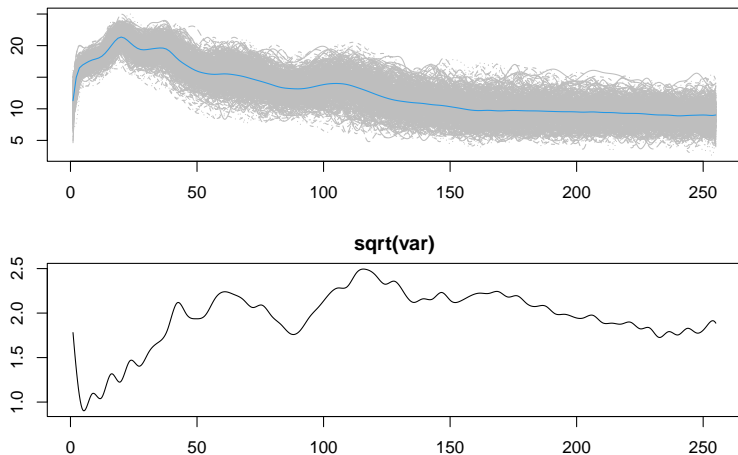
- Sample:  $X_1, \dots, X_n$ , assume that these are independent and identically distributed as  $X$ .
- Sample mean:  $\bar{X}(t) = \frac{1}{n} \sum_{i=1}^n X_i(t)$  for  $t \in \mathcal{I}$
- Sample cov: for  $(s, t) \in \mathcal{I} \times \mathcal{I}$

$$\hat{\Gamma}(s, t) = \frac{1}{n} \sum_{i=1}^n (X_i(s) - \bar{X}(s))(X_i(t) - \bar{X}(t))$$

- Sample variance:  $\hat{v}(t) = \hat{\Gamma}(t, t)$
- Sample standard deviation:  $\hat{s}(t) = \sqrt{\hat{v}(t)}$



# Phoneme: mean and std function



# Multivariate: Eigen-decomposition of covariance matrix

Let  $\lambda_k$  be the eigenvalue of  $\Sigma$  with the corresponding eigenvector  $\mathbf{u}_k$ , that is,  $\Sigma \mathbf{u}_k = \lambda_k \mathbf{u}_k$ . Then

$$\begin{aligned}\Sigma &= UDU^\top = [\mathbf{u}_1, \dots, \mathbf{u}_p] \text{diag}(\lambda_1, \dots, \lambda_p) \begin{pmatrix} \mathbf{u}_1^\top \\ \vdots \\ \mathbf{u}_p^\top \end{pmatrix} \\ &= \sum_{k=1}^p \lambda_k \mathbf{u}_k \mathbf{u}_k^\top\end{aligned}$$

where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$  and  $U^\top U = UU^\top = I_p$ .

For sample covariance matrix:

$$\Sigma_n = U_n D_n U_n^\top = \sum_{k=1}^p \hat{\lambda}_k \hat{\mathbf{u}}_k \hat{\mathbf{u}}_k^\top$$

# Variance decomposition

Total variation = Mean variation + Mean residual variation:

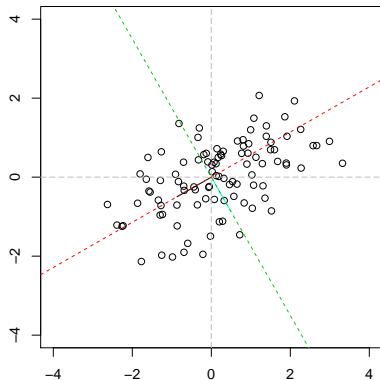
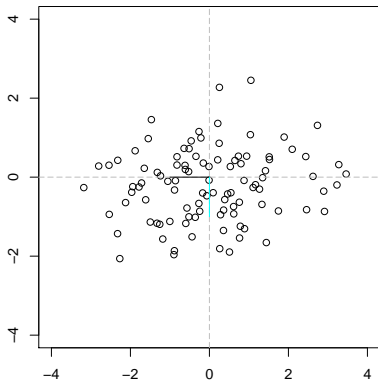
$$\sum_{i=1}^n \|X_i\|^2 = \sum_{i=1}^n \|\bar{X}\|^2 + \sum_{i=1}^n \|X_i - \bar{X}\|^2$$

Decomposition of mean residual variation:

$$\begin{aligned}\sum_{i=1}^n \|X_i - \bar{X}\|^2 &= \sum_{i=1}^n (X_i - \bar{X})^\top (X_i - \bar{X}) \\ &= \text{trace}\left(\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^\top\right) \\ \frac{1}{n} \sum_{i=1}^n \|X_i - \bar{X}\|^2 &= \text{trace}(\Sigma_n) \\ &= \text{trace}(U_n D_n U_n^\top) = \text{trace}(D_n) \\ &= \sum_{j=1}^p \hat{\lambda}_j\end{aligned}$$

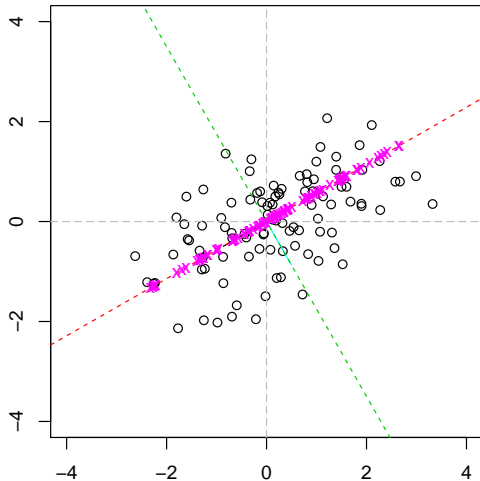
# PCA - toy example

$$p = 2$$



# PCA

Find *direction of greatest variability*



# PCA data representation

- ▶ Eigenvectors are orthonormal basis:

$$X_i = \boldsymbol{\mu} + \sum_{k=1}^p \langle X_i - \boldsymbol{\mu}, \mathbf{u}_k \rangle \mathbf{u}_k = \boldsymbol{\mu} + \sum_{k=1}^p a_{ik} \mathbf{u}_k$$

where  $a_{ik} \sim N(0, \lambda_k)$ .

- ▶ Sample representation:

$$X_i = \bar{X} + \sum_{k=1}^p \langle X_i - \bar{X}, \hat{\mathbf{u}}_k \rangle \hat{\mathbf{u}}_k = \bar{X} + \sum_{k=1}^p \hat{a}_{ik} \hat{\mathbf{u}}_k$$

where  $\hat{a}_{ik} \sim (0, \hat{\lambda}_k)$

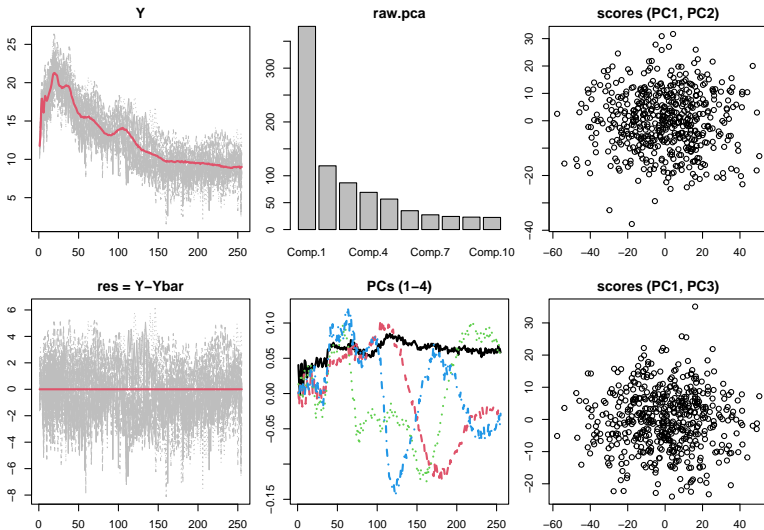
- ▶ Dimension reduction: use  $K < p$  components

$$\hat{X}_i = \bar{X} + \sum_{k=1}^K \langle X_i - \bar{X}, \hat{\mathbf{u}}_k \rangle \hat{\mathbf{u}}_k = \bar{X} + \sum_{k=1}^K \hat{a}_{ik} \hat{\mathbf{u}}_k$$

- ▶ Estimation of the inverse of covariance matrix:

$$\Sigma_n \approx \sum_{k=1}^K \hat{\lambda}_k \hat{\mathbf{u}}_k \hat{\mathbf{u}}_k^\top, \quad \widehat{\Sigma}^{-1} = \sum_{k=1}^K \frac{1}{\hat{\lambda}_k} \hat{\mathbf{u}}_k \hat{\mathbf{u}}_k^\top$$

# Phoneme: (multivariate) PCA



## Karhunen-Lóeve decomposition:

Assume that  $E[X_i(t)] = \mu(t)$ ,  $\text{cov}(X_i(s), X_i(t)) = \Gamma(s, t)$

$$\Gamma(s, t) = \sum_{k=1}^{\infty} \lambda_k \phi_k(s) \phi_k(t)$$

where  $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$  are eigenvalues and  $\phi_1, \phi_2, \dots$  are the corresponding eigenfunctions.

Eigen-decomposition of the variance-covariance function:

$$\int_{\mathcal{I}} \Gamma(s, t) \phi_k(s) ds = \lambda_k \phi_k(t)$$

subject to  $\int_{\mathcal{I}} \phi_k^2(t) = 1$  and  $\int_{\mathcal{I}} \phi_k(t) \phi_{\ell}(t) = 0$  for  $k \neq \ell$ .



# Functional PCA

Functional variable  $X_i$  can be expressed as

$$X_i(t) = \mu(t) + \sum_{k=1}^{\infty} \langle X_i, \phi_k \rangle \phi_k(t) = \mu(t) + \sum_{k=1}^{\infty} \xi_{ik} \phi_k(t)$$

where  $\xi_{ik} \sim (0, \lambda_k)$ .

Approximation based on FPCA:

$$X_i^K(t) = \mu(t) + \sum_{k=1}^K \xi_{ik} \phi_k(t)$$

Dimension reduction based on FPCA:

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

where  $\hat{\xi}_{ik} \sim (0, \hat{\lambda}_k)$ .

# Functional PCA with noisy data

$$Y_{ij} = X_i(t_{ij}) + \varepsilon_i(t_{ij}), \quad X_i(t) = \mu(t) + \sum_{k=1}^K \xi_{ik} \phi_k(t)$$

*Method 1:* discretization on common grid

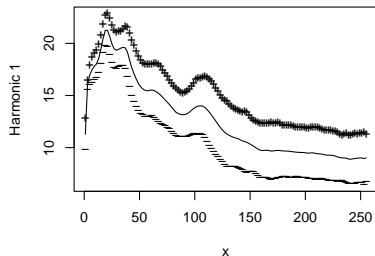
- ▶ Individual smoothing:  $\hat{X}_i(t) = \text{Smooth}(t, Y_{ij}, j = 1, \dots, n_i)$
- ▶ Interpolate at common grid points at  $t_1, \dots, t_m$ :  
 $X_i = (\hat{X}_i(t_1), \dots, \hat{X}_i(t_m))^T$
- ▶ Apply multivariate PCA (up to scaling factors correction)

*Method 2:* FPCA by regularization

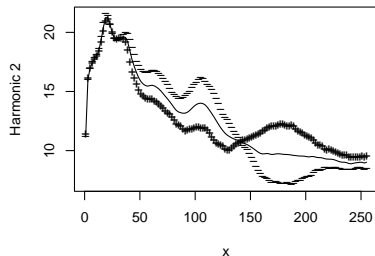
- ▶ Single smoothing for mean:  $\hat{\mu}(t) = \text{Smooth}(\frac{1}{n} \sum_{i=1}^n Y_{ij})$
- ▶ Residuals:  $r_{ij} = Y_{ij} - \hat{\mu}(t_{ij})$ 
  - ▶ Smooth covariance  $(r_{ij} r_{ik})$  before eigen-decomposition
  - ▶ Smooth eigenfunction by regularization

# Phoneme: FPCA eigenfun

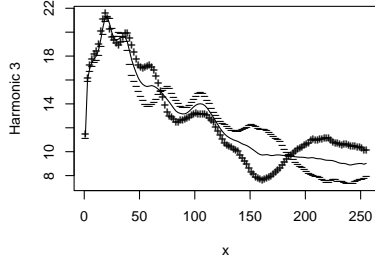
PCA function 1 (Percentage of variability 37.7 )



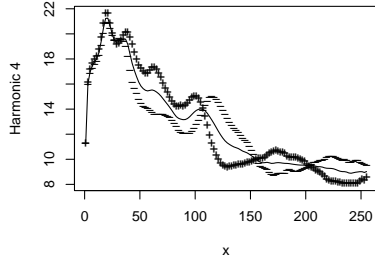
PCA function 2 (Percentage of variability 11.8 )



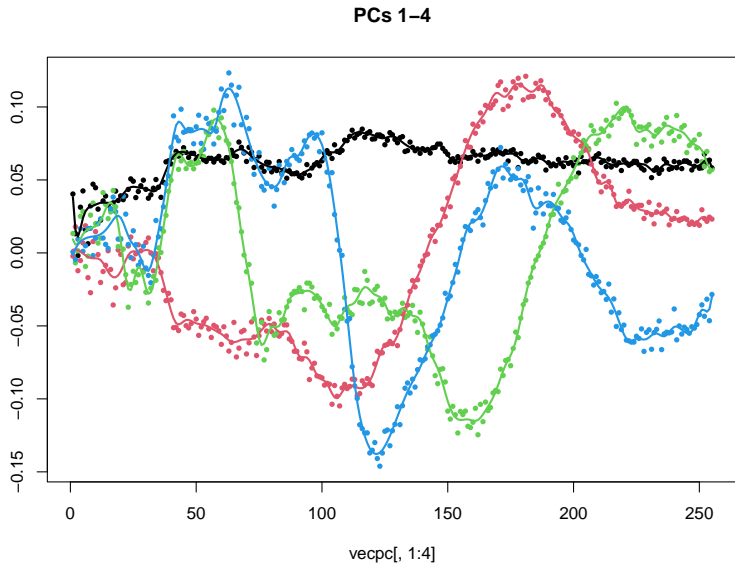
PCA function 3 (Percentage of variability 8.6 )



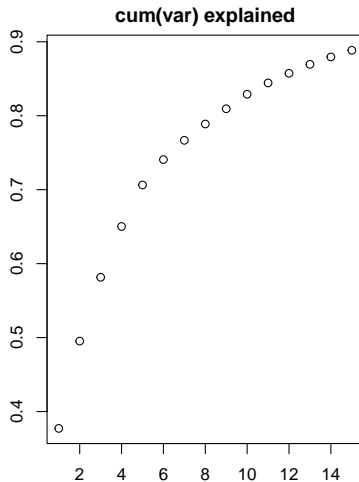
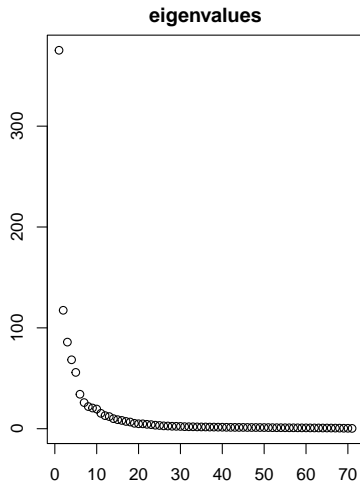
PCA function 4 (Percentage of variability 6.9 )



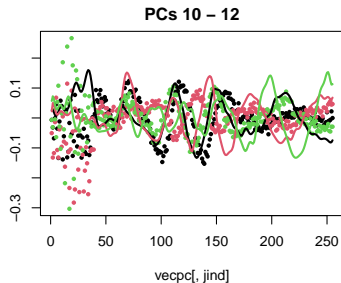
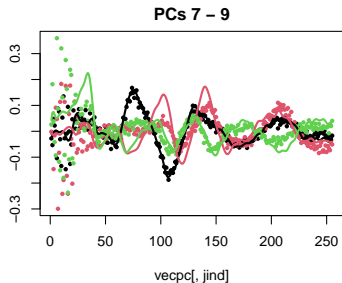
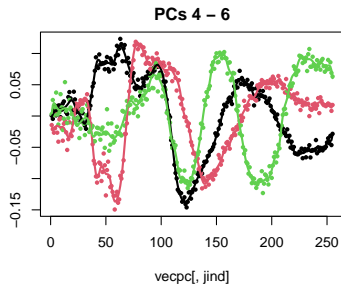
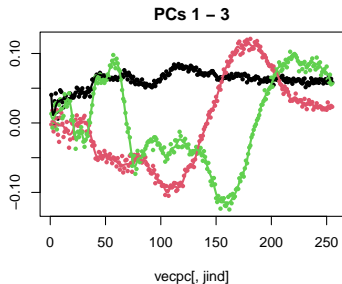
# Phoneme: FPCA vs PCA eigenfunc



# Phoneme: FPCA eigenvalues



# Phoneme: FPCA eigenfunctions (cont)



## Section 4

### Smoothing methods

# Representing functions with noisy data

- ▶ Data:  $n$  noisy observations  $(Y_i, i = 1, \dots, n)$  available at fixed or random design points  $(T_i \in [a, b])$
- ▶ Assumption: data represent unknown smooth function  $(f)$
- ▶ Statistical Model:

$$Y_i = f(T_i) + \varepsilon_i \quad i = 1, \dots, n$$

where the error has mean zero and finite variance.

- ▶ The underlying function represents

$$f(t) = E[Y|T = t] \quad f(T) = E[Y|T]$$

the conditional mean of  $Y$  given  $T = t$ .

- ▶ Aim: recover  $f(t), t \in [a, b]$  from finite number of noisy data
- ▶ Point estimation at  $t$ :  $\hat{X}(t) = \hat{f}(t), t \in [a, b]$

⇒ *Nonparametric regression problem*



# Nonparametric regression

Observed data:  $(t_i, y_i) : i = 1, \dots, n$

Find  $f(\cdot) \in \mathcal{F} = \{f : \text{continuous}\}$  that minimizes the squared error

$$\sum_{i=1}^n \{y_i - f(t_i)\}^2$$

- ▶ feasible set too large: can find an exact solution (*overfit*)
- ▶ need to impose constraints: use *smoothness* constraints
  - ▶ Local (polynomial) approximation: for  $t$  in the neighborhood  $t_0$

$$\begin{aligned} f(t) &\approx f(t_0) + f'(t_0)(t - t_0) + \dots + \frac{f^{(p)}(t_0)}{p!}(t - t_0)^p \\ &= \beta_0 + \beta_1(t - t_0) + \dots + \beta_p(t - t_0)^p \end{aligned}$$

- ▶ Global approximation: choose basis functions  $\{\phi_1, \dots, \phi_k\}$  for all  $t \in [a, b]$

$$f(t) \approx \alpha_1 \phi_1(t) + \alpha_2 \phi_2(t) + \dots + \alpha_k \phi_k(t)$$

# Standard smoothing methods

Different ways to control smoothness in the function:

- ▶ *Kernel smoothing* or *Local polynomial regression*: minimize

$$\sum_{i=1}^n w_i \{y_i - \beta_0 - \beta_1(t_i - t) - \dots - \beta_p(t_i - t)^p\}^2$$

- smoothing parameter: size of neighborhood ( $h$ )

- ▶ *Regression splines*: minimize

$$\sum_{i=1}^n \{y_i - a_1\phi_1(t_i) - \dots - a_k\phi_k(t_i)\}^2$$

- smoothing parameter: number of basis functions  $k$

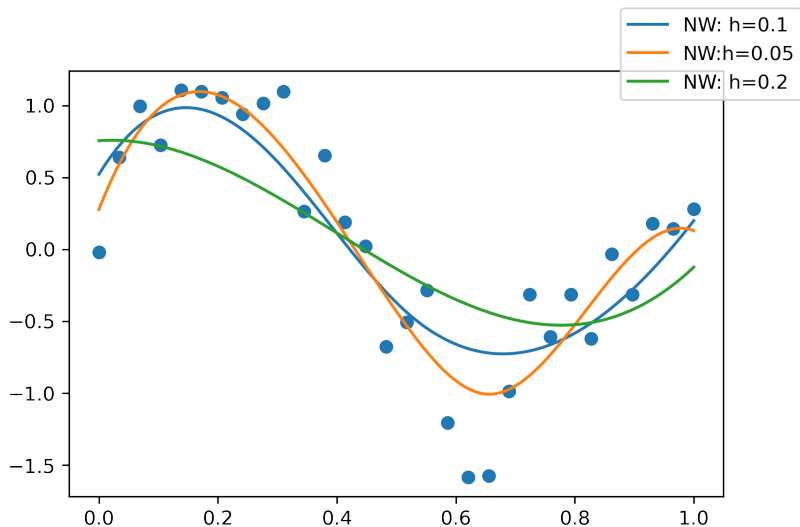
- ▶ *Smoothing splines*: minimize

$$\sum_{i=1}^n \{y_i - a_1\phi_1(t_i) - \dots - a_k\phi_k(t_i)\}^2 + \lambda P(f)$$

where  $P(f)$  is smoothness penalty, often  $P(f) = \int_a^b \{f''(t)\}^2 dt$

- smoothing parameter:  $\lambda$

## Kernel smoothing: smoothing parameter



## Smoothing splines: smoothing parameter

