

# Unsupervised learning from high-dimensional and functional data

FAICEL CHAMROUKHI

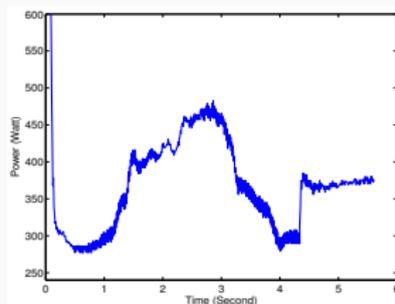
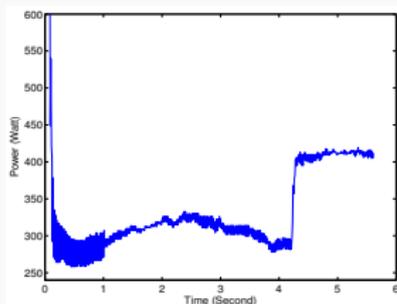


Research Summer School on Statistics for Data Science S4D 2018

June 22, 2018

# Temporal data

## Temporal data with regime changes



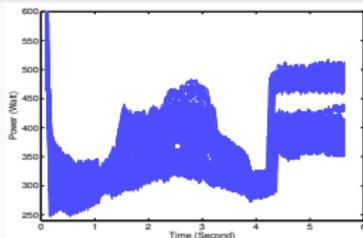
- Data with regime changes over time
- Abrupt and/or smooth regime changes

## Objectives

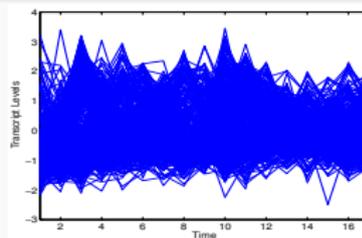
Temporal data modeling and segmentation

# Functional data

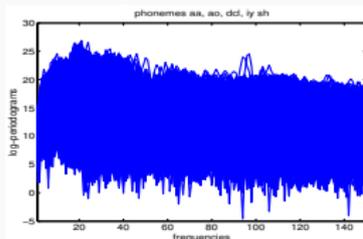
Many curves to analyze



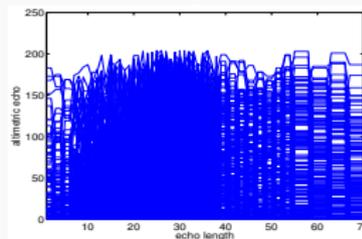
Railway switch curves



Yeast cell cycle curves



Phonemes curves



Satellite waveforms

## Objectives

- Curve clustering/classification (functional data analysis framework)
- Deal with the problem of regime changes  $\leftrightarrow$  Curve segmentation

## Scientific context

- The area of **statistical learning** and **analysis of complex data**.
- **Data** : Complex data  $\leftrightarrow$  *heterogeneous, temporal/dynamical, high-dimensional/functional, incomplete,...*
- **Objective**: Transform the data into knowledge :  
 $\leftrightarrow$  **Reconstruct hidden structure/information, groups/hierarchy of groups, summarizing prototypes, underlying dynamical processes, etc**

## Modeling framework

- **Latent variable** models :  $f(x|\boldsymbol{\theta}) = \int_z f(x, z|\boldsymbol{\theta})dz$

**Generative** formulation :

$$z \sim q(z|\boldsymbol{\theta})$$

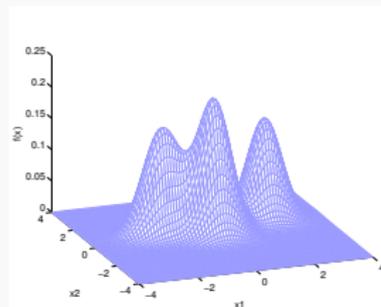
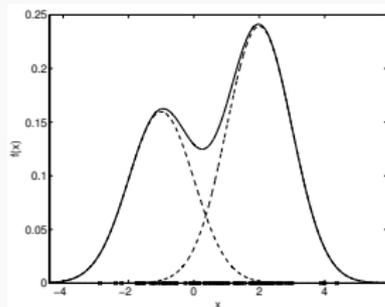
$$x|z \sim f(x|z, \boldsymbol{\theta})$$

- $\leftrightarrow$  **Mixture models** :  $f(x|\boldsymbol{\theta}) = \sum_{k=1}^K \mathbb{P}(z = k) f(x|z = k, \boldsymbol{\theta}_k)$  and extensions

# Mixture modeling framework

## Mixture modeling framework

- Mixture density:  $f(x|\theta) = \sum_{k=1}^K \pi_k f_k(x|\theta_k)$



- Generative model

$$z \sim \mathcal{M}(1; \pi_1, \dots, \pi_K)$$

$$x|z \sim f(x|\theta_z)$$

↪ Algorithms for inferring  $\theta$  from the data

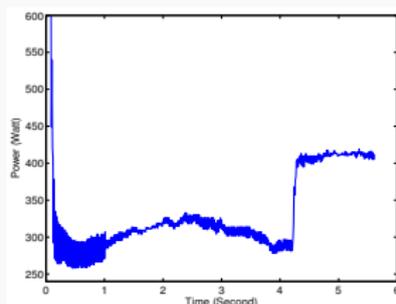
# Outline

- 1 Mixture models for temporal data segmentation
- 2 Mixture models for functional data analysis

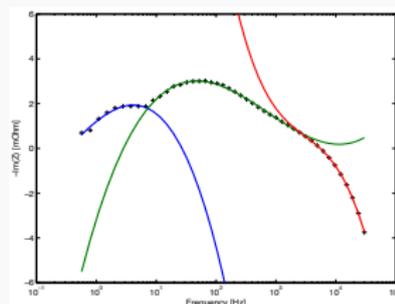
# Outline

- 1 Mixture models for temporal data segmentation
  - Regression with hidden logistic process
- 2 Mixture models for functional data analysis

## Temporal data with regime changes



Railway data



Energy data

# Mixture models for temporal data segmentation

$\mathbf{y} = (y_1, \dots, y_n)$  a time series of  $n$  univariate observations  $y_i \in \mathbb{R}$  observed at the time points  $\mathbf{t} = (t_1, \dots, t_n)$

## Times series segmentation context

- Time series segmentation is a popular problem with a broad literature
- Common problem for different communities, including statistics, detection, signal processing, machine learning, finance
- The observed time series is generated by an underlying process  
↔ segmentation  $\equiv$  recovering the parameters the process' states.
- Conventional solutions are subject to limitations in the control of the transitions between these states
- ↔ Propose generative latent data modeling for segmentation and approximation
- ↔ segmentation  $\equiv$  inferring the model parameters and the underlying

# Regression with hidden logistic process

Let  $\mathbf{y} = (y_1, \dots, y_n)$  be a time series of  $n$  univariate observations  $y_i \in \mathbb{R}$  observed at the time points  $\mathbf{t} = (t_1, \dots, t_n)$  governed by  $K$  regimes.

## The Regression model with Hidden Logistic Process (RHLP) [1]

$$y_i = \beta_{z_i}^T \mathbf{x}_i + \sigma_{z_i} \epsilon_i \quad ; \quad \epsilon_i \sim \mathcal{N}(0, 1), \quad (i = 1, \dots, n)$$
$$Z_i \sim \mathcal{M}(1, \pi_1(t_i; \mathbf{w}), \dots, \pi_K(t_i; \mathbf{w}))$$

Polynomial segments  $\beta_{z_i}^T \mathbf{x}_i$  with  $\mathbf{x}_i = (1, t_i, \dots, t_i^p)^T$  with logistic probabilities

$$\pi_k(t_i; \mathbf{w}) = \mathbb{P}(Z_i = k | t_i; \mathbf{w}) = \frac{\exp(w_{k1}t_i + w_{k0})}{\sum_{\ell=1}^K \exp(w_{\ell 1}t_i + w_{\ell 0})}$$

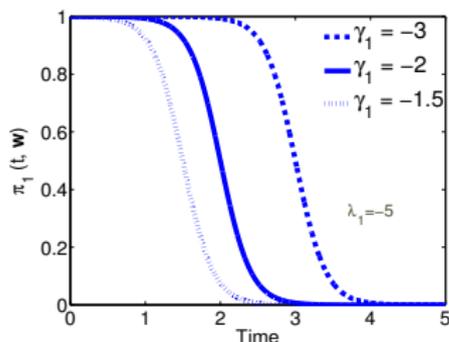
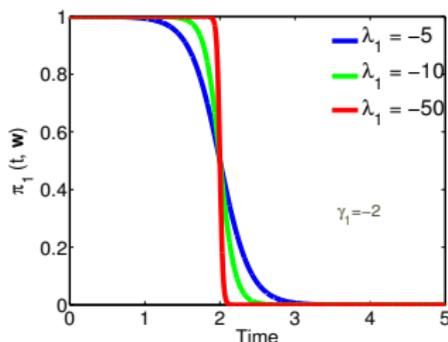
$$f(y_i | t_i; \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k(t_i; \mathbf{w}) \mathcal{N}(y_i; \beta_k^T \mathbf{x}_i, \sigma_k^2)$$

- Both the mixing proportions and the component parameters are time-varying
- Parameter vector of the model :  $\boldsymbol{\theta} = (\mathbf{w}^T, \beta_1^T, \dots, \beta_K^T, \sigma_1^2, \dots, \sigma_K^2)^T$

# Illustration

- Modeling with the logistic distribution allows activating simultaneously and preferentially several regimes during time

$$\pi_k(t_i; \mathbf{w}) = \frac{\exp(\lambda_k(t_i + \gamma_k))}{\sum_{\ell=1}^K \exp(\lambda_\ell(t_i + \gamma_\ell))}$$

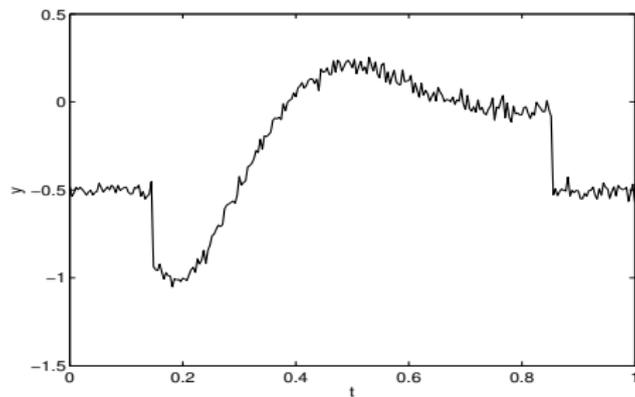


⇒ The parameter  $w_{k1}$  controls the quality of transitions between regimes

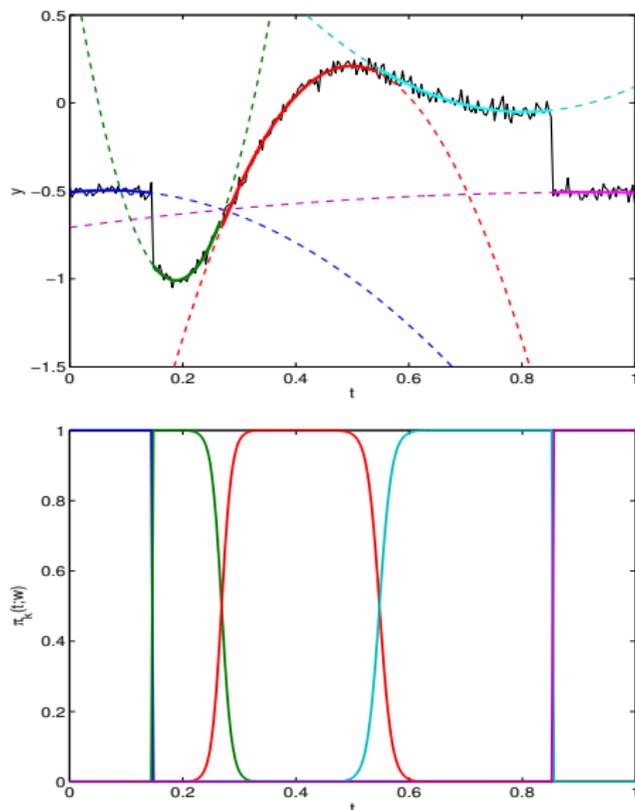
⇒ The parameter  $w_{k0}$  is related to the transition time point

- Ensure time series segmentation into contiguous segments

# Illustration



# Illustration



$K = 5$  polynomial components of degree  $p = 2$

# Parameter estimation: MLE via EM: EM-RHLP

- Parameter vector:  $\boldsymbol{\theta} = (\mathbf{w}^T, \boldsymbol{\beta}_1^T, \dots, \boldsymbol{\beta}_K^T, \sigma_1^2, \dots, \sigma_K^2)^T$
- Maximize the observed-data log-likelihood:

$$\log L(\boldsymbol{\theta}; \mathbf{y}, \mathbf{t}) = \sum_{i=1}^n \log \sum_{k=1}^K \pi_k(t_i; \mathbf{w}) \mathcal{N}(y_i; \boldsymbol{\beta}_k^T \mathbf{x}_i, \sigma_k^2)$$

- Complete-data log-likelihood

$$\log L_c(\boldsymbol{\theta}; \mathbf{y}, \mathbf{t}, \mathbf{z}) = \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log[\pi_k(t_i; \mathbf{w}) \mathcal{N}(y_i; \boldsymbol{\beta}_k^T \mathbf{x}_i, \sigma_k^2)]$$

$Z_{ik} = 1$  if  $Z_i = k$  (i.e., when  $y_i$  belongs to the  $k$ th component)

- The  $Q$ -function

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)}) &= \mathbb{E} \left[ \log L_c(\boldsymbol{\theta}; \mathbf{y}, \mathbf{t}, \mathbf{z}) \mid \mathbf{y}, \mathbf{t}; \boldsymbol{\theta}^{(q)} \right] \\ &= \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}^{(q)} \left[ \log \pi_k(t_i; \mathbf{w}) \mathcal{N}(y_i; \boldsymbol{\beta}_k^T \mathbf{x}_i, \sigma_k^2) \right] \end{aligned}$$

- **E-Step:** compute the posterior component memberships:

$$\tau_{ik}^{(q)} = \mathbb{P}(Z_i = k | y_i, t_i; \boldsymbol{\theta}^{(q)}) = \frac{\pi_k(t_i; \mathbf{w}^{(q)}) \mathcal{N}(y_i; \boldsymbol{\beta}_k^{T(q)} \mathbf{x}_i, \sigma_k^{2(q)})}{\sum_{\ell=1}^K \pi_\ell(t_i; \mathbf{w}^{(q)}) \mathcal{N}(y_i; \boldsymbol{\beta}_\ell^{T(q)} \mathbf{x}_i, \sigma_\ell^{2(q)})}.$$

- **M-Step:** compute the parameter update  $\boldsymbol{\theta}^{(q+1)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(q)})$

$$\boldsymbol{\beta}_k^{(q+1)} = \left[ \sum_{i=1}^n \tau_{ik}^{(q)} \mathbf{x}_i \mathbf{x}_i^T \right]^{-1} \sum_{i=1}^n \tau_{ik}^{(q)} y_i \mathbf{x}_i \quad \text{weighted polynomial regression}$$

$$\sigma_k^{2(q+1)} = \frac{1}{\sum_{i=1}^n \tau_{ik}^{(q)}} \sum_{i=1}^n \tau_{ik}^{(q)} (y_i - \boldsymbol{\beta}_k^{T(q+1)} \mathbf{x}_i)^2$$

$$\mathbf{w}^{(q+1)} = \arg \max_{\mathbf{w}} \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}^{(q)} \log \pi_k(t_i; \mathbf{w}) \quad \text{weighted logistic regression}$$

# EM-RHLP algorithm

## M-Step: Weighted multi-class logistic regression

$$\mathbf{w}^{(q+1)} = \arg \max_{\mathbf{w}} \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}^{(q)} \log \pi_k(t_i; \mathbf{w})$$

- A convex optimization problem
- Solved with a multi-class Iteratively Reweighted Least Squares (IRLS) algorithm (Newton-Raphson)

$$\mathbf{w}^{(l+1)} = \mathbf{w}^{(l)} - \left[ \frac{\partial^2 Q_{\mathbf{w}}(\mathbf{w}, \boldsymbol{\theta}^{(q)})}{\partial \mathbf{w} \partial \mathbf{w}^T} \right]_{\mathbf{w}=\mathbf{w}^{(l)}}^{-1} \left. \frac{\partial Q_{\mathbf{w}}(\mathbf{w}, \boldsymbol{\theta}^{(q)})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^{(l)}}$$

- Analytic calculation of the Hessian and the gradient
- EM-RHLP algorithm complexity:  $\mathcal{O}(I_{\text{EM}} I_{\text{IRLS}} K^3 p^3 n)$  (more advantageous than dynamic programming).

# Time series approximation and segmentation

## 1 Approximation: a prototype mean curve

$$\hat{y}_i = \mathbb{E}[y_i | t_i; \hat{\boldsymbol{\theta}}] = \sum_{k=1}^K \pi_k(t_i; \hat{\mathbf{w}}) \hat{\boldsymbol{\beta}}_k^T \mathbf{x}_i$$

↪ A smooth and flexible approximation thanks to the the logistic weights

↪ The RHLP can be used as nonlinear regression model  $y_i = f(t_i; \boldsymbol{\theta}) + \epsilon_i$  by covering functions of the form  $f(t_i; \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k(t_i; \mathbf{w}) \boldsymbol{\beta}_k^T \mathbf{x}_i$  [3]

## 2 Curve segmentation:

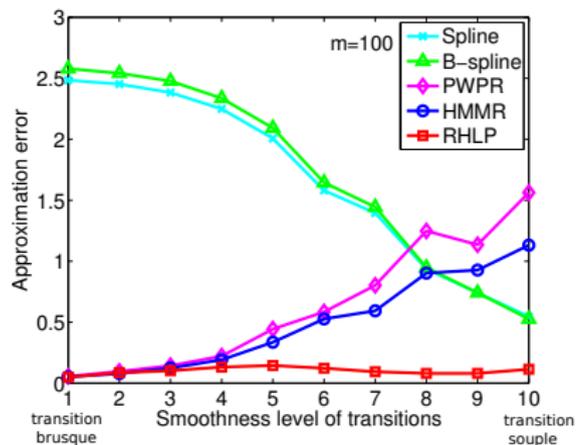
$$\hat{z}_i = \arg \max_{1 \leq k \leq K} \mathbb{E}[z_i | t_i; \hat{\mathbf{w}}] = \arg \max_{1 \leq k \leq K} \pi_k(t_i; \hat{\mathbf{w}})$$

## 3 Model selection Application of BIC, ICL

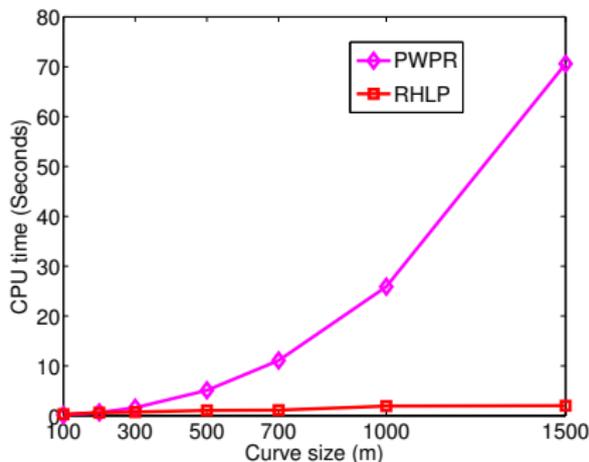
$\text{BIC}(K, p) = \log L(\hat{\boldsymbol{\theta}}) - \frac{\nu_{\boldsymbol{\theta}} \log(n)}{2}$ ;  $\text{ICL}(K, p) = \log L_c(\hat{\boldsymbol{\theta}}) - \frac{\nu_{\boldsymbol{\theta}} \log(n)}{2}$  where  $\nu_{\boldsymbol{\theta}} = K(p + 4) - 2$ .

# Evaluation in modeling and segmentation

Approximation error as a function of the speed of transitions

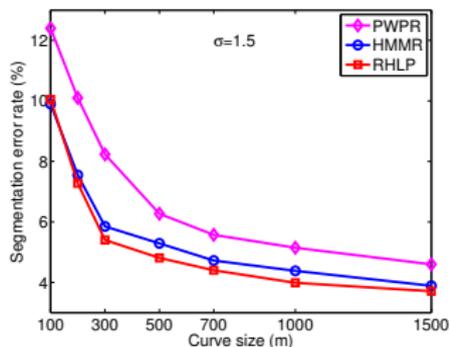
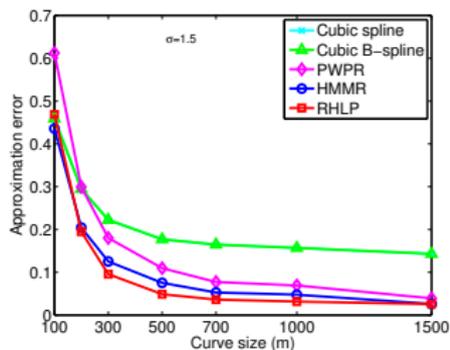


Computing time

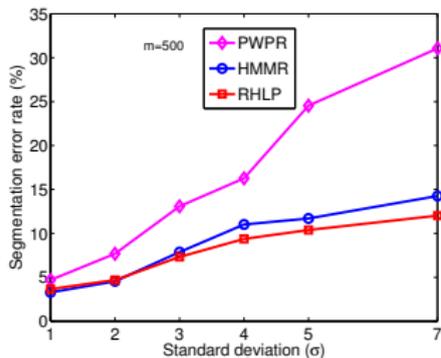
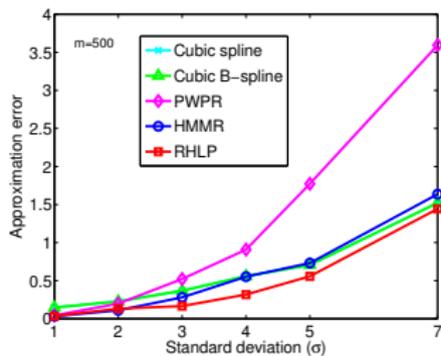


# Evaluation in approximation and segmentation

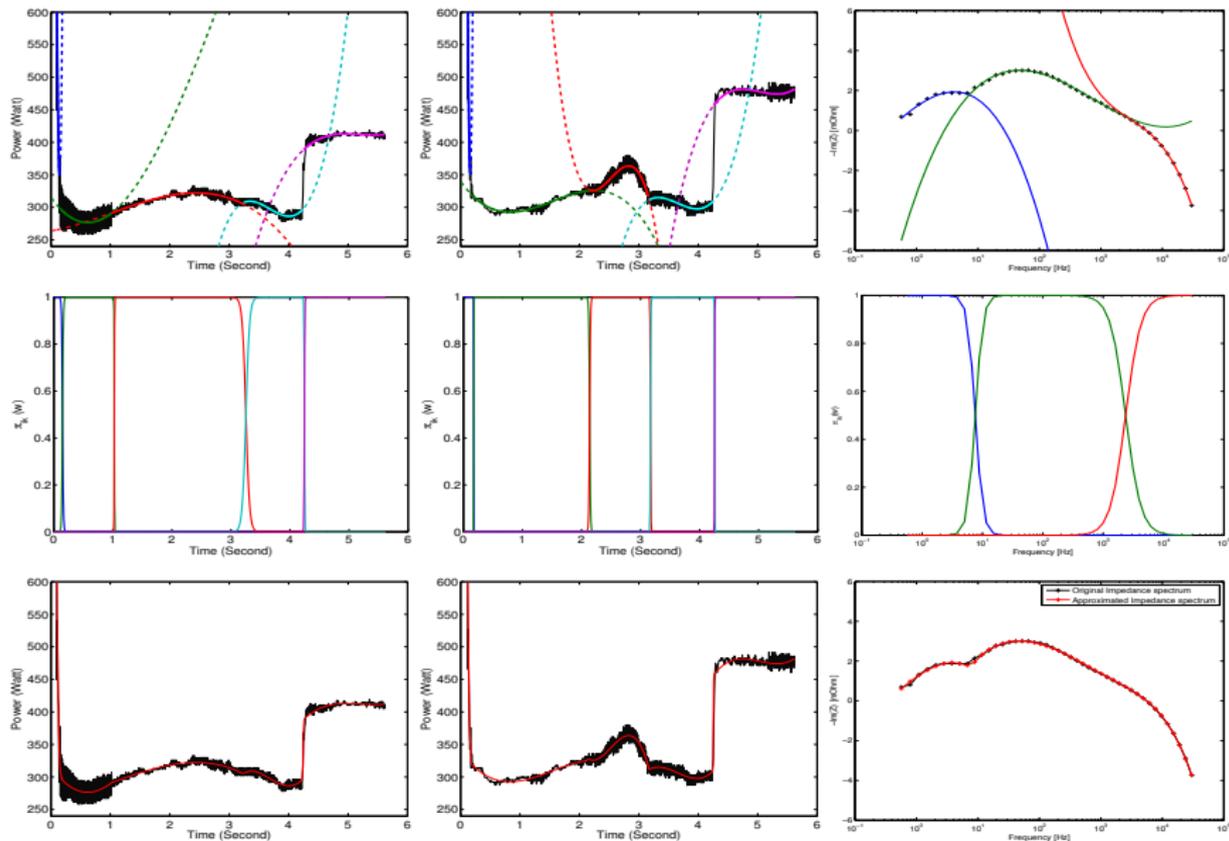
## varying $m$



## varying $\sigma$



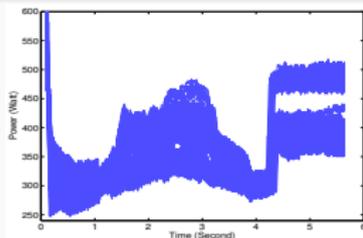
# Application to real data



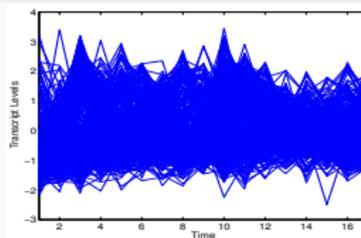
- 1 Mixture models for temporal data segmentation
- 2 Mixture models for functional data analysis
  - Mixture of piecewise regressions
  - Mixture of hidden logistic process regressions
  - Functional discriminant analysis

# Functional data analysis context

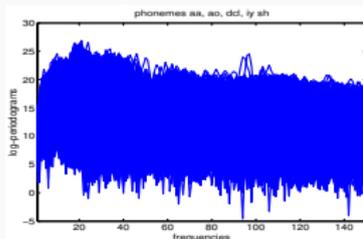
Many curves to analyze



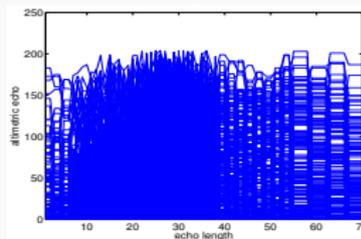
Railway switch curves



Yeast cell cycle curves



Phonemes curves



Satellite waveforms

## Objectives

- Curve clustering/classification (functional data analysis framework)
- Deal with the problem of regime changes  $\leftrightarrow$  Curve segmentation

# Functional data analysis context

## Data

- The individuals are entire functions (e.g., curves, surfaces)
- A set of  $n$  univariate curves  $((\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n))$
- $(\mathbf{x}_i, \mathbf{y}_i)$  consists of  $m_i$  observations  $\mathbf{y}_i = (y_{i1}, \dots, y_{im_i})$  observed at the independent covariates, (e.g., time  $t$  in time series),  $(x_{i1}, \dots, x_{im_i})$

## Objectives: exploratory or decisional

- 1 Unsupervised classification (clustering, segmentation) of functional data, particularly curves with regime changes: [4] [9], [C11] [16]
- 2 Discriminant analysis of functional data: [2], [5]

## Functional data clustering/classification tools

- A broad literature (Kmeans-type, Model-based, etc)  
⇒ Mixture-model based cluster and discriminant analyzes

# Mixture modeling framework for functional data

- The functional mixture model:

$$f(\mathbf{y}|\mathbf{x}; \Psi) = \sum_{k=1}^K \alpha_k f_k(\mathbf{y}|\mathbf{x}; \Psi_k)$$

- $f_k(\mathbf{y}|\mathbf{x})$  are tailored to functional data: can be polynomial (B-)spline regression, regression using wavelet bases etc, or Gaussian process regression, functional PCA

↔ more tailored to approximate smooth functions

↔ do not account for segmentation

Here  $f_k(\mathbf{y}|\mathbf{x})$  itself exhibits a clustering property via hidden variables (regimes):

- 1 Riecewise regression model (PWR)
- 2 Regression model with a hidden process (RHLP)

# Piecewise regression mixture model (PWRM) [9]

- A probabilistic version of the  $K$ -means-like approach of (?)

$$f(\mathbf{y}_i | \mathbf{x}_i; \Psi) = \sum_{k=1}^K \alpha_k \underbrace{\prod_{r=1}^{R_k} \prod_{j \in I_{kr}} \mathcal{N}(y_{ij}; \beta_{kr}^T \mathbf{x}_{ij}, \sigma_{kr}^2)}_{\text{PWR}}$$

$I_{kr} = (\xi_{kr}, \xi_{k,r+1}]$  are the element indexes of segment  $r$  for component  $k$

- $\hookrightarrow$  Simultaneously accounts for curve clustering and segmentation
- Parameter vector  $\Psi = (\alpha_1, \dots, \alpha_{K-1}, \boldsymbol{\theta}_1^T, \dots, \boldsymbol{\theta}_K^T, \boldsymbol{\xi}_1^T, \dots, \boldsymbol{\xi}_K^T)^T$  with  $\boldsymbol{\theta}_k = (\beta_{k1}^T, \dots, \beta_{kR_k}^T, \sigma_{k1}^2, \dots, \sigma_{kR_k}^2)^T$  and  $\boldsymbol{\xi}_k = (\xi_{k1}, \dots, \xi_{k,R_k+1})^T$

## Parameter estimation

- 1 Maximum likelihood estimation: EM-PWRM
- 2 Maximum classification likelihood estimation: CEM-PWRM

# Maximum likelihood estimation via EM: EM-PWRM

- Maximize the observed-data log-likelihood:

$$\log L(\Psi) = \sum_{i=1}^n \log \sum_{k=1}^K \alpha_k \prod_{r=1}^{R_k} \prod_{j \in I_{kr}} \mathcal{N}(y_{ij}; \beta_{kr}^T \mathbf{x}_{ij}, \sigma_{kr}^2)$$

- The complete-data log-likelihood

$$\log L_c(\Psi, \mathbf{z}) = \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log \alpha_k + \sum_{i=1}^n \sum_{k=1}^K \sum_{r=1}^{R_k} \sum_{j \in I_{kr}} Z_{ik} \log \mathcal{N}(y_{ij}; \beta_{kr}^T \mathbf{x}_{ij}, \sigma_{kr}^2)$$

- The conditional expected complete-data log-likelihood

$$Q(\Psi, \Psi^{(q)}) = \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}^{(q)} \log \alpha_k + \sum_{i=1}^n \sum_{k=1}^K \sum_{r=1}^{R_k} \sum_{j \in I_{kr}} \tau_{ik}^{(q)} \log \mathcal{N}(y_{ij}; \beta_{kr}^T \mathbf{x}_{ij}, \sigma_{kr}^2)$$

# EM-PWRM algorithm

## E-step: Compute the $Q$ -function

$\hookrightarrow$  Compute the posterior probability that the  $i$ th curve belongs to component  $k$ :

$$\tau_{ik}^{(q)} = \mathbb{P}(Z_i = k | \mathbf{y}_i, \mathbf{x}_i; \Psi^{(q)}) = \frac{\alpha_k^{(q)} f_k(\mathbf{y}_i | \mathbf{x}_i; \Psi_k^{(q)})}{\sum_{k'=1}^K \alpha_{k'}^{(q)} f_{k'}(\mathbf{y}_i | \mathbf{x}_i; \Psi_{k'}^{(q)})}$$

## M-step: Compute the update $\Psi^{(q+1)} = \arg \max_{\Psi} Q(\Psi, \Psi^{(q)})$

- $\alpha_k^{(q+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(q)}}{n}$ , ( $k = 1, \dots, K$ )
- maximization w.r.t the piecewise regression parameters  $\{\xi_{kr}, \beta_{kr}, \sigma_{kr}^2\} \hookrightarrow$  a weighted piecewise regression problem  $\hookrightarrow$  dynamic programming:

$$\beta_{kr}^{(q+1)} = \left[ \sum_{i=1}^n \tau_{ik}^{(q)} \mathbf{X}_{ir}^T \mathbf{X}_{ir} \right]^{-1} \sum_{i=1}^n \tau_{ik}^{(q)} \mathbf{X}_{ir} \mathbf{y}_{ir}$$
$$\sigma_{kr}^{2(q+1)} = \frac{1}{\sum_{i=1}^n \sum_{j \in I_{kr}^{(q)}} \tau_{ik}^{(q)}} \sum_{i=1}^n \tau_{ik}^{(q)} \|\mathbf{y}_{ir} - \mathbf{X}_{ir} \beta_{kr}^{(q+1)}\|^2$$

$\mathbf{y}_{ir}$  are the observations of segment  $r$  of the  $i$ th curve and  $\mathbf{X}_{ir}$  its design matrix

## Maximum classification likelihood estimation: CEM-PWRM

- Maximize the complete-data log-likelihood w.r.t  $(\Psi, \mathbf{z})$  simultaneously
- C-step: Bayes' optimal allocation rule:  $\hat{z}_i = \arg \max_{1 \leq k \leq K} \tau_{ik}(\hat{\Psi})$

CEM-PWRM is equivalent to the  $K$ -means-like algorithm of ?:

$$\log L_c(\mathbf{z}, \Psi) \propto \mathcal{J}(\mathbf{z}, \{\mu_{kr}, I_{kr}\}) = \sum_{k=1}^K \sum_{r=1}^{R_k} \sum_{i|Z_i=k} \sum_{j \in I_{kr}} (y_{ij} - \mu_{kr})^2$$

if the following conditions hold:

- $\alpha_k = \frac{1}{K} \forall K$  (identical mixing proportions);
  - $\sigma_{kr}^2 = \sigma^2 \forall r$  and  $\forall k$ ; (isotropic and homoskedastic model);
  - $\mu_{kr}$ : piecewise *constant* regime approximation
- 
- Curve clustering:  $\hat{z}_i = \arg \max_k \tau_{ik}(\hat{\Psi})$  with  $\tau_{ik}(\hat{\Psi}) = \mathbb{P}(Z_i | \mathbf{x}_i, \mathbf{y}_i; \hat{\Psi})$
  - Model selection: Application of BIC, ICL
  - Complexity in  $\mathcal{O}(I_{EM} K R n m^2 p^3)$ : Significant computational load for large  $m$

# Simulation results

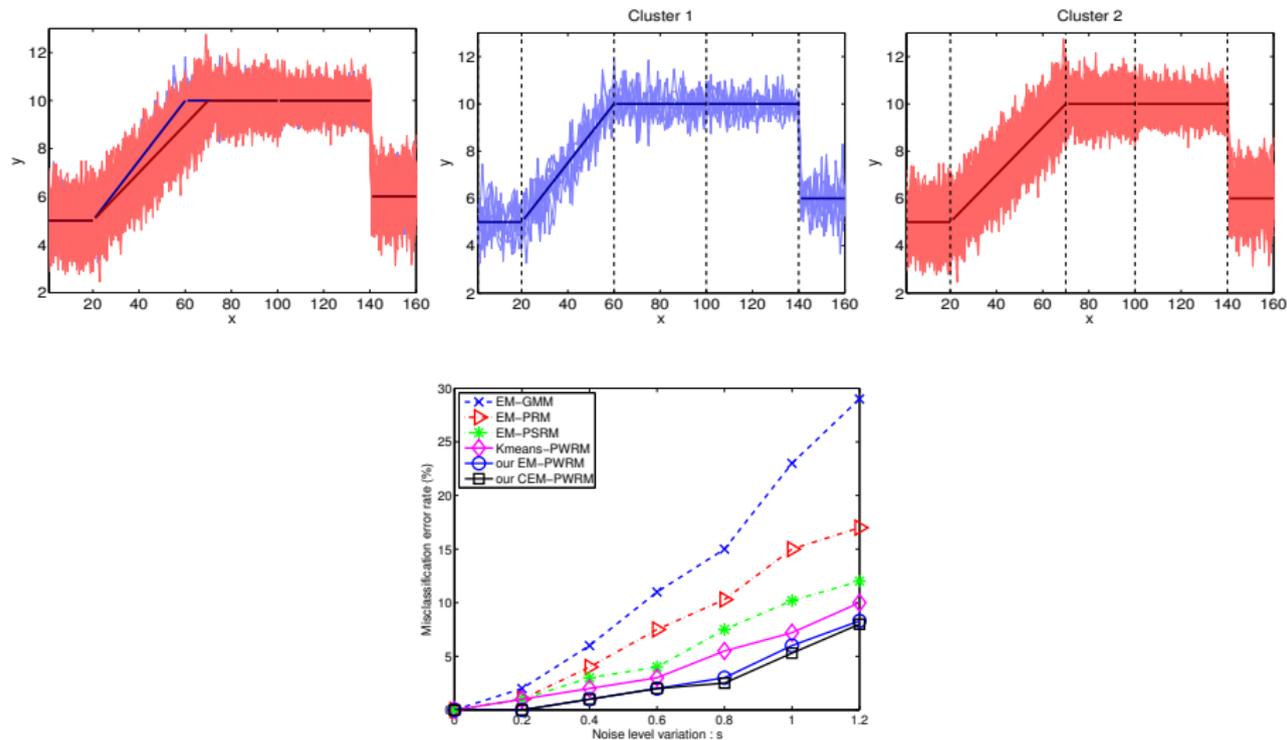
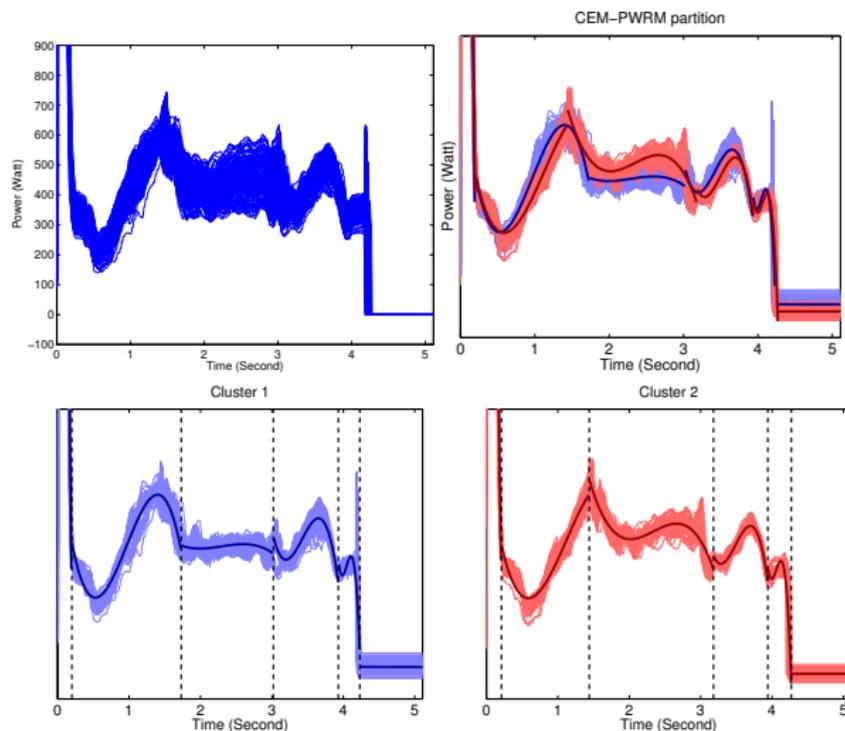


Figure: Misclassification error rate versus the noise level variation.

# Application to switch operation curves

Data set:  $n = 146$  real curves of  $m = 511$  observations.

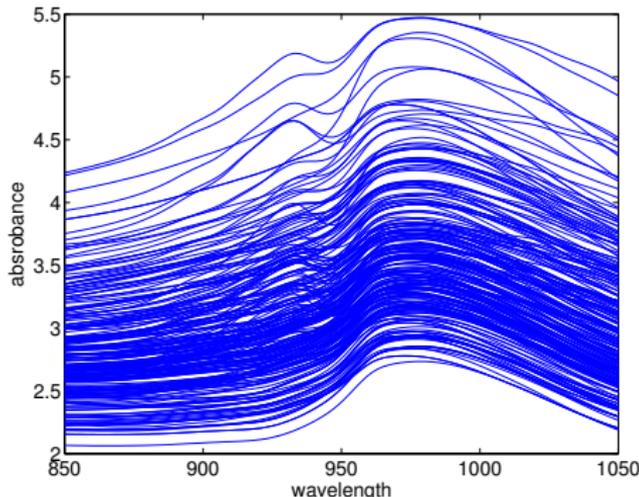
Each curve is composed of  $R = 6$  electromechanical phases (regimes)



# Application to Tecator data

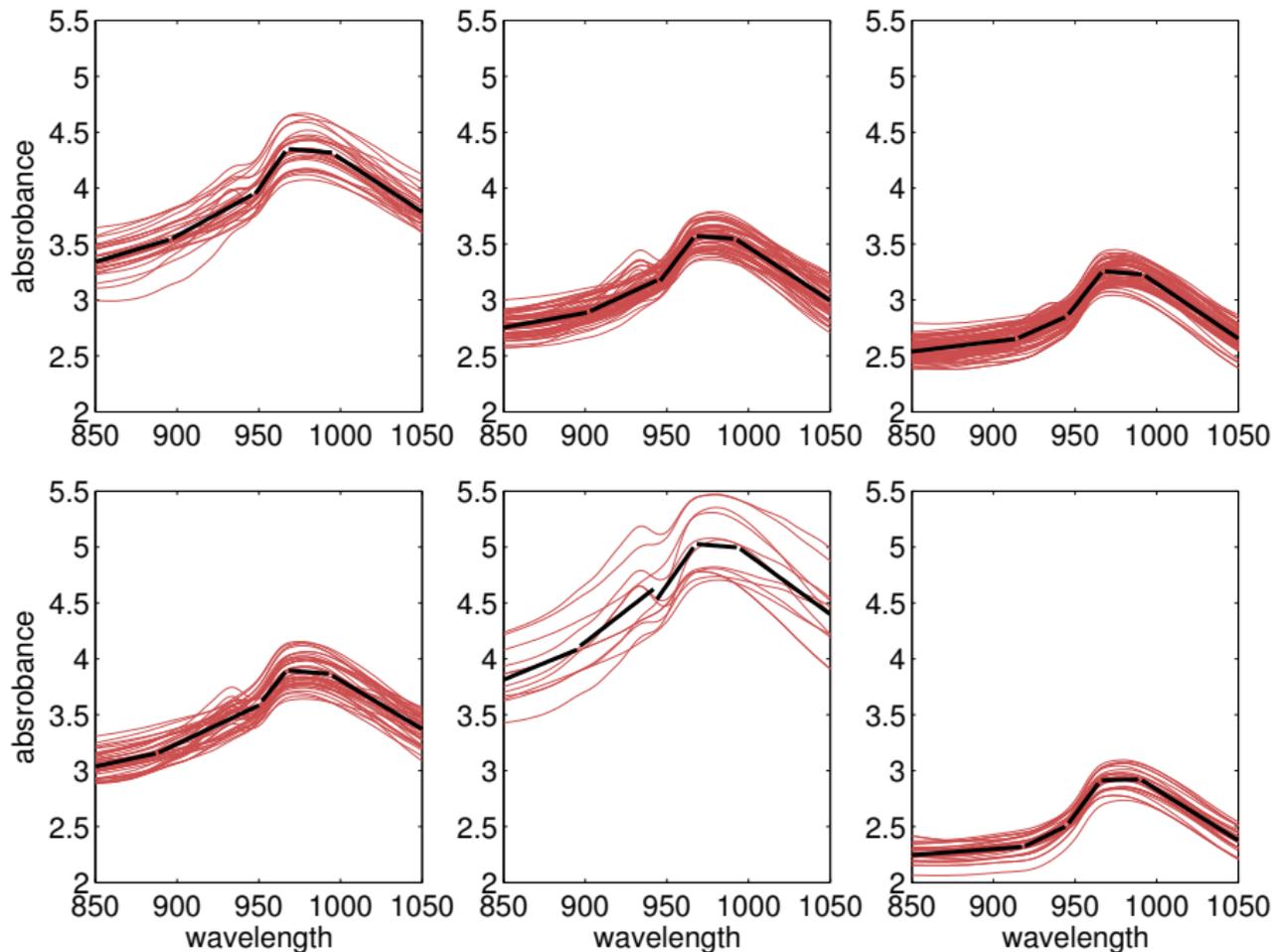
The Tecator data set<sup>1</sup> contains  $n = 240$  spectra with  $m = 100$  observations for each spectrum

Data considered in the same setting as in ? (six clusters, each cluster is approximated by five linear segments ( $R = 5, p = 1$ ))



---

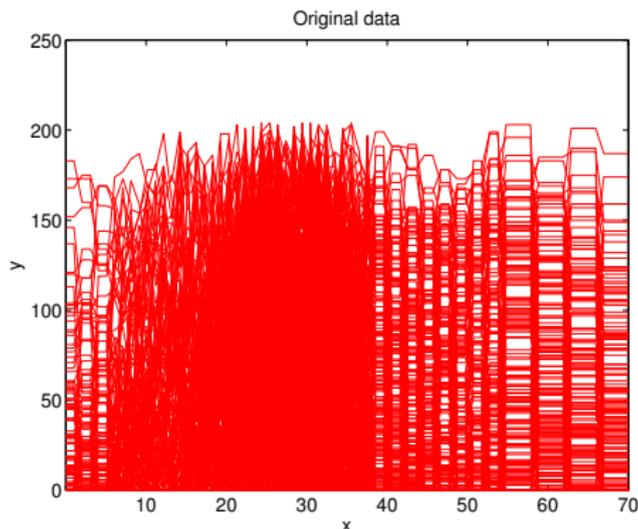
<sup>1</sup>Tecator data are available at <http://lib.stat.cmu.edu/datasets/tecator>.



# Topex/Poseidon satellite data

The Topex/Poseidon radar satellite data<sup>2</sup> contains  $n = 472$  waveforms of the measured echoes, sampled at  $m = 70$  (number of echoes)

We considered the same number of clusters (twenty) and a piecewise linear approximation of four segments per cluster as in ?.

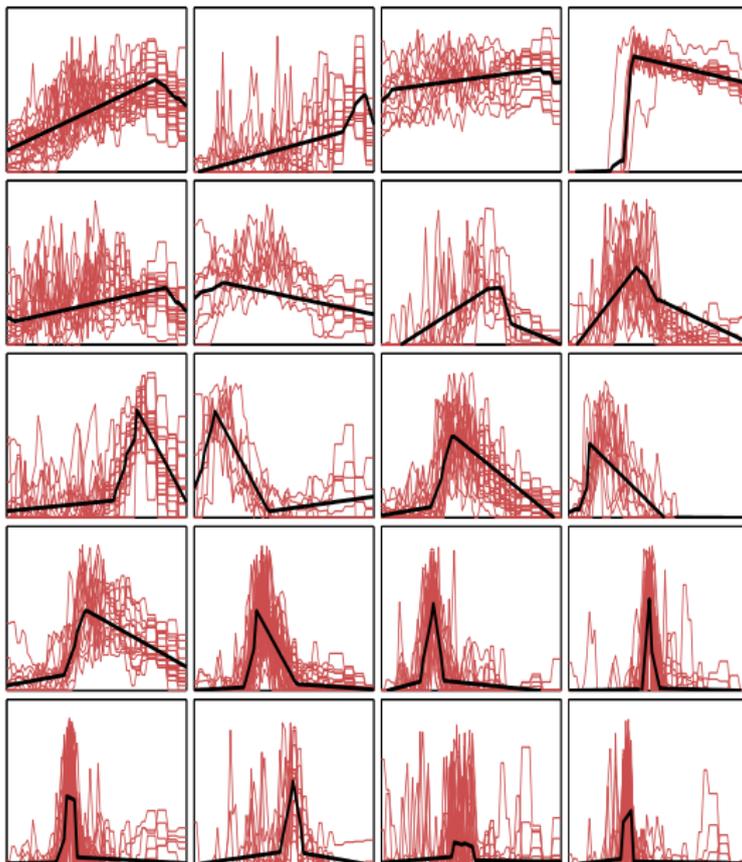


---

<sup>2</sup>Satellite data are available at

<http://www.lsp.ups-tlse.fr/staph/npfda/npfda-datasets.html>.

# CEM-PWRM clustering



# Summary

- Probabilistic approach to the simultaneous curve clustering and optimal segmentation
  - Two algorithms: EM-PWRM and CEM-PWRM
  - CEM-PWRM is a probabilistic-based version of the  $K$ -means-like algorithm ?
- 
- If the aim is density estimation, the EM version is suggested (CEM provides biased estimators but is well-tailored to the segmentation/clustering end)
  - For continuous functions the PWRM in its current formulation, may lead to discontinuities between segments for the piecewise approximation.
  - This may be avoided by posterior interpolation as in ?.
  - May lead to significant computational load especially for large time series. However, for quite reasonable dimensions, the algorithms remain usable

# Mixture of hidden logistic process regressions [4]

- The mixture of regressions with hidden logistic processes (MixRHLP):

$$f(\mathbf{y}_i | \mathbf{x}_i; \Psi) = \sum_{k=1}^K \alpha_k \underbrace{\prod_{j=1}^{m_i} \sum_{r=1}^{R_k} \pi_{kr}(x_j; \mathbf{w}_k) \mathcal{N}(y_{ij}; \beta_{kr}^T \mathbf{x}_j, \sigma_{kr}^2)}_{\text{RHLP}}$$

$$\pi_{kr}(x_j; \mathbf{w}_k) = \mathbb{P}(H_{ij} = r | Z_i = k, x_j; \mathbf{w}_k) = \frac{\exp(w_{kr0} + w_{kr1}x_j)}{\sum_{r'=1}^{R_k} \exp(w_{kr'0} + w_{kr'1}x_j)},$$

- Two types of component memberships:

↪ cluster memberships (global)  $Z_{ik} = 1$  iff  $Z_i = k$

↪ regime memberships for a given cluster (local):  $H_{ijr} = 1$  iff  $H_{ij} = r$

MixRHLP deals better with the quality of regime changes

- Parameter estimation via the EM algorithm: EM-MixRHLP

# MLE estimation via the EM algorithm

- The observed-data log-likelihood

$$\log L(\Psi) = \sum_{i=1}^n \log \sum_{k=1}^K \alpha_k \prod_{j=1}^{m_i} \sum_{r=1}^{R_k} \pi_{kr}(x_j; \mathbf{w}_k) \mathcal{N}(y_{ij}; \beta_{kr}^T \mathbf{x}_j, \sigma_{kr}^2)$$

- The complete-data log-likelihood:

$$\log L_c(\Psi) = \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log \alpha_k + \sum_{i,j} \sum_{k=1}^K \sum_{r=1}^{R_k} Z_{ik} H_{ijr} \log \left[ \pi_{kr}(x_j; \mathbf{w}_k) \mathcal{N}(y_{ij}; \beta_{kr}^T \mathbf{x}_j, \sigma_{kr}^2) \right]$$

- The conditional expected complete-data log-likelihood

$$\begin{aligned} Q(\Psi, \Psi^{(q)}) &= \mathbb{E} \left[ \log L_c(\Psi) \mid \mathcal{D}; \Psi^{(q)} \right] \\ &= \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}^{(q)} \log \alpha_k + \sum_{i,j} \sum_{k=1}^K \sum_{r=1}^{R_k} \tau_{ik}^{(q)} \gamma_{ijr}^{(q)} \log \left[ \pi_{kr}(x_j; \mathbf{w}_k) \mathcal{N}(y_{ij}; \beta_{kr}^T \mathbf{x}_j, \sigma_{kr}^2) \right]. \end{aligned}$$

# EM-MixRHLP algorithm

## E-step

- The posterior cluster memberships:

$$\tau_{ik}^{(q)} = \mathbb{P}(Z_i = k | \mathbf{y}_i, \mathbf{x}_i; \Psi_k^{(q)}) = \frac{\alpha_k^{(q)} f(\mathbf{y}_i | Z_i = k, \mathbf{x}_i; \Psi_k^{(q)})}{\sum_{k'=1}^K \alpha_{k'}^{(q)} f(\mathbf{y}_i | Z_i = k', \mathbf{x}_i; \Psi_{k'}^{(q)})}$$

- the posterior regime memberships:

$$\gamma_{ijr}^{(q)} = \mathbb{P}(H_{ij} = r | Z_i = k, y_{ij}, t_j; \Psi_k^{(q)}) = \frac{\pi_{kr}(x_j; \mathbf{w}_k^{(q)}) \mathcal{N}(y_{ij}; \boldsymbol{\beta}_{kr}^T(x_j), \sigma_{kr}^2(x_j))}{\sum_{r'=1}^{R_k} \pi_{kr'}(x_j; \mathbf{w}_k^{(q)}) \mathcal{N}(y_{ij}; \boldsymbol{\beta}_{kr'}^T(x_j), \sigma_{kr'}^2(x_j))}$$

Computed directly (i.e, without a forward-backward recursion as in the Markovian model).

# M-step of the EM-MixRHLP

**M-step:** calculate the update  $\Psi^{(q+1)} = \arg \max_{\Psi} Q(\Psi, \Psi^{(q)})$ .

- Mixing proportions update: standard

$$\alpha_k^{(q+1)} = \frac{1}{n} \sum_{i=1}^n \tau_{ik}^{(q)}, \quad (k = 1, \dots, K).$$

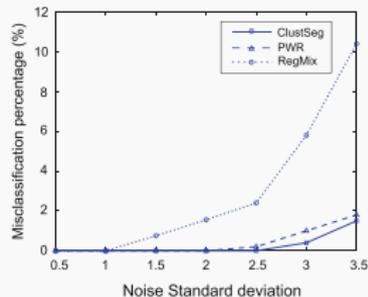
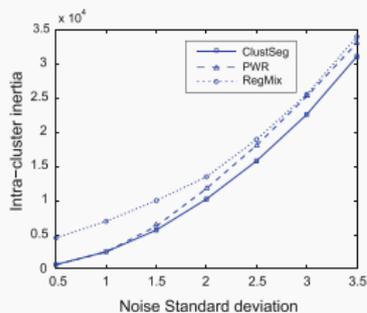
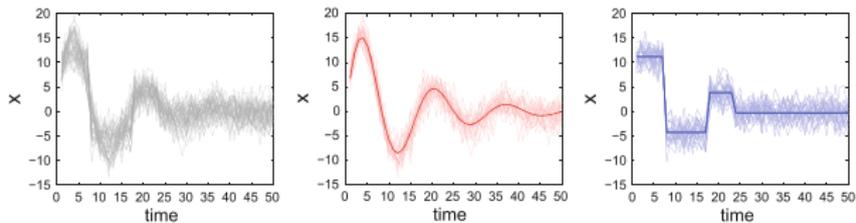
- Regression parameters update: Analytic weighted least-squares problems

$$\beta_{kr}^{(q+1)} = \left[ \sum_{i=1}^n \tau_{ik}^{(q)} \mathbf{X}_i^T \mathbf{W}_{ikr}^{(q)} \mathbf{X}_i \right]^{-1} \sum_{i=1}^n \tau_{ik}^{(q)} \mathbf{X}_i^T \mathbf{W}_{ikr}^{(q)} \mathbf{y}_i,$$
$$\sigma_{kr}^2{}^{(q+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(q)} \|\sqrt{\mathbf{W}_{ikr}^{(q)}} (\mathbf{y}_i - \mathbf{X}_i \beta_{kr}^{(q+1)})\|^2}{\sum_{i=1}^n \tau_{ik}^{(q)} \text{trace}(\mathbf{W}_{ikr}^{(q)})},$$

where  $\mathbf{W}_{ikr}^{(q)} = \text{diag}(\gamma_{ijr}^{(q)}; j = 1, \dots, m_i)$ .

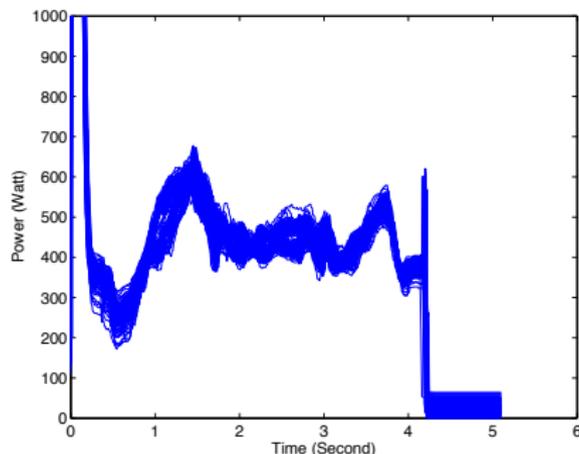
- Maximization w.r.t the logistic processes' parameters  $\{\mathbf{w}_k\}$ : solving multinomial logistic regression problems  $\Rightarrow$  IRLS
- $\hookrightarrow$  EM-MixRHLP has complexity in  $\mathcal{O}(I_{EM} I_{IRLS} K R^3 n m p^3)$  ( $K$ -means like algo. for PWR is in  $\mathcal{O}(I_{KM} K R n m^2 p^3)$ )  $\hookrightarrow$  computationally attractive for large  $m$  with moderate value of  $R$ .

# EM-MixRHLP clustering of simulated data



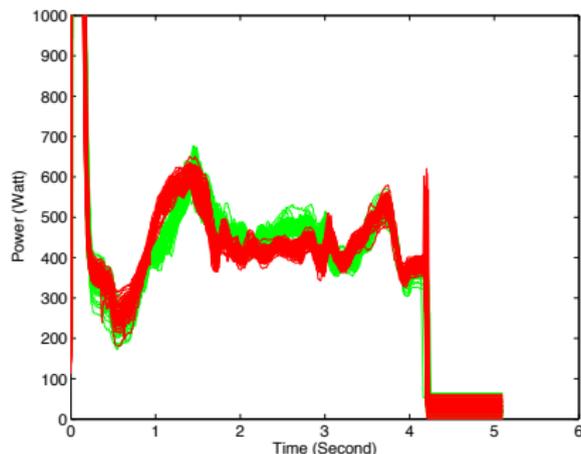
# Clustering switch operations

**Clustering real curves of switch operations** The data set contains 115 curves of  $R = 6$  operations electromechanical process  
 $K = 2$  clusters: operating state without/with possible defect



# Clustering switch operations

**Clustering real curves of switch operations** The data set contains 115 curves of  $R = 6$  operations electromechanical process  
 $K = 2$  clusters: operating state without/with possible defect



# Functional discriminant analysis

## Supervised classification context

- Data: a training set of labeled functions  $((\mathbf{x}_1, \mathbf{y}_1, c_1), \dots, (\mathbf{x}_n, \mathbf{y}_n, c_n))$  where  $c_i \in \{1, \dots, G\}$  is the class label of the  $i$ th curve
- Problem: predict the class label  $c_i$  for a new unlabeled function  $(\mathbf{x}_i, \mathbf{y}_i)$

## Tool: Discriminant analysis

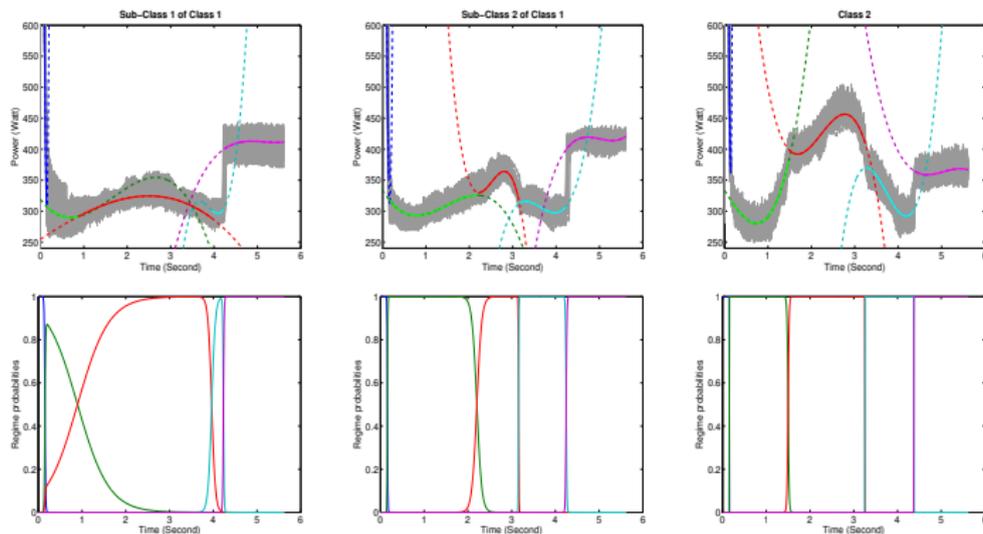
Use the Bayes' allocation rule

$$\hat{c}_i = \arg \max_{1 \leq g \leq G} \frac{\mathbb{P}(C_i = g) f(\mathbf{y}_i | \mathbf{x}_i; \Psi_g)}{\sum_{g'=1}^G \mathbb{P}(C_i = g') f(\mathbf{y}_i | \mathbf{x}_i; \Psi_{g'})},$$

based on a generative model  $f(\mathbf{y}_i | \mathbf{x}_i; \Psi_g)$  for each group  $g$

- Homogeneous classes: Functional Linear Discriminant Analysis [8]
- Dispersed classes: Functional Mixture Discriminant Analysis [5]

# Applications to switch curves



Approach	Classification error rate (%)	Intra-class inertia
FLDA-PR	11.5	$10.7350 \times 10^9$
FLDA-SR	9.53	$9.4503 \times 10^9$
FLDA-RHLP	8.62	$8.7633 \times 10^9$
FMDA-PRM	9.02	$7.9450 \times 10^9$
FMDA-SRM	8.50	$5.8312 \times 10^9$
<b>FMDA-MixRHLP</b>	<b>6.25</b>	<b><math>3.2012 \times 10^9</math></b>

# Summary

- A full generative model for curve clustering and segmentation
- The segmentation is smoothly controlled by logistic functions
- An alternative to the previously described mixture of piecewise regressions
- more advantageous compared to approaches involving dynamic programming namely when using piecewise regression especially for large samples.
- Could be extended to the multivariate case without a major effort

# Some ongoing research and perspectives

- Model-based co-clustering for high-dimensional functional data

## Functional latent block model (FLBM) available soon on arXiv

Data:  $\mathbf{Y} = (\mathbf{y}_{ij})$ :  $n$  individuals defined on a set  $\mathcal{I}$  with  $d$  continuous functional variables defined on a set  $\mathcal{J}$  where  $y_{ij}(t) = \mu(x_{ij}(t); \boldsymbol{\beta}) + \epsilon(t)$ ,  $t$  defined on  $\mathcal{T}$ .

- FLBM model:

$$\begin{aligned} f(\mathbf{Y}|\mathbf{X};\boldsymbol{\Psi}) &= \sum_{(z,w) \in \mathcal{Z} \times \mathcal{W}} \mathbb{P}(\mathbf{Z}, \mathbf{W}) f(\mathbf{Y}|\mathbf{X}, \mathbf{Z}, \mathbf{W}; \boldsymbol{\theta}) \\ &= \sum_{(z,w) \in \mathcal{Z} \times \mathcal{W}} \prod_{i,k} \pi_k^{z_{ik}} \prod_{j,\ell} \rho_\ell^{w_{j\ell}} \prod_{i,j,k,\ell} f(\mathbf{y}_{ij}|\mathbf{x}_{ij}; \boldsymbol{\theta}_{k\ell})^{z_{ik}w_{j\ell}}. \end{aligned}$$

- An RHLP is used as a conditional block distribution  $f(\mathbf{y}_{ij}|\mathbf{x}_{ij}; \boldsymbol{\theta}_{k\ell})$
- Model inference using Stochastic EM

## Some references

- Faicel Chamroukhi and Hien D. Nguyen. Model-based clustering and classification of functional data. 2018. URL <https://chamroukhi.users.lmno.cnrs.fr/papers/MBCC-FDA.pdf>. arXiv:1803.00276v2
- F. Chamroukhi and C. Biernacki. Model-Based Co-Clustering of Multivariate Functional Data. In *ISI 2017 - 61st World Statistics Congress*, Marrakech, Morocco, Jul 2017 . URL <https://hal.archives-ouvertes.fr/hal-01653782>
- F. Chamroukhi. Skew  $t$  mixture of experts. *Neurocomputing - Elsevier*, 266: 390–408, 2017. URL <https://chamroukhi.users.lmno.cnrs.fr/papers/STMoE.pdf>
- F. Chamroukhi. Robust mixture of experts modeling using the  $t$ -distribution. *Neural Networks - Elsevier*, 79:20–36, 2016b. URL <https://chamroukhi.users.lmno.cnrs.fr/papers/TMoE.pdf>
- Faicel Chamroukhi. Piecewise regression mixture for simultaneous functional data clustering and optimal segmentation. *Journal of Classification*, 33(3):374–411, 2016c. URL <https://chamroukhi.users.lmno.cnrs.fr/papers/Chamroukhi-PWRM-JournalClassif-2016.pdf>
- F. Chamroukhi. Unsupervised learning of regression mixture models with unknown number of components. *Journal of Statistical Computation and Simulation*, 86: 2308 – 2334, 2016a. URL <https://chamroukhi.users.lmno.cnrs.fr/papers/Chamroukhi-JSCS-2015.pdf>

- F. Chamroukhi, D. Trabelsi, S. Mohammed, L. Oukhellou, and Y. Amirat. Joint segmentation of multivariate time series with hidden process regression for human activity recognition. *Neurocomputing*, 120:633–644, November 2013b. URL [https://chamroukhi.users.lmno.cnrs.fr/papers/chamroukhi\\_et\\_al\\_neucomp2013b.pdf](https://chamroukhi.users.lmno.cnrs.fr/papers/chamroukhi_et_al_neucomp2013b.pdf)
- F. Chamroukhi, H. Glotin, and A. Samé. Model-based functional mixture discriminant analysis with hidden process regression for curve classification. *Neurocomputing*, 112:153–163, 2013a. URL [https://chamroukhi.users.lmno.cnrs.fr/papers/chamroukhi\\_et\\_al\\_neucomp2013a.pdf](https://chamroukhi.users.lmno.cnrs.fr/papers/chamroukhi_et_al_neucomp2013a.pdf)
- A. Samé, F. Chamroukhi, Gérard Govaert, and P. Aknin. Model-based clustering and segmentation of time series with changes in regime. *Advances in Data Analysis and Classification*, 5:301–321, 2011. URL <https://chamroukhi.users.lmno.cnrs.fr/papers/adac-2011.pdf>
- F. Chamroukhi, A. Samé, G. Govaert, and P. Aknin. A hidden process regression model for functional data description. application to curve discrimination. *Neurocomputing*, 73(7-9):1210–1221, 2010. URL [https://chamroukhi.users.lmno.cnrs.fr/papers/chamroukhi\\_neucomp\\_2010.pdf](https://chamroukhi.users.lmno.cnrs.fr/papers/chamroukhi_neucomp_2010.pdf)
- F. Chamroukhi, A. Samé, G. Govaert, and P. Aknin. Time series modeling by a regression approach based on a latent process. *Neural Networks*, 22(5-6):593–602, 2009. URL [https://chamroukhi.users.lmno.cnrs.fr/papers/Chamroukhi\\_Neural\\_Networks\\_2009.pdf](https://chamroukhi.users.lmno.cnrs.fr/papers/Chamroukhi_Neural_Networks_2009.pdf)

# Unsupervised learning from high-dimensional and functional data

FAICEL CHAMROUKHI



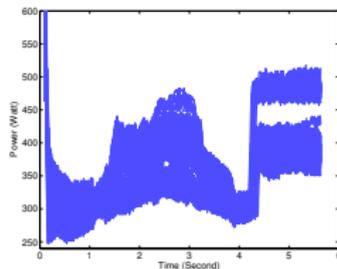
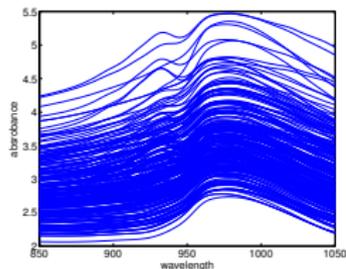
Research Summer School on Statistics for Data Science S4D 2018

June 22, 2018

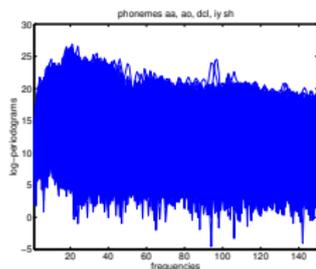
- 1 Motivation
- 2 Model-based co-clustering
- 3 Temporal curve segmentation (RHLP)
- 4 Model-based co-clustering embedding RHLP
- 5 Conclusion and perspectives

# Functional data are increasingly frequent

Chamroukhi and Nguyen [2018]

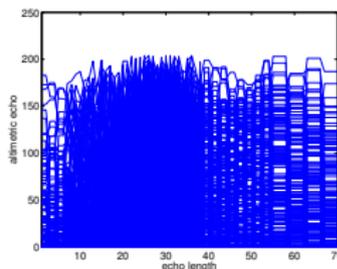


Tecator data



Phonemes curves

Railway switch curves



Satellite waveforms

# Clustering of functional data

Tecator data set<sup>1</sup>:  $n = 240$  spectra with  $m = 100$  observations for each spectrum

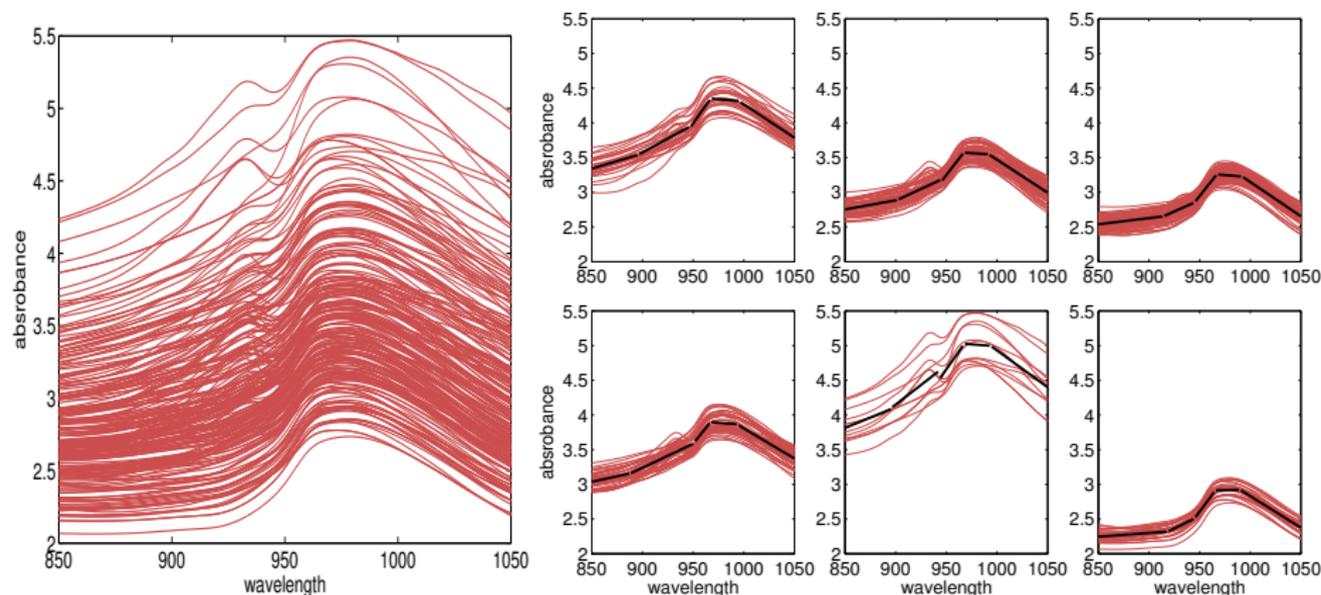


Figure: Original data and clustering results from Chamroukhi [2016] for the data considered in the same setting as in Hébrail et al. [2010] (six clusters, each cluster is approximated by five linear segments ( $R = 5, p = 1$ ))

<sup>1</sup>Tecator data are available at <http://lib.stat.cmu.edu/datasets/tecator>.

# Multivariate functional data are increasingly present

Measurements collected from different network elements (transceivers, cells, sites. . .)

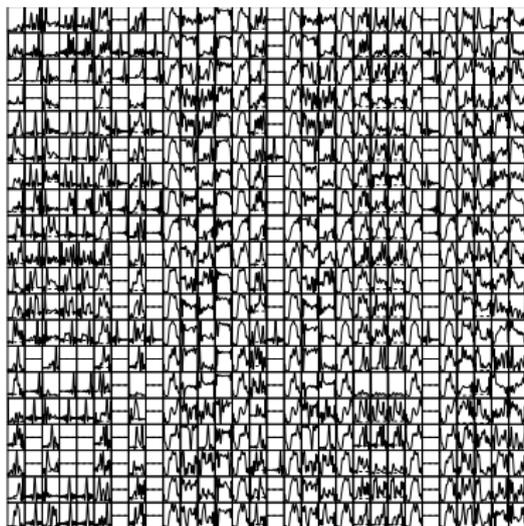


Figure: An example with  $d = 30$  and  $n = 20$  daily observations, from [Ben Slimen et al., 2016].

# This talk

## Questioning

Clustering of highly multivariate functional data with two guidelines:

- (1) Mathematical guideline: warranty for estimation and selection
- (2) User guideline: keep a user-friendly meaning of the process

Both are important because clustering is a highly risky task. . .

## Proposed answering

(1) Model-based co-clustering with (2) temporal curve segmentation

Novelty corresponds to combining both (1) and (2)

- 1 Motivation
- 2 Model-based co-clustering
- 3 Temporal curve segmentation (RHLP)
- 4 Model-based co-clustering embedding RHLP
- 5 Conclusion and perspectives

# Difference between clustering and co-clustering

- Simultaneous clustering of lines/individ. ( $Z$ ) and columns/var. ( $W$ )
- Can be used as a way to reduce dimensionality (var.  $\rightarrow W$ )

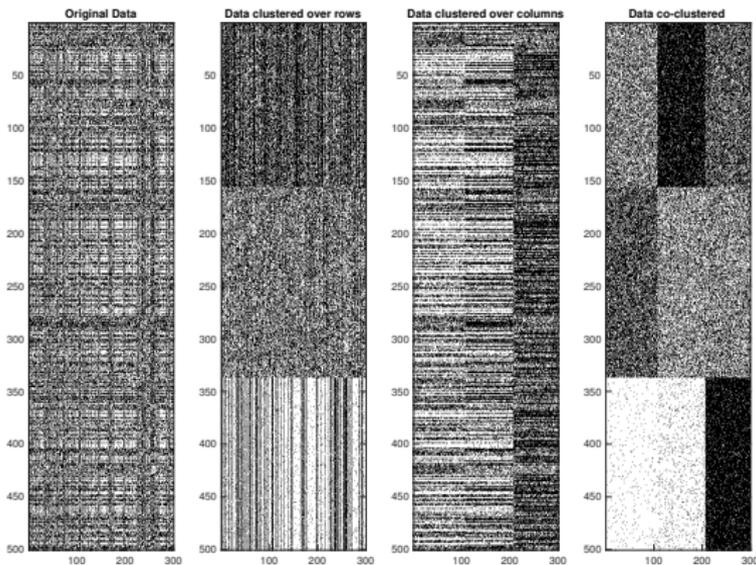


Figure: Binary data set with  $n = 500$ ,  $d = 300$ ,  $K = M = 3$

# Latent block model for co-clustering

$$f(\mathbf{X}; \Psi) = \sum_{(z,w) \in \mathcal{Z} \times \mathcal{W}} \mathbb{P}(\mathbf{Z}, \mathbf{W}; \pi, \rho) \underbrace{f(\mathbf{X} | \mathbf{Z}, \mathbf{W}; \theta)}_{\text{data kind dependent}}$$

- The latent variables  $\mathbf{Z}$  and  $\mathbf{W}$  are independent:  $\mathbb{P}(\mathbf{Z}, \mathbf{W}) = \mathbb{P}(\mathbf{Z})\mathbb{P}(\mathbf{W})$  and iid:  
 $\mathbb{P}(\mathbf{Z}) = \prod_i \mathbb{P}(z_i)$  with  $z_i \sim \text{Multinomial}(\pi_1, \dots, \pi_K)$  where  $\pi_k = \mathbb{P}(z_k = k)$   
 $\mathbb{P}(\mathbf{W}) = \prod_j \mathbb{P}(w_j)$  with  $w_j \sim \text{Multinomial}(\rho_1, \dots, \rho_M)$  where  $\rho_\ell = \mathbb{P}(w_j = \ell)$
- Conditional independence:  $x_{ij} | (z_i, w_j) \perp x_{i'j'} | (z_{i'}, w_{j'})$

↔ binary data: binary [Govaert and Nadif, 2003, 2008; Keribin et al., 2012],

↔ categorical data: multinomial [Keribin et al., 2014]

↔ contingency table: Poisson [Govaert and Nadif, 2003, 2006, 2008]

↔ continuous data: Gaussian [Lomet, 2012; Govaert and Nadif, 2013]

↔ functional data: functional PCA + Gaussian, see further [Ben Slimen et al., 2016]

# Inference for the latent block model

- Parameter estimation: maximum likelihood with variational block EM algorithm (VBEM) [Govaert and Nadif, 2006, 2008]
- Number of blocks estimation: ICL criterion

Recent or ongoing theoretical warranties for consistency

- 1 Motivation
- 2 Model-based co-clustering
- 3 Temporal curve segmentation (RHLP)
- 4 Model-based co-clustering embedding RHLP
- 5 Conclusion and perspectives

# Functional data notation

- Data: (discretized) values of underlying smooth functions, not just vectors
- Data: A sample of  $n$  heterogeneous univariate curves  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$
- $(\mathbf{x}_i, \mathbf{y}_i)$  consists of  $m_i$  observations  $\mathbf{y}_i = (y_{i1}, \dots, y_{im_i})$  observed at the independent covariates, (e.g., time  $t$  in time series),  $(x_{i1}, \dots, x_{im_i})$

# Functional data modeling: “classical” approach

[Ramsay and Silverman, 2005] and many others

- Step 1:  $(\mathbf{x}, \mathbf{y})$  decomposed into a finite basis of function (B-spline. . .)
- Step 2: functional principal components analysis (PCA) which is performed as a usual PCA of the basis expansion coefficients  $\mathbf{c}$  using a metric defined by the inner products between the basis functions
- Step 3: set a distribution of probability on  $\mathbf{c}$ , typically Gaussian

It defines a distribution on  $\mathbf{c}$  instead of  $\mathbf{y}$ . . .

# Functional data modeling: regression RHLP

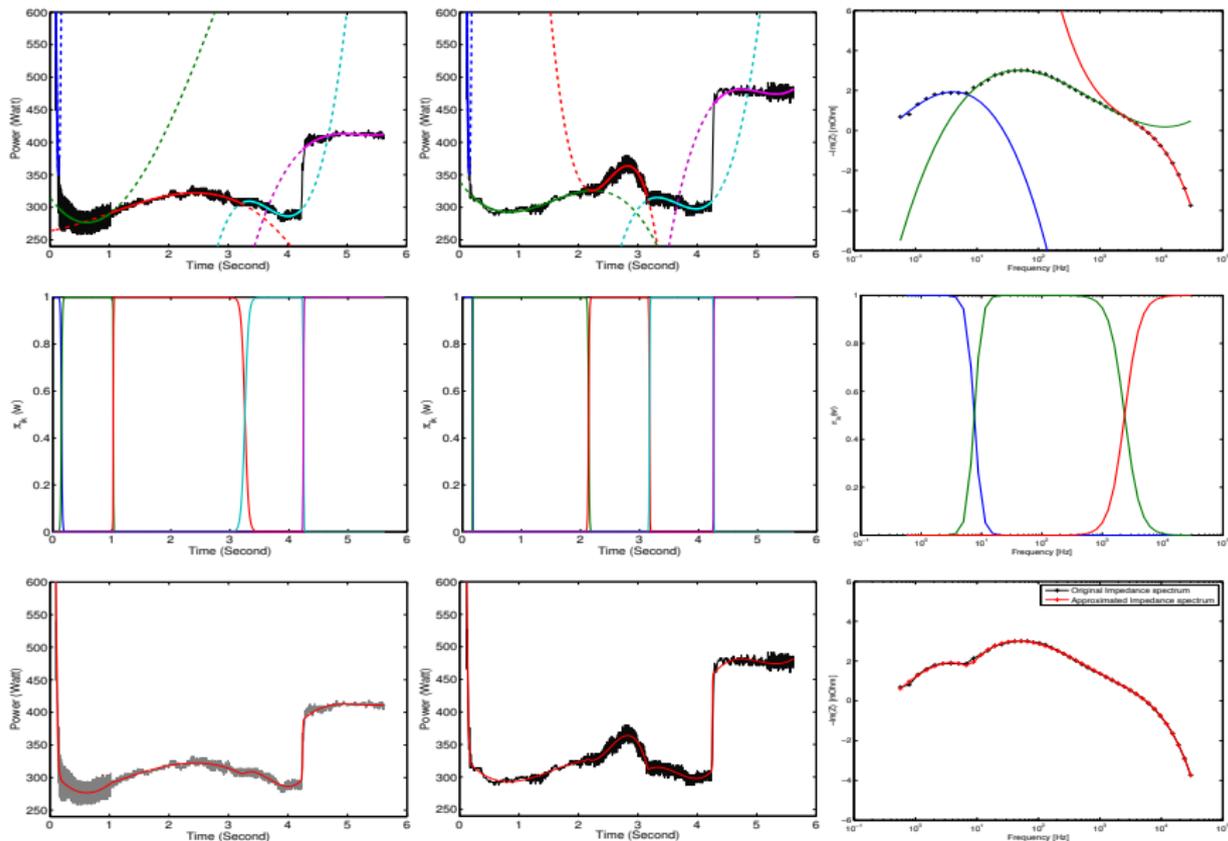
Alternatively, use a segmentation via generative piecewise polynomial regression modeling of  $f(\mathbf{y}|\mathbf{x})$  [Chamroukhi et al.]

↔ Regression with Hidden Logistic Process (RHLP)

↔ See formula later

It gives a distribution on  $\mathbf{y}$  and also a meaningful segmentation of the curve

# RHLP for modeling different kinds of functions

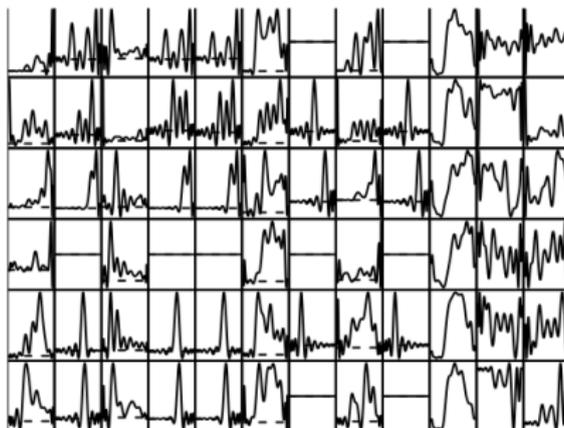


- 1 Motivation
- 2 Model-based co-clustering
- 3 Temporal curve segmentation (RHLP)
- 4 Model-based co-clustering embedding RHLP
- 5 Conclusion and perspectives

# Multivariate functional data

Chamroukhi and Biernacki [2017]

- Data:  $\mathbf{Y} = (\mathbf{y}_{ij})$  a data sample matrix of  $n$  individuals defined on a set  $\mathcal{I}$  and  $d$  continuous functional variables defined on a set  $\mathcal{J}$ .
- Each variable  $\mathbf{y}_{ij}$  is an univariate curve  $\mathbf{y}_{ij} = (y_{ij}(t_1), \dots, y_{ij}(t_{T_{ij}}))$  of  $T_{ij}$  observations  $y(t) \in \mathbb{R}$  linked to covariates  $\mathbf{x}_{ij} = (x_{ij}(t_1), \dots, x_{ij}(t_{T_{ij}}))$  at the points  $(t_1, \dots, t_{T_{ij}})$ , typically a sampling time



# Embedding RHLP in co-clustering

Chamroukhi and Biernacki [2017]

■ Co-clustering:

$$\begin{aligned} f(\mathbf{Y}|\mathbf{X}; \Psi) &= \sum_{(z,w) \in \mathcal{Z} \times \mathcal{W}} \mathbb{P}(\mathbf{Z}; \boldsymbol{\pi}) \mathbb{P}(\mathbf{W}; \boldsymbol{\rho}) f(\mathbf{Y}|\mathbf{X}, \mathbf{Z}, \mathbf{W}; \boldsymbol{\theta}) \\ &= \sum_{(z,w) \in \mathcal{Z} \times \mathcal{W}} \prod_{i,k} \pi_k^{z_{ik}} \prod_{j,\ell} \rho_\ell^{w_{j\ell}} \prod_{i,j,k,\ell} \underbrace{f(\mathbf{y}_{ij}|\mathbf{x}_{ij}; \boldsymbol{\theta}_{k\ell})}_{\text{RHLP}}^{z_{ik}w_{j\ell}}. \end{aligned}$$

with parameter vector  $\Psi = (\boldsymbol{\pi}^T, \boldsymbol{\rho}^T, \boldsymbol{\theta}^T)^T$ , where  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)^T$ ,  $\boldsymbol{\rho} = (\rho_1, \dots, \rho_M)^T$ , and  $\boldsymbol{\theta} = (\boldsymbol{\theta}_{11}^T, \dots, \boldsymbol{\theta}_{k\ell}^T, \dots, \boldsymbol{\theta}_{KM}^T)^T$ .

# Embedding RHLF in co-clustering

- RHLF: model the conditional data distribution for each block  $kl$ , assuming that each functional variable  $\mathbf{y}_{ij}$  is governed by an  $S_{kl}$ -state hidden process of  $y_{ij}$ :

$$f(\mathbf{y}_{ij} | \mathbf{x}_{ij}; \boldsymbol{\theta}_{kl}) = \prod_{t=1}^{T_{ij}} \sum_{r=1}^{S_{kl}} \alpha_{klr}(t; \boldsymbol{\xi}_{kl}) \mathcal{N}(y_{ij}(t); \boldsymbol{\beta}_{klr}^T \mathbf{x}_{ij}(t), \sigma_{klr}^2)$$

where the dynamical weights  $\alpha$ 's are given by the multinomial logistic:

$$\alpha_{klr}(t; \boldsymbol{\xi}_{kl}) = \frac{\exp(\xi_{klr0} + \xi_{klr1}t)}{\sum_{r'=1}^{S_{kl}} \exp(\xi_{klr'0} + \xi_{klr'1}t)}.$$

↔ Can be seen as a generative piecewise polynomial regression model where the transition points are smoothly controlled by logistic weights

# Parameter estimation: EM not feasible

- The complete-data log-likelihood:

$$\begin{aligned}\log L_c(\Psi) &= \log f(\mathbf{Y}, \mathbf{Z}, \mathbf{W}, \mathbf{H} | \mathbf{X}; \Psi) \\ &= \sum_{i,k} z_{ik} \log \pi_k + \sum_{j,\ell} w_{j\ell} \log \rho_\ell \\ &\quad + \sum_{i,j,k,\ell,t,r} z_{ik} w_{j\ell} h_{tr} \log \left[ \alpha_{k\ell r}(t; \boldsymbol{\xi}_{k\ell}) \mathcal{N} \left( y_{ij}(t); \boldsymbol{\beta}_{k\ell r}^T \mathbf{x}_{ij}(t), \sigma_{k\ell r}^2 \right) \right]\end{aligned}$$

where  $(h_{tr}; t = 1, \dots, T_{ij}, r = 1, \dots, S_{k\ell})$  is a binary variable indicating from which state the observation  $y_{ij}(t)$  within the block cluster  $k\ell$  is originated

# Parameter estimation: EM not feasible

- The E-Step computes the expected complete-data log-likelihood, given the observed curves  $(\mathbf{X}, \mathbf{Y})$ , and the current parameter estimation  $\Psi^{(q)}$

$$\begin{aligned} Q(\Psi, \Psi^{(q)}) &= \mathbb{E} \left[ \log L_c(\Psi) \mid \mathbf{X}, \mathbf{Y}; \Psi^{(q)} \right] \\ &= \sum_{i,k} \mathbb{P}(z_{ik} = 1 \mid \mathbf{y}_{ij}, \mathbf{x}_{ij}) \log \pi_k + \sum_{j,\ell} \mathbb{P}(w_{j\ell} = 1 \mid \mathbf{y}_{ij}, \mathbf{x}_{ij}) \log \rho_\ell \\ &\quad + \sum_{i,j,k,\ell,t,r} \mathbb{P}(z_{ik} w_{j\ell} = 1 \mid \mathbf{y}_{ij}, \mathbf{x}_{ij}) \mathbb{P}(h_{tr} = 1 \mid z_{ik}, w_{j\ell}, y_{ij}(t), x_{ij}(t)) \times \\ &\quad \log \left[ \alpha_{k\ell r}(t; \boldsymbol{\xi}_{k\ell}) \mathcal{N} \left( y_{ij}(t); \boldsymbol{\beta}_{k\ell r}^T \mathbf{x}_{ij}(t), \sigma_{k\ell r}^2 \right) \right] \end{aligned}$$

- ↪ Requires the calculation of the posterior joint distribution  $\mathbb{P}(z_{ik} w_{j\ell} = 1 \mid \mathbf{y}_{ij}, \mathbf{x}_{ij})$
- ↪ does not factorize due to the conditional dependence on the observed curves of the row and the column labels
- ⇒ [Govaert and Nadif, 2008, 2013] proposed a variational approximation by relying on the Neal and Hinton's interpretation of the EM algorithm [Neal and Hinton, 1998].
- ↪ We adopt this variational approximation in our context

# Variational block EM algorithm

$$\mathbb{P}(z_{ik}w_{j\ell} = 1|\mathbf{y}_{ij}, \mathbf{x}_{ij}) \approx \mathbb{P}(z_{ik} = 1|\mathbf{y}_{ij}, \mathbf{x}_{ij}) \times \mathbb{P}(w_{j\ell} = 1|\mathbf{y}_{ij}, \mathbf{x}_{ij})$$

**Initialization:** start from an initial solution at iteration  $q = 0$ , and then alternate at the  $(q + 1)$ th iteration between the following variational E- and M- steps until convergence:

**VE Step** Estimate the variational approximated posterior memberships:

- 1  $\tilde{z}_{ik}^{(q+1)} \propto \pi_k^{(q)} \exp\left(\sum_{j,\ell,t,r} \tilde{w}_{j\ell}^{(q)} \tilde{h}_{tr}^{(q)} \log\left[\alpha_{k\ell r}(t; \boldsymbol{\xi}_{k\ell}^{(q)}) \mathcal{N}\left(y_{ij}(t); \boldsymbol{\beta}_{k\ell r}^{T(q)} \mathbf{x}_{ij}(t), \sigma_{k\ell r}^{(q)2}\right)\right]\right)$
- 2  $\tilde{w}_{j\ell}^{(q+1)} \propto \rho_\ell^{(q)} \exp\left(\sum_{i,k,t,r} \tilde{z}_{ik}^{(q)} \tilde{h}_{tr}^{(q)} \log\left[\alpha_{k\ell r}(t; \boldsymbol{\xi}_{k\ell}^{(q)}) \mathcal{N}\left(y_{ij}(t); \boldsymbol{\beta}_{k\ell r}^{T(q)} \mathbf{x}_{ij}(t), \sigma_{k\ell r}^{(q)2}\right)\right]\right)$
- 3  $\tilde{h}_{tr}^{(q+1)} \propto \alpha_{k\ell r}^{(q)}(t; \boldsymbol{\xi}_{k\ell}^{(q)}) \mathcal{N}\left(y_{ij}(t); \boldsymbol{\beta}_{k\ell r}^{(q)T} \mathbf{x}_{ij}(t), \sigma_{k\ell r}^{(q)2}\right)$

where:  $\tilde{z}_{ik} = \mathbb{P}(z_{ik} = 1|\mathbf{y}_{ij}, \mathbf{x}_{ij})$ ,  $\tilde{w}_{j\ell} = \mathbb{P}(w_{j\ell} = 1|\mathbf{y}_{ij}, \mathbf{x}_{ij})$ ,  
 $\tilde{h}_{tr} = \mathbb{P}(h_{tr} = 1|z_i, w_j, y_{ij}(t), x_{ij}(t))$

# Variational block EM algorithm

**M Step** update the parameters estimates  $\theta^{(q+1)}$  given the estimated posterior memberships at the current iteration  $q + 1$ :

$$1 \quad \pi_k^{(q+1)} = \frac{\sum_i \tilde{z}_{ik}^{(q+1)}}{n}$$

$$2 \quad \rho_\ell^{(q+1)} = \frac{\sum_j \tilde{w}_{j\ell}^{(q+1)}}{d}$$

The update of each block parameters  $\theta_{kl}$  consists in a weighted version of the RHLF updating rules:

$$3 \quad \xi_{kl}^{(new)} = \xi_{kl}^{(old)} - \left[ \frac{\partial^2 F(\xi_{kl})}{\partial \xi_{kl} \partial \xi_{kl}^T} \right]_{\xi_{kl} = \xi_{kl}^{(old)}}^{-1} \frac{\partial F(\xi_{kl})}{\partial \xi_{kl}} \Big|_{\xi_{kl} = \xi_{kl}^{(old)}} \text{ which is the IRLS maximisation of } F(\xi_{kl}) = \sum_{i,j,t} \tilde{z}_{ik}^{(q)} \tilde{w}_{j\ell}^{(q)} \tilde{h}_{tr}^{(q)} \log \alpha_{klr}(t; \xi_{kl}) \text{ w.r.t } \xi_{kl}.$$

The regression parameters updates consist in analytic WLS problems:

$$4 \quad \beta_{klr}^{(q+1)} = \left[ \sum_{i,j} \tilde{z}_{ik}^{(q)} \tilde{w}_{j\ell}^{(q)} \mathbf{X}_{ij}^T \Lambda_{ijk_r}^{(q)} \mathbf{X}_{ij} \right]^{-1} \sum_{i,j} \tilde{z}_{ik}^{(q)} \tilde{w}_{j\ell}^{(q)} \mathbf{X}_{ij}^T \Lambda_{ijk_r}^{(q)} \mathbf{y}_{ij}$$

$$5 \quad \sigma_{klr}^{2(q+1)} = \frac{\sum_{i,j} \tilde{z}_{ik}^{(q)} \tilde{w}_{j\ell}^{(q)} \left\| \sqrt{\Lambda_{ijk_r}^{(q)}} (\mathbf{y}_{ij} - \mathbf{X}_{ij} \beta_{klr}^{(q+1)}) \right\|^2}{\sum_{i,j} \tilde{z}_{ik}^{(q)} \tilde{w}_{j\ell}^{(q)} \text{trace}(\Lambda_{ijk_r}^{(q)})} \text{ where } \mathbf{X}_{ij} \text{ is the design matrix for}$$

the  $i$ th curve,  $\Lambda_{ijk_r}^{(q)}$  is the diagonal matrix whose diagonal elements are the posterior segment memberships  $\{\tilde{h}_{ijtr}^{(q)}; t = 1, \dots, T_{ij}\}$ .

- 1 Motivation
- 2 Model-based co-clustering
- 3 Temporal curve segmentation (RHLP)
- 4 Model-based co-clustering embedding RHLP
- 5 Conclusion and perspectives

# Conclusion and perspectives

## Conclusion

- A full generative framework for the cluster analysis and segmentation of high-dimensional non-stationary functional data
- The model inference can be performed by a variational EM algorithm

## Perspectives

- Replace Variational EM by an SEM algorithm, which does not use approximation
- Numerical experiments
- Package

# References I

- Y. Ben Slimen, S. Allio, and J. Jacques. Model-Based Co-clustering for Functional Data. HAL preprint hal-01422756, December 2016. URL <https://hal.inria.fr/hal-01422756>.
- F. Chamroukhi. Piecewise regression mixture for simultaneous curve clustering and optimal segmentation. *Journal of Classification - Springer*, 33(3):374–411, 2016.
- F. Chamroukhi and C. Biernacki. Model-Based Co-Clustering of Multivariate Functional Data. In *ISI 2017 - 61st World Statistics Congress*, Marrakech, Morocco, Jul 2017. URL <https://hal.archives-ouvertes.fr/hal-01653782>.
- Faïcel Chamroukhi and Hien D. Nguyen. Model-based clustering and classification of functional data. 2018. URL <https://chamroukhi.users.lmno.cnrs.fr/papers/MBCC-FDA.pdf>. arXiv:1803.00276v2.
- G. Govaert and M. Nadif. Clustering with block mixture models. *Pattern Recognition*, 36(2):463 – 473, 2003. Biometrics.
- G. Govaert and M. Nadif. Fuzzy clustering to estimate the parameters of block mixture models. *Soft Computing*, 10(5): 415–422, 2006.
- G. Govaert and M. Nadif. Block clustering with Bernoulli mixture models: Comparison of different approaches. *Computational Statistics and Data Analysis*, 52(6):3233 –3245, 2008.
- G. Govaert and M. Nadif. *Co-Clustering*. Computer engineering series. Wiley-ISTE, November 2013. 256 pages.
- G. Hébrail, B. Huguene, Y. Lechevallier, and F. Rossi. Exploratory analysis of functional data via clustering and optimal segmentation. *Neurocomputing*, 73(7-9):1125–1141, March 2010.
- C. Keribin, V. Brault, G. Celeux, and G. Govaert. Model selection for the binary latent block model. In *Proceedings of COMPSTAT*, 2012.
- C. Keribin, V. Brault, G. Celeux, and G. Govaert. Estimation and selection for the latent block model on categorical data. *Statistics and Computing*, pages 1–16, 2014. ISSN 0960-3174. doi: 10.1007/s11222-014-9472-2. URL <http://dx.doi.org/10.1007/s11222-014-9472-2>.
- A. Lomet. *Sélection de modèle pour la classification croisée de données continues*. Ph.D. thesis, Université de Technologie de Compiègne, 2012.
- R. Neal and G. E. Hinton. *A view of the EM algorithm that justifies incremental, sparse, and other variants*, pages 355–368. Dordrecht: Kluwer Academic Publishers, 1998.
- J. O. Ramsay and B. W. Silverman. *Functional Data Analysis*. Springer Series in Statistics. Springer, June 2005.

Thank you for your attention!