Statistical data science and some unsupervised learning problems

FAICEL CHAMROUKHI

https://chamroukhi.com/



Workshop DSCI, 22-24 dec 2018





Statistical data science and some unsupervised learning problems

Outline

- 1 Statistics and Data Science
- 2 Clustering of multivariate data
- 3 Unsupervised learning for dimensionality reduction
- 4 Time series segmentaion
- 5 Clustering of functional data
- 6 Unsupervised Bayesian (non-)parametric learning
- 7 Model-Based Co-Clustering of Multivariate Functional Data

- The term "Data Science" has surged in popularity
- Data science is increasingly commonly used with "big data."
- Data science, including Big Data has recently attracted an enormous interest from the scientific community







Data Scientist: The Sexiest Job of the 21st Century



Faicel Chamroukh

- What does Data Science mean?
- What about Statistics in the Data Science "area" ?
- There is not yet a consensus on what precisely constitutes Data Science





Serge Abiteboul, directeur de recherche Inria, École normale supérieure de Cachan, membre de l'Académie des sciences et Patrick Flandrin, directeur de recherche CNRS, École normale supérieure de Lvon, membre de l'Académie des sciences.



À la découverte des connaissances massives de la Toile Serge Abiteboul, directeur de recherche Inria, École normale supérieure de Cachan, membre de l'Académie des sciences

Des mathématiques pour l'analyse de données massives Stéphane Mallat, professeur à l'École normale supérieure, Paris





Big Data et Relation Client : quel impact sur les industries et activités de services traditionnelles ?

François Bourdoncle, co-fondateur et CTO d'Exalead, filiale de Dassault Systèmes





Vidéos réalisées par la cellule Webcast CC-IN2P3 du CNRS Stovers 600

- There is not yet a consensus on what precisely constitutes Data Science, but
- Data Science can be seen (defined ?) as ^a :
 - ▶ the study of the generalizable extraction of knowledge from data.
 - requires an integrated skill set spanning mathematics, machine learning, artificial intelligence, statistics, databases, and optimization
- a. Vasant Dhar (2013) : Communications of the ACM, Vol. 56 No. 12 : 64-73
 - Data Science clearly has an interdisciplinary nature and requires substantial collaborative effort
 - Databases, statistics and machine learning, and distributed systems are emerging as foundational to data science
- (i) Databases : organization of data resources,
- (ii) Statistics and Machine Learning : convert data into knowledge,
- (iii) Distributed and Parallel Systems : computational infrastructure

Statistics and Data Science

- \hookrightarrow Statistics play a central role in data science
 - Allow to quantify the randomness component in the data
 - A well-established background to deal with uncertainty (probabilistic frame- work) and to establish generizable methods for prediction and estimation
 - allow soft decision : e.g. confidence interval in regression and posterior probabilities in classification
 - help for understanding the underlying generative process

Statistical modeling for data science

- The observed data (x_1, \ldots, x_n) where $x_i \in \mathcal{X} \subseteq \mathbb{R}^d$ are assumed to represent samples from random variables X with unknown probability distribution f
- The main questions are i) how to define flexible and generic models for f ii) construct estimators with desirable properties to learn f from the data iii) to deal with the computational and practical issues for "complex" data
- The area of statistical learning for the analysis of complex data.

Context and Objectives

- **Context**: Large-scale data are increasingly frequent : Complex data *heterogeneous, dynamical (temporal, functional), incomplete, high-dimension, and possibly massive*
- **Objectives :** learn/discern useful information in an unsupervised way from raw data :

 \hookrightarrow Reconstruct/reveal hidden structures, i.e, (hierarchy of) groups; learn/select relevant features, etc









Outline

1 Statistics and Data Science

- 2 Clustering of multivariate data
- 3 Unsupervised learning for dimensionality reduction
- 4 Time series segmentaion
- **5** Clustering of functional data
- 6 Unsupervised Bayesian (non-)parametric learning
- 7 Model-Based Co-Clustering of Multivariate Functional Data

Clustering of multivariate data



FAICEL CHAMROUKHI Statistical data science and some unsupervised learning problem

Clustering of multivariate data



FAICEL CHAMROUKHI Statistical data science and some unsupervised learning probler

K-means

- a straightforward and widely used clustering algorithm, is one of the most important algorithms in unsupervised learning.
- Each cluster is represented by its mean (cluster centroid) μ_k in \mathbb{R}^d .

K-means MacQueen [1967]

$$(\widehat{\boldsymbol{\mu}}_1,\ldots,\widehat{\boldsymbol{\mu}}_K,\widehat{\mathbf{z}}) \in \arg\min_{\boldsymbol{\mu}_1,\ldots,\boldsymbol{\mu}_K,\mathbf{z}} \mathcal{J}(\boldsymbol{\mu}_1,\ldots,\boldsymbol{\mu}_K,\mathbf{z})$$

objective function : $\mathcal{J}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \mathbf{z}) = \sum_{k=1}^K \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{z_i}\|^2$

K-means

- a straightforward and widely used clustering algorithm, is one of the most important algorithms in unsupervised learning.
- Each cluster is represented by its mean (cluster centroid) μ_k in \mathbb{R}^d .

K-means MacQueen [1967]

$$(\widehat{oldsymbol{\mu}}_1,\ldots,\widehat{oldsymbol{\mu}}_K,\widehat{\mathbf{z}})\inrg\min_{oldsymbol{\mu}_1,\ldots,oldsymbol{\mu}_K,\mathbf{z}}\mathcal{J}(oldsymbol{\mu}_1,\ldots,oldsymbol{\mu}_K,\mathbf{z})$$

objective function :
$$\mathcal{J}(\boldsymbol{\mu}_1,\ldots,\boldsymbol{\mu}_K,\mathbf{z}) = \sum_{k=1}^K \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_{z_i}\|^2$$

Initialization : $(\boldsymbol{\mu}_1^{(0)},\ldots,\boldsymbol{\mu}_K^{(0)})$ (eg, randomly chosen data points)

1 Assignment step : $z_i^{(q)} = \arg\min_{z \in \mathcal{Z}} \|\mathbf{x}_i - \boldsymbol{\mu}_z\|^2$ 2 Relocation step : $\boldsymbol{\mu}_k^{(q+1)} = \frac{\sum_{i=1}^n z_{ik}^{(q)} \mathbf{x}_i}{\sum_{i=1}^n z_{ik}^{(q)}},$

 \Rightarrow The K-means algorithm is simple to implement and relatively fast.

















K-means

How to measure uncertainty?



FIGURE – K-means partition (left) vs GMM-EM partition (right).

Finite Mixture Models [McLachlan and Peel., 2000]

 $f(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_{k=1}^{K} \pi_k f_k(\boldsymbol{x}; \boldsymbol{\theta}_k)$ with $\pi_k > 0 \ \forall k$ and $\sum_{k=1}^{K} \pi_k = 1$.

Finite Mixture Models [McLachlan and Peel., 2000]

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_{k=1}^{K} \pi_k f_k(\boldsymbol{x}; \boldsymbol{\theta}_k)$$
 with $\pi_k > 0 \ \forall k$ and $\sum_{k=1}^{K} \pi_k = 1$.

Maximum-Likelihood Estimation

 $\widehat{\boldsymbol{\theta}} \in \arg \max_{\boldsymbol{\theta}} \ln L(\boldsymbol{\theta}) \\ \text{log-likelihood} : \ln L(\boldsymbol{\theta}) = \sum_{i=1}^{n} \log \sum_{k=1}^{K} \pi_k f_k(\boldsymbol{x}_i; \boldsymbol{\theta}_k).$

Finite Mixture Models [McLachlan and Peel., 2000]

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_{k=1}^{K} \pi_k f_k(\boldsymbol{x}; \boldsymbol{\theta}_k)$$
 with $\pi_k > 0 \ \forall k$ and $\sum_{k=1}^{K} \pi_k = 1$.

Maximum-Likelihood Estimation

$$\widehat{\boldsymbol{\theta}} \in \arg \max_{\boldsymbol{\theta}} \ln L(\boldsymbol{\theta})$$
log-likelihood : $\ln L(\boldsymbol{\theta}) = \sum_{i=1}^{n} \log \sum_{k=1}^{K} \pi_k f_k(\boldsymbol{x}_i; \boldsymbol{\theta}_k).$

The EM algorithm [Dempster et al., 1977]

$$oldsymbol{ heta}^{new} \in rg\max_{oldsymbol{ heta}\in\Omega} \mathbb{E}[\ln L_c(oldsymbol{ heta}) | \mathcal{D}, oldsymbol{ heta}^{old}]$$

complete log-likelihood : $\ln L_c(\boldsymbol{\theta}) = \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log [\pi_k f_k(\boldsymbol{x}_i; \boldsymbol{\theta}_k)]$ where Z_{ik} is such that $Z_{ik} = 1$ if $Z_i = k$ and $Z_{ik} = 0$ otherwise.

Finite Mixture Models [McLachlan and Peel., 2000]

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_{k=1}^{K} \pi_k f_k(\boldsymbol{x}; \boldsymbol{\theta}_k)$$
 with $\pi_k > 0 \; \forall k \text{ and } \sum_{k=1}^{K} \pi_k = 1.$

Maximum-Likelihood Estimation

$$\widehat{\boldsymbol{\theta}} \in \arg \max_{\boldsymbol{\theta}} \ln L(\boldsymbol{\theta})$$
log-likelihood : $\ln L(\boldsymbol{\theta}) = \sum_{i=1}^{n} \log \sum_{k=1}^{K} \pi_k f_k(\boldsymbol{x}_i; \boldsymbol{\theta}_k).$

The EM algorithm [Dempster et al., 1977]

$$oldsymbol{ heta}^{new} \in rg\max_{oldsymbol{ heta}\in\Omega} \mathbb{E}[\ln L_c(oldsymbol{ heta}) | \mathcal{D}, oldsymbol{ heta}^{old}]$$

complete log-likelihood : $\ln L_c(\boldsymbol{\theta}) = \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log [\pi_k f_k(\boldsymbol{x}_i; \boldsymbol{\theta}_k)]$ where Z_{ik} is such that $Z_{ik} = 1$ if $Z_i = k$ and $Z_{ik} = 0$ otherwise.

Clustering

$$\widehat{z}_i = \arg \max_{1 \le k \le K} \mathbb{P}(Z_i = k | \boldsymbol{x}_i; \widehat{\boldsymbol{\theta}}), \quad (i = 1, \dots, n)$$

Gaussian mixture models (GMMs)

The finite Gaussian mixture density is defined as :

$$f(\mathbf{x}_i; \mathbf{\Psi}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



 Figure – An example of a three-component Gaussian mixture density in $\mathbb{R}^2.$

EM for Gaussian mixture models

1 E-Step : calculates the posterior component memberships :

$$\tau_{ik}^{(q)} = \mathbb{P}(Z_i = k | \mathbf{x}_i, \boldsymbol{\Psi}^{(q)}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_k^{(q)}, \boldsymbol{\Sigma}_k^{(q)})}{\sum_{\ell=1}^K \pi_\ell \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_\ell^{(q)}, \boldsymbol{\Sigma}_\ell^{(q)})}$$

that \mathbf{x}_i originates from the *k*th component density.

2 M-Step : parameter updates :

$$\begin{aligned} \pi_k^{(q+1)} &= \frac{\sum_{i=1}^n \tau_{ik}^{(q)}}{n} = \frac{n_k^{(q)}}{n}, \\ \mu_k^{(q+1)} &= \frac{1}{n_k^{(q)}} \sum_{i=1}^n \tau_{ik}^{(q)} \mathbf{x}_i, \\ \mathbf{\Sigma}_k^{(q+1)} &= \frac{1}{n_k^{(q)}} \sum_{i=1}^n \tau_{ik}^{(q)} (\mathbf{x}_i - \boldsymbol{\mu}^{(q+1)}) (\mathbf{x}_i - \boldsymbol{\mu}^{(q+1)})^T. \end{aligned}$$



 $\rm FIGURE$ – A three-class example of a real data set : Iris data of Fisher.



 $\rm FIGURE$ – Iris data : Clustering results with EM for a GMM and AIC.



 $\rm FIGURE$ – Iris data of Fisher : The data are colored according to the true partition.














































Parsimonious GMMs for high-dimensional data

- Parsimonious Gaussian mixture models¹are statistical models that allow for capturing a specific cluster shapes (e.g., clusters having the same shape or different shapes, spherical or elliptical clusters, etc).
- Eigenvalue decomposition of the cluster covariance matrices :

$$\mathbf{\Sigma}_k = \lambda_k \mathbf{D}_k \mathbf{A}_k \mathbf{D}_k^T$$

where

- ► λ_k represents the volume of the kth cluster (the amount of space of the cluster).
- ► D_k is a matrix with columns corresponding to the eigenvectors of Σ_k that determines the orientation of the cluster.
- A_k is a diagonal matrix, whose diagonal entries are the normalized eigenvalues of Σ_k arranged in a decreasing order and its determinant is 1. This matrix is associated with the shape of the cluster.

1. Banfield and Raftery [1993], ?

Parsimonious GMMs for high-dimensional data



Model selection

- The problem of choosing the number of clusters can be seen as a model selection problem.
- The model selection task consists of choosing a suitable compromise between flexibility so that a reasonable fit to the available data is obtained, and over-fitting.
- A common way is to use a criterion (score function) that ensure the compromise.

score(model) = error(model) + penalty(model complexity)

which will be minimized.

 Here the complexity of a model *M* is related to the number of its (free) parameters *v*

Model selection

• Akaike Information Criterion (AIC) :

$$\operatorname{AIC}(\mathcal{M}_m) = \ln L(\widehat{\boldsymbol{\theta}}_m) - \nu_m$$

• Bayesian Information Criterion (BIC) :

$$\mathsf{BIC}(\mathcal{M}_m) = \ln L(\widehat{\boldsymbol{\theta}}_m) - \frac{\nu_m \log(n)}{2}$$

• Integrated Classification Likelihood (ICL) :

$$\mathsf{ICL}(\mathcal{M}_m) = \ln L_c(\widehat{\theta}_m) - \frac{\nu_m \log(n)}{2}$$

where $\ln L_c(\widehat{\theta}_m)$ is the complete-data log-likelihood for the model \mathcal{M}_m and ν_m denotes the number of free model parameters. For example, in the case of a *d*-dimensional Gaussian mixture model we have :

$$\nu = \underbrace{(K-1)}_{\pi_k'\mathsf{s}} + \underbrace{K \times d}_{\{\boldsymbol{\mu}_k\}} + \underbrace{K \times \frac{d \times (d+1)}{2}}_{\{\boldsymbol{\Sigma}_k\}} = \frac{K \times (d+1) \times (d+2)}{2} - 1.$$



FIGURE – Clustering results obtained with K-means algorithm (left) with K = 2 and the EM algorithm (right). The cluster centers are shown by the red and blue crosses and the ellipses are the contours of the Gaussian component densities at level 0.4 estimated by EM. The number of clusters for EM have been chosen by BIC for $K = 1, \ldots, 4$.



 $\rm FIGURE$ – A three-class example of a real data set : Iris data of Fisher.



 $\rm FIGURE$ – Iris data : Clustering results with EM for a GMM and AIC.



 $\rm FIGURE$ – Iris data of Fisher : The data are colored according to the true partition.

Latent data models for dimensionality reduction

- Dimensionality reduction for high dimensional data (for representation/visualization etc)
- Principal Component Analysis (PCA) [Pearson, 1901, Hotelling, 1933],
- Probabilistic PCA [Tipping and Bishop, 1997, 1999, Roweis, 1998]
- Factor Analysis (FA)[Spearman, 1904, Thurstone, 1947],

Principal Component Analysis (PCA)

 PCA is a linear projection which maximizes the variance in the projected space [Hotelling, 1933].

Consider a sample $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ with $\mathbf{x}_i \in \mathbb{R}^d$.

 \Rightarrow The aim is to project the data onto a space having dimensionality M < d while maximizing the variance of the projected data.

Consider the sample mean vector and the sample covariance matrix : $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i$ and $\mathbf{S} = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T$. \Rightarrow The variance of the projected data is therefore given by the scalar :

$$v(\mathbf{u}) = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{u}^{T} \mathbf{x}_{i} - \mathbf{u}^{T} \bar{\mathbf{x}}) (\mathbf{u}^{T} \mathbf{x}_{i} - \mathbf{u}^{T} \bar{\mathbf{x}})^{T} = \mathbf{u}^{T} \mathbf{S} \mathbf{u} \cdot$$
(1)

The principal axes (the direction vectors) are then given by :

$$\mathbf{u} = \arg \max_{\mathbf{u} \in \mathbb{R}^d} \mathbf{u}^T \mathbf{S} \mathbf{u}$$
(2)

subject to
$$\mathbf{u}^T \mathbf{u} = 1$$
 and $\mathbf{u}_j^T \mathbf{u}_k = 0$ for $j \neq k$



Disadvantage :

The absence of a probability density model and associated likelihood measure.

Deriving PCA from the perspective of density estimation would offer a number of important advantages, including the following :

- The likelihood measure allows comparison with other density models
- We can derive EM for PCA and hence deal with possible missing values
- Possibility to perform Bayesian inference (e.g. for model selection)
- Possibility of computing the the posterior class probabilities if PCA is used to model the class-conditional densities in classification,
- The value of the probability density function would give a measure of the novelty of a new data point.
- PCA model could be extended to a mixture framework.
- ⇒ Use Probabilistic Principal Component Analysis (PPCA)

Probabilistic Principal Component Analysis (PPCA)

Latent variable model for PPCA [Tipping and Bishop, 1997, 1999, Roweis, 1998] :



$$\begin{split} \mathbf{x}_i &= \mathbf{W} \mathbf{z}_i + \boldsymbol{\mu} + \boldsymbol{\epsilon}_i \text{ Observed data} = \text{linear transf. of } \mathbf{z} + \text{additive Gaussian noise} \\ \mathbf{z}_i &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \text{ latent variables of the principal component subspace} \\ \boldsymbol{\epsilon} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \text{ zero-mean Gaussian noise} \\ \hline \mathbf{x}_i | \mathbf{z}_i &\sim \mathcal{N}(\mathbf{W} \mathbf{z}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \text{ conditional density for the observed data} \end{split}$$

 $\mathbf{x}_i | \mathbf{z}_i \sim \mathcal{N}(\mathbf{w} \mathbf{z}_i + \boldsymbol{\mu}, \sigma^{-1})$ conditional density for the observed data $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})$ marginal density for the observed data

 $\hookrightarrow (\mathbf{W}, \mu, \sigma^2) : \text{Parameter estimation using EM [Tipping and Bishop, 1997, 1999, Roweis, 1998]}$
EM for PPCA

NB : for μ , we get its closed form solution : $\hat{\mu} = \bar{\mathbf{x}}$ Only \mathbf{W} and σ^2 are computed in an iterative way by EM

1 E-step : By using the old parameters values, compute

$$\mathbb{E}[\mathbf{z}_i] = (\mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^T (\mathbf{x}_i - \bar{\mathbf{x}})$$
(3)

$$\mathbb{E}[\mathbf{z}_i \mathbf{z}_i^T] = \sigma^2 (\mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I})^{-1} + \mathbb{E}[\mathbf{z}_i] \mathbb{E}[\mathbf{z}_i]^T$$
(4)

2 M-step

$$\mathbf{W}_{\mathsf{new}} = \left[\sum_{i=1}^{n} (\mathbf{x}_{i} - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_{i}]^{T}\right] \left[\sum_{i=1}^{n} \mathbb{E}[\mathbf{z}_{i}\mathbf{z}_{i}^{T}]\right]^{-1}$$
(5)
$$\sigma_{\mathsf{new}}^{2} = \frac{1}{nd} \sum_{i=1}^{n} \left\{ \|\mathbf{x}_{i} - \bar{\mathbf{x}}\|^{2} - 2\mathbb{E}[\mathbf{z}_{i}]^{T} \mathbf{W}_{\mathsf{new}}^{T} (\mathbf{x}_{i} - \bar{\mathbf{x}}) + trace(\mathbb{E}[\mathbf{z}_{i}\mathbf{z}_{i}^{T}] \mathbf{W}_{\mathsf{new}} \mathbf{W}_{\mathsf{new}}^{T} \right\}$$
(6)

NB. Here $\mathbb{E}[.]$ is actually $\mathbb{E}[.|\mathbf{X}, {\{\mathbf{W}, \boldsymbol{\mu}, \sigma^2\}}_{old}]$

Ι

Factor Analysis (FA) I

Factor Analysis (FA) [Spearman, 1904, Thurstone, 1947] FA is closely related to PPCA The only difference is

 $\mathbf{x}_i | \mathbf{z}_i \sim \mathcal{N}(\mathbf{W} \mathbf{z}_i + oldsymbol{\mu}, oldsymbol{\Psi})$ conditional density for the observed data

 ${f \Psi}$ is a d imes d digonal matrix; rather than

 $\mathbf{x}_i | \mathbf{z}_i \sim \mathcal{N}(\mathbf{W} \mathbf{z}_i + oldsymbol{\mu}, \sigma^2 \mathbf{I})$ conditional density for the observed data

(isotropic covariance matrix).

Factor Analysis (FA) II

Generative model

- $\mathbf{x}_i = \mathbf{W}\mathbf{z}_i + \boldsymbol{\mu} + \boldsymbol{\epsilon}_i$ Observed data = linear transf. of \mathbf{z} + additive Gaussian noise $\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$ latent variables of the principal component subspace
 - $\epsilon~\sim~\mathcal{N}(\mathbf{0},\mathbf{I})$ zero-mean Gaussian noise

 $\begin{array}{rcl} \mathbf{x}_i | \mathbf{z}_i & \sim & \mathcal{N}(\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) \text{ conditional density for the observed data} \\ \mathbf{x}_i & \sim & \mathcal{N}(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\Psi}) \text{ marginal density for the observed data} \end{array}$



EM for Facotr Analysis

1 E-step

$$\mathbb{E}[\mathbf{z}_i] = (\mathbf{I} + \mathbf{W}^T \mathbf{\Psi}^{-1} \mathbf{W})^{-1} \mathbf{W}^T (\mathbf{\Psi}^{-1} \mathbf{x}_i - \bar{\mathbf{x}})$$
(7)

$$\mathbb{E}[\mathbf{z}_i \mathbf{z}_i^T] = (\mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} + \mathbb{E}[\mathbf{z}_i] \mathbb{E}[\mathbf{z}_i]^T$$
(8)

2 M-step

$$\mathbf{W}_{\mathsf{new}} = \left[\sum_{i=1}^{n} (\mathbf{x}_{i} - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_{i}]^{T}\right] \left[\sum_{i=1}^{n} \mathbb{E}[\mathbf{z}_{i} \mathbf{z}_{i}^{T}]\right]^{-1}$$
(9)
$$\mathbf{\Psi}_{\mathsf{new}} = \mathsf{diag} \left\{ \mathbf{S} - \mathbf{W}_{\mathsf{new}} \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}[\mathbf{z}_{i}] (\mathbf{x}_{i} - \bar{\mathbf{x}})^{T} \right\}$$
(10)

NB. Here $\mathbb{E}[.]$ is actually $\mathbb{E}[.|\mathbf{X}, \{\mathbf{W}, \boldsymbol{\mu}, \boldsymbol{\Psi}\}_{\mathsf{old}}]$

Time series segmentation

- 1 Statistics and Data Science
- 2 Clustering of multivariate data
- 3 Unsupervised learning for dimensionality reduction
- 4 Time series segmentaion
- **5** Clustering of functional data
- 6 Unsupervised Bayesian (non-)parametric learning
- 7 Model-Based Co-Clustering of Multivariate Functional Data

Regression with hidden logistic process

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

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

$$y_i = \boldsymbol{\beta}_{\boldsymbol{z}_i}^T \boldsymbol{x}_i + \sigma_{\boldsymbol{z}_i} \boldsymbol{\epsilon}_i \quad ; \quad \boldsymbol{\epsilon}_i \sim \mathcal{N}(0, 1), \quad (i = 1, \dots, n)$$

$$Z_i \sim \mathcal{M}(1, \pi_1(t_i; \mathbf{w}), \dots, \pi_K(t_i; \mathbf{w}))$$

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

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

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

Both the mixing proportions and the component parameters are time-varying

Parameter estimation via a the EM algorithm : EM-RHLP

EM-RHLP

Parameter estimation via a the EM algorithm : EM-RHLP

Parameter estimation via a the EM algorithm (EM-RHLP)

 $\label{eq:M-Step:includes a weighted logistic regression problem} \hookrightarrow \mathsf{IRLS} (and weighted polynomial regressions)$

• EM-RHLP algorithm complexity : $\mathcal{O}(I_{\text{EM}}I_{\text{IRLS}}K^3p^3n)$ (more advantageous than dynamic programming).

Time series approximation and segmentation

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

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

2 Curve segmentation : $\hat{z}_i = \arg \max_k \mathbb{E}[z_i|t_i; \hat{\mathbf{w}}] = \arg \max_k \pi_k(t_i; \hat{\mathbf{w}})$

Model selection : Application of BIC, ICL ($\nu_{\theta} = K(p+4) - 2$.)

Application to temporal data modeling and segmentation



Joint segmentation of multivariate time series

Multiple hidden process regression

- Data : $(\boldsymbol{y}_1, \dots, \boldsymbol{y}_n)$ a time series of n multidimensional observations $\boldsymbol{y}_i = (y_i^{(1)}, \dots, y_i^{(d)})^T \in \mathbb{R}^d$ observed at instants $\mathbf{t} = (t_1, \dots, t_n)$.
- Model $\boldsymbol{y}_i = \mathbf{B}_{\boldsymbol{z}_i}^T \boldsymbol{x}_i + \mathbf{e}_i$; $\mathbf{e}_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{z}_i}), \quad (i = 1, \dots, n)$

 $\mathbf{z} = (z_1, \dots, z)$ A latent process generating the data

 \hookrightarrow Multiple Hidden Markov model regression (MHMMR) [7]

Application to human activity time series

Outline

- 1 Statistics and Data Science
- 2 Clustering of multivariate data
- 3 Unsupervised learning for dimensionality reduction
- 4 Time series segmentaion
- 5 Clustering of functional data
- 6 Unsupervised Bayesian (non-)parametric learning
- 7 Model-Based Co-Clustering of Multivariate Functional Data

Functional data are increasingly frequent

[James and Hastie, 2001, James and Sugar, 2003] [Ramsay and Silverman, 2005] [Chamroukhi et al., 2010] [Bouveyron and Jacques, 2011] [Samé et al., 2011] [Jacques and Preda, 2014] [Bouveyron et al., 2018] [Chamroukhi and Nguyen, 2018]



Statistical data science and some unsupervised learning problems

High-dimensional FDA by clustering/segmentation

Non-stationary time series/functions



Railway curves

Satellite waveforms

Objectives

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

Functional data clustering



Functional data clustering



Functional data analysis context

Data

- The individuals are entire functions (e.g., curves, surfaces)
- A set of n univariate curves $(({\boldsymbol{x}}_1, {\boldsymbol{y}}_1), \ldots, ({\boldsymbol{x}}_n, {\boldsymbol{y}}_n)$

• (x_i, y_i) consists of m_i observations $y_i = (y_{i1}, \ldots, y_{im_i})$ observed at the independent covariates, (e.g., time t in time series), $(x_{i1}, \ldots, x_{im_i})$

Objectives : exploratory or decisional

- 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)
 - \Rightarrow Mixture-model based cluster and discriminant analyzes

Mixture modeling framework for functional data

The functional mixture model :

$$f(oldsymbol{y}|oldsymbol{x};oldsymbol{\Psi}) = \sum_{k=1}^K lpha_k f_k(oldsymbol{y}|oldsymbol{x};oldsymbol{\Psi}_k)$$

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

 \hookrightarrow more tailored to approximate smooth functions

 \hookrightarrow do not account for segmentation

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

- 1 Riecewise regression model (PWR)
- 2 Regression model with a hidden Markov process (HMMR)
- 3 Regression model with hidden logistic process (RHLP)

Piecewise regression mixture model (PWRM) [9]

A probabilistic version of the K-means-like approach of [Hébrail et al., 2010]

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

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

 $\blacksquare\,\hookrightarrow\,\mathsf{Simultaneously}$ accounts for curve clustering and segmentation

Parameter estimation

1 Maximum likelihood estimation : EM-PWRM

2 Maximum classification likelihood estimation : CEM-PWRM

 \hookrightarrow a generalization of the K-means-like algorithm of Hébrail et al. [2010] :

 $\textbf{M-step}: \textsf{includes wighted piecewise regressions} \hookrightarrow \textsf{dynamic programming}$

Complexity in $\mathcal{O}(I_{\text{EM}}KRnm^2p^3)$: An issue for large m

Curve clustering : $\hat{z}_i = \arg \max_k \tau_{ik}(\hat{\Psi})$ with $\tau_{ik}(\hat{\Psi}) = \mathbb{P}(Z_i | \boldsymbol{x}_i, \boldsymbol{y}_i; \hat{\Psi})$

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)



EL CHAMROUKHI Statistical data science and some unsupervised learning problems

Application to Topex/Poseidon satellite data

The Topex/Poseidon radar satellite data ² 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 Hébrail et al. [2010].



2. Satellite data are available at

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

FAICEL CHAMROUKHI Statistical data science and some unsupervised learning problems

CEM-PWRM clustering of the satellite data



Statistical data science and some unsupervised learning problems

Mixture of hidden logistic process regressions [4]

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

$$f(\boldsymbol{y}_{i}|\boldsymbol{x}_{i};\boldsymbol{\Psi}) = \sum_{k=1}^{K} \alpha_{k} \underbrace{\prod_{j=1}^{m_{i}} \sum_{r=1}^{R_{k}} \pi_{kr}(\boldsymbol{x}_{j}; \mathbf{w}_{k}) \mathcal{N}(y_{ij}; \boldsymbol{\beta}_{kr}^{T} \boldsymbol{x}_{j}, \sigma_{kr}^{2})}_{\text{RHLP}}$$

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

Two types of component memberships :

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

 \hookrightarrow 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
- EM-MixRHLP has complexity in $\mathcal{O}(I_{\mathsf{EM}}I_{\mathsf{IRLS}}KR^3nmp^3)$ (K-means type for piecewise regression is in $\mathcal{O}(I_{\mathsf{KM}}KRnm^2p^3) \hookrightarrow \mathsf{EM}$ -MixRHLP is computationally attractive for large values of m and moderate values of R.

EM-MixRHLP clustering of simulated data



Functional Linear Discriminant Analysis [8] Functional Mixture Discriminant Analysis [5]



tatistical data science and some unsupervised learning problems

Phonemes data

Phonemes data set used in Ferraty and Vieu [2003]³ 1000 log-periodograms (200 per cluster)



FIGURE - Original phoneme data and curves of the five classes : "ao", "aa", "yi", "dcl", "sh".

3. Data from http://www.math.univ-toulouse.fr/staph/npfda/

EM-like clustering results for Phonemes

Phonemes data set used in Ferraty and Vieu [2003]⁴ 1000 log-periodograms (200 per cluster)



4 Data from http://www.math.univ-toulouse fr/stanh/nnfda/

CHAMROUKHI Statistical data science and some unsupervised learning problems

EM-like clustering results for yeast cell cycle data

- Time course Gene expression data as in Yeung et al. [2001]⁵
- 384 genes expression levels over 17 time points.



FIGURE – EM-like clustering results with the bSRM model.

Rand index : 0.7914 which indicates that the partition is quite well defined.

5. http://faculty.washington.edu/kayee/model/

Outline

- 1 Statistics and Data Science
- 2 Clustering of multivariate data
- 3 Unsupervised learning for dimensionality reduction
- 4 Time series segmentaion
- **5** Clustering of functional data
- 6 Unsupervised Bayesian (non-)parametric learning

7 Model-Based Co-Clustering of Multivariate Functional Data

Bayesian spatial spline regression with mixed-effects

- Data : $((\boldsymbol{x}_1, \boldsymbol{y}_1), \dots, (\boldsymbol{x}_n, \boldsymbol{y}_n))$ a sample of n surfaces $\boldsymbol{y}_i = (y_{i1}, \dots, y_{im_i})^T$ and their spatial coordinates $\boldsymbol{x}_i = ((x_{i11}, x_{i12}), \dots, (x_{im_i1}, x_{im_i2}))^T$.
- Propose regression and regression mixtures, with three additional features :
- 1 Include random effects
- 2 Models for spatial functional data
- 3 A full Bayesian inference

Bayesian spatial spline regression with mixed-effects [Esann 2016, 13]

$$\boldsymbol{y}_i = \mathbf{S}_i(\boldsymbol{\beta} + \mathbf{b}_i) + \mathbf{e}_i, \ \mathbf{e}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_{m_i}), \ (i = 1, \dots, n)$$

- β : fixed-effects regression coefficients
- **b**_i : random subject-specific regression coefficients $\mathbf{b}_i \perp \mathbf{e}_i \sim \mathcal{N}(\mathbf{0}, \xi^2 \mathbf{I}_{m_i})$
- **S**_i is a spatial design matrix.

Bayesian mixture of spatial spline regressions

Data : A sample of n surfaces (y_1, \ldots, y_n) and their spatial covariates (S_1, \ldots, S_n) issued from K sub-populations

 Bayesian mixture of spatial spline regression models with mixed-effects (BMSSR) :

$$f(\boldsymbol{y}_i|\mathbf{S}_i;\boldsymbol{\Psi}) = \sum_{k=1}^{K} \pi_k \ \mathcal{N}\left(\boldsymbol{y}_i; \mathbf{S}_i(\boldsymbol{\beta}_k + \mathbf{b}_{ik}), \sigma_k^2 \mathbf{I}_{m_i}\right)$$

 \hookrightarrow Useful for density estimation and model-based clustering of heterogeneous surfaces

Hierarchical prior from for the BMSSR

$$\begin{array}{lll} \boldsymbol{\pi} & \sim & \mathcal{D}(\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_K) \\ \boldsymbol{\beta}_k & \sim & \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \\ \mathbf{b}_{ik} | \boldsymbol{\xi}_k^2 & \sim & \mathcal{N}(\mathbf{0}_d, \boldsymbol{\xi}_k^2 \mathbf{I}_d) \\ \boldsymbol{\xi}_k^2 & \sim & \mathcal{I}\mathcal{G}(a_0, b_0) \\ \boldsymbol{\sigma}_k^2 & \sim & \mathcal{I}\mathcal{G}(g_0, h_0). \end{array}$$

Bayesian inference of the BMSSR

For the BMSSR, the parameter Ψ is augmented by the unknown components labels $\mathbf{z} = (z_1, \dots, z_n)$

Bayesian inference of the BMSSR using Gibbs sampling

Sample from the analytic full conditional distributions :

$$egin{aligned} &Z_i|... \sim \mathcal{M}(1; au_{i1}, \ldots, au_{iK}) ext{ with } au_{ik}(1 \leq k \leq K) = \mathbb{P}(Z_i = k | oldsymbol{y}_i, oldsymbol{S}_i; oldsymbol{\Psi}) \ &\pi|... \sim \mathcal{D}\left(lpha_1 + n_1, \ldots, lpha_K + n_K
ight) \ &oldsymbol{eta}_k|... \sim \mathcal{N}(oldsymbol{
u}_0, oldsymbol{V}_0) \ &oldsymbol{b}_{ik}|... \sim \mathcal{N}(oldsymbol{
u}_1, oldsymbol{V}_1) \ &\sigma_k^2|... \sim \mathcal{IG}(g_1, h_1) \ &\xi_k^2|... \sim \mathcal{IG}\left(a_1, b_1
ight) \end{aligned}$$

 relabel the obtained posterior parameter samples if label switching by the K-means-like algorithm of [Celeux, 1999, Celeux et al., 2000].

Handwritten digit clustering using the BMSSR

- BMSSR applied on a subset of the ZIPcode data set (issued from MNIST)
- Each individual y_i contains $m_i = 256$ observations A subset of 1000 digits randomly chosen from the test set



 FIGURE – Cluster mean images obtained by the BMSSR model with 12 mixture components.

The best solution is selected in terms of the Adjusted Rand Index (ARI) values, which promotes a partition with K = 12 clusters (ARI : 0.5238).

Outline

- 1 Statistics and Data Science
- 2 Clustering of multivariate data
- 3 Unsupervised learning for dimensionality reduction
- 4 Time series segmentaion
- **5** Clustering of functional data
- 6 Unsupervised Bayesian (non-)parametric learning

7 Model-Based Co-Clustering of Multivariate Functional Data

Dirichlet Process Parsimonious Mixtures

- Bayesian parametric inference : [Bensmail, 1995, Bensmail and Celeux, 1996, Bensmail et al., 1997, Bensmail and Meulman, 2003]
- $\blacksquare \, \hookrightarrow \, \mathsf{Mixture}$ models for multivariate data in a fully Bayesian framework
- $\blacksquare \hookrightarrow$ Dirichlet Process and Parsimonious Mixtures [C5,6,8], [11]

Dirichlet Processes (DP)

 $\mathsf{DP}(\alpha, G_0)$ [Ferguson, 1973] is a distribution over distributions : $\tilde{\theta}_i | G \sim G$; $G | \alpha, G_0 \sim \mathsf{DP}(\alpha, G_0)$, i = 1, 2, ...

Pólya urn representation [Blackwell and MacQueen, 1973]

$$\tilde{\boldsymbol{\theta}}_i | \tilde{\boldsymbol{\theta}}_1, \dots \tilde{\boldsymbol{\theta}}_{i-1} \sim \frac{\alpha}{\alpha+i-1} G_0 + \sum_{k=1}^{K_{i-1}} \frac{n_k}{\alpha+i-1} \delta_{\boldsymbol{\theta}_k}$$

DP places its probability mass on an infinite mixture of Dirac deltas

$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\boldsymbol{\theta}_k} \quad \boldsymbol{\theta}_k | G_0 \sim G_0, \ k = 1, 2, ..., \text{ with } \sum_{k=1}^{\infty} \pi_k = 1$$

 \hookrightarrow The generated parameters $ilde{m{ heta}}_i$ for a DP process exhibit a clustering property

DPM : Generative model

Chinese Restaurant Process mixtures (Pitman, 2002; Samuel and Blei, 2012)

- Latent variables (z_1, \ldots, z_n)
- Predictive distribution :

$$p(z_i = k | z_1, ..., z_{i-1}; \alpha) = \frac{\alpha}{\alpha + i - 1} \delta(z_i, K_{i-1} + 1) + \sum_{k=1}^{K_{i-1}} \frac{n_k}{\alpha + i - 1} \delta(z_i, k) \cdot$$



Generative model :

 $\begin{array}{rcl} z_i | \alpha & \sim & \mathsf{CRP}(\mathbf{z}_{\backslash i}; \alpha) \\ \boldsymbol{\theta}_{z_i} | G_0 & \sim & G_0 \\ \mathbf{x}_i | \boldsymbol{\theta}_{z_i} & \sim & f(.| \boldsymbol{\theta}_{z_i}) \end{array}$

FAICEL CHAMROUKHI Statistical data science and some unsupervised learning

Implemented parsimonious models

Decomposition	Model-Type	Prior	Applied to
λΙ	Spherical	IG	λ
$\lambda_k \mathbf{I}$	Spherical	\mathcal{IG}	λ_k
$\lambda \mathbf{A}$	Diagonal	\mathcal{IG}	each diagonal element of $\lambda {f A}$
$\lambda_k \mathbf{A}$	Diagonal	\mathcal{IG}	each diagonal element of $\lambda_k \mathbf{A}$
$\lambda \mathbf{D} \mathbf{A} \mathbf{D}^T$	General	IW	$\mathbf{\Sigma} = \lambda \mathbf{D} \mathbf{A} \mathbf{D}^T$
$\lambda_k \mathbf{D} \mathbf{A} \mathbf{D}^T$	General	$\mathcal{I}\mathcal{G}$ and $\mathcal{I}\mathcal{W}$	λ_k and $oldsymbol{\Sigma} = \mathbf{D} \mathbf{A} \mathbf{D}^T$
$\lambda \mathbf{D} \mathbf{A}_k \mathbf{D}^T *$	General	\mathcal{IG}	each diagonal element of $\lambda \mathbf{A}_k$
$\lambda_k \mathbf{D} \mathbf{A}_k \mathbf{D}^T *$	General	\mathcal{IG}	each diagonal element of $\lambda_k \mathbf{A}_k$
$\lambda \mathbf{D}_k \mathbf{A} \mathbf{D}_k^T$	General	\mathcal{IG}	each diagonal element of $\lambda {f A}$
$\lambda_k \mathbf{D}_k \mathbf{A} \mathbf{D}_k^T$	General	\mathcal{IG}	each diagonal element of $\lambda_k \mathbf{A}$
$\lambda \mathbf{D}_k \mathbf{A}_k \mathbf{D}_k^T *$	General	$\mathcal{I}\mathcal{G}$ and $\mathcal{I}\mathcal{W}$	λ and $\mathbf{\Sigma}_k = \mathbf{D}_k \mathbf{A}_k \mathbf{D}_k^T$
$\lambda_k \mathbf{D}_k \mathbf{A}_k \mathbf{D}_k^T$	General	\mathcal{IW}	$\mathbf{\Sigma}_k = \lambda_k \mathbf{D}_k \mathbf{A}_k \mathbf{D}_k^T$

Bayesian inference using Gibbs sampling

- Posterior distribution for the component labels : $p(z_i = k | \mathbf{z}_{-i}, \mathbf{X}, \mathbf{\Theta}, \alpha) \propto p(\mathbf{x}_i | z_i; \mathbf{\Theta}) p(z_i | \mathbf{z}_{-i}; \alpha)$ with $p(z_i | \mathbf{z}_{-i}; \alpha)$ the CRP prior
- Posterior distribution for the component parameters : $p(\boldsymbol{\theta}_k | \mathbf{z}, \mathbf{X}, \boldsymbol{\Theta}_{-k}, \alpha; H) \propto \prod_{i \mid z_i = k} p(\mathbf{x}_i | z_i = k; \boldsymbol{\theta}_k) p(\boldsymbol{\theta}_k; H)$ with $p(\boldsymbol{\theta}_k; H)$: Prior distribution over $\boldsymbol{\theta}_k$

Bayesian model comparison by using Bayes Factors $BF_{12} = \frac{p(\mathbf{X}|M_1)p(M_1)}{p(\mathbf{X}|M_2)p(M_2)} \approx \frac{p(\mathbf{X}|M_1)}{p(\mathbf{X}|M_2)} \text{ with the Laplace-Metropolis approximation}$ $p(\mathbf{X}|M_m) = \int p(\mathbf{X}|\boldsymbol{\theta}_m, M_m)p(\boldsymbol{\theta}_m|M_m) d\boldsymbol{\theta}_m \approx (2\pi)^{\frac{\nu_m}{2}} |\hat{\mathbf{H}}|^{\frac{1}{2}} p(\mathbf{X}|\hat{\boldsymbol{\theta}}_m, M_m)p(\hat{\boldsymbol{\theta}}_m|M_m)$

Clustering of benchmarks

Diabetes data set, Geyser data set, Crabs data set



Statistical data science and some unsupervised learning problem
Humpback whale song decomposition

- Real fully unsupervised problem
- Data : 8.6 minutes of a Humpback whale song recording (with MFCC)



Objectives

- Discovering "call units", which can be considered as a whale "alphabet"
- Find a partition of the whale song into clusters (segments), and automatically infer the unknown number of clusters from the data.



Sound demo of Unit 5 DPPM λI : (sec. 0) (sec. 12)



Sound demo of Unit 8 DPPM λI : (sec. 8) (sec. 10)



Sound demo of Unit 4 DPPM $\lambda_k \mathbf{A}$: (sec. 1) (sec. 7)



Sound demo of Unit 8 DPPM $\lambda_k \mathbf{A}$: (sec. 6) (sec. 12)

Outline

- 1 Statistics and Data Science
- 2 Clustering of multivariate data
- 3 Unsupervised learning for dimensionality reduction
- 4 Time series segmentaion
- 5 Clustering of functional data
- 6 Unsupervised Bayesian (non-)parametric learning

7 Model-Based Co-Clustering of Multivariate Functional Data

- Model-based co-clustering
- Co-clustering of multivariate functional data

Faicel Chamroukhi

istical data science and some unsupervised learning problems

Outline

- 1 Statistics and Data Science
- 2 Clustering of multivariate data
- 3 Unsupervised learning for dimensionality reduction
- 4 Time series segmentaion
- 5 Clustering of functional data
- 6 Unsupervised Bayesian (non-)parametric learning

7 Model-Based Co-Clustering of Multivariate Functional Data

- Model-based co-clustering
- Co-clustering of multivariate functional data

Faicel Chamroukhi

istical data science and some unsupervised learning problems

High-dimensional functional data clustering

- Multivariate functional data are increasingly present
- e.g : Data continuously recorded for different subjects from multiple subject' sensors

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



FIGURE – An example with d = 30 and n = 20 daily observations [Ben Slimen et al., 2016].

High-dimensional functional data clustering

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)

Clustering VS Co-Clustering

Simultaneous clustering of lines/indiv. (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

The Latent Block Model [Govaert and Nadif, 2013]

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

Hypotheses

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

Latent block model for co-clustering

The Latent Block Model [Govaert and Nadif, 2013]

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

Hypotheses

- The latent variables Z and W are independent : $\mathbb{P}(Z, W) = \mathbb{P}(Z)\mathbb{P}(W)$ and iid : $\mathbb{P}(Z) = \prod_i \mathbb{P}(z_i)$ with $z_i \sim \text{Multinomial}(\pi_1, \dots, \pi_K)$ where $\pi_k = \mathbb{P}(z_k = k)$ $\mathbb{P}(W) = \prod_j \mathbb{P}(w_j)$ with $w_j \sim \text{Multinomial}(\rho_1, \dots, \rho_M)$ where $\rho_\ell = \mathbb{P}(w_j = \ell)$
- Conditional independence : $x_{ij}|(z_i, w_j) \perp x_{i\prime j\prime}|(z_{i\prime}, w_{j\prime})$
- \hookrightarrow binary data : binary [Govaert and Nadif, 2003, 2008, Keribin et al., 2012],
- \hookrightarrow categorical data : multinomial [Keribin et al., 2014]
- \hookrightarrow contingency table : Poisson [Govaert and Nadif, 2003, 2006, 2008]
- \hookrightarrow continuous data : Gaussian [Lomet, 2012, Govaert and Nadif, 2013]
- \hookrightarrow functional data : functional PCA + Gaussian, see further [Ben Slimen et al., 2016]

Inference for the latent block model

Inference of the latent block model

- variational block EM (VBEM) for maximum likelihood estimation and fuzzy co-clustering [Govaert and Nadif, 2006, 2008].
- block classification EM (CEM) algorithm for maximum classification likelihood and hard co-clustering [Govaert and Nadif, 2003, 2006, 2008]
- Bayesian inference [Keribin et al., 2012, 2014] : Bayesian latent block mixtures for binary data and categorical data & a variational Bayesian inference and Gibbs sampling.
- Number of blocks estimation : ICL criterion [Lomet, 2012, Keribin et al., 2014]

Functional data modeling : "classical" approach

[Ramsay and Silverman, 2005] and many others

- Step 1 : (x, y) decomposed into a finite basis of function (B-spline...) : $Y_i(t) \approx \sum_{r=1}^d c_{ir} \phi_r(x_i(t))$ with c estimated by OLS
- Step 2 : functional principal components analysis (PCA) which is performed as a usual PCA of the basis expansion coefficients c using a metric defined by the inner products between the basis functions
- Step 3 : set a probability distribution on \mathbf{c} , typically Gaussian

It defines a distribution on ${f c}$ instead of y_{\cdots}

Functional data modeling : regression RHLP

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

 $\label{eq:response} \hookrightarrow \mathsf{Regression} \ \mathsf{with} \ \mathsf{Hidden} \ \mathsf{Logistic} \ \mathsf{Process} \ (\mathsf{RHLP}) \\ \hookrightarrow \mathsf{See} \ \mathsf{formula} \ \mathsf{later}$

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

RHLP for modeling different types of functions



Faicel Chamroukhi

Statistical data science and some unsupervised learning problems

Multivariate functional data co-clustering

[Chamroukhi and Biernacki, 2017]

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



Embedding RHLP in co-clustering

[Chamroukhi and Biernacki, 2017]

Functional Latent Block Model for Co-clustering :

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

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

Parameter estimation : EM not feasible

- \hookrightarrow Requires the calculation of the posterior joint distribution $\mathbb{P}(z_{ik}w_{j\ell}=1|m{y}_{ij},m{x}_{ij})$
- \hookrightarrow does not factorize due to the conditional dependence on the observed curves of the row and the column labels
- \Rightarrow Variational block EM algorithm : [Govaert and Nadif, 2008, 2013]
- $\,\hookrightarrow\,$ We adopt this variational approximation in our context

Variational block EM algorithm

variational approximation

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

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

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

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 :

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

where :

$$\tilde{z}_{ik} = \mathbb{P}(z_{ik} = 1 | \boldsymbol{y}_{ij}, \boldsymbol{x}_{ij}),$$

$$\tilde{w}_{j\ell} = \mathbb{P}(w_{j\ell} = 1 | \boldsymbol{y}_{ij}, \boldsymbol{x}_{ij}),$$

$$\tilde{h}_{tr} = \mathbb{P}(h_{tr} = 1 | z_i, w_j, y_{ij}(t), x_{ij}(t))$$

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

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

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

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

$$\begin{array}{l} 1 \quad \pi_k^{(q+1)} = \frac{\sum_i \tilde{z}_{ik}^{(q+1)}}{n} \\ 2 \quad \rho_\ell^{(q+1)} = \frac{\sum_j \tilde{w}_{j\ell}^{(q+1)}}{d} \end{array}$$

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

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

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

$$\begin{array}{l} 1 \quad \pi_k^{(q+1)} = \frac{\sum_i \tilde{z}_{ik}^{(q+1)}}{n} \\ 2 \quad \rho_\ell^{(q+1)} = \frac{\sum_j \tilde{w}_{j\ell}^{(q+1)}}{d} \end{array} \end{array}$$

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

$$\begin{array}{l} \textbf{3} \hspace{0.5cm} \boldsymbol{\xi}_{k\ell}^{(new)} = \boldsymbol{\xi}_{k\ell}^{(old)} - \left[\frac{\partial^{2}F(\boldsymbol{\xi}_{k\ell})}{\partial \boldsymbol{\xi}_{k\ell}\partial \boldsymbol{\xi}_{k\ell}^{T}}\right]_{\boldsymbol{\xi}_{k\ell} = \boldsymbol{\xi}_{k\ell}^{(old)}}^{-1} \frac{\partial F(\boldsymbol{\xi}_{k\ell})}{\partial \boldsymbol{\xi}_{k\ell}} \Big|_{\boldsymbol{\xi}_{k\ell} = \boldsymbol{\xi}_{k\ell}^{(old)}} \hspace{0.5cm} \text{which is the IRLS} \\ \text{maximisation of } F(\boldsymbol{\xi}_{k\ell}) = \sum_{i,j,t} \tilde{z}_{ik}^{(q)} \tilde{w}_{j\ell}^{(q)} \tilde{h}_{tr}^{(q)} \log \alpha_{k\ell r}(t; \boldsymbol{\xi}_{k\ell}) \hspace{0.5cm} \text{w.r.t} \hspace{0.5cm} \boldsymbol{\xi}_{k\ell}. \\ \text{The regression parameters updates consist in analytic WLS problems :} \\ \textbf{4} \hspace{0.5cm} \boldsymbol{\beta}_{k\ell r}^{(q+1)} = \left[\sum_{i,j} \tilde{z}_{ik}^{(q)} \tilde{w}_{j\ell}^{(q)} \mathbf{X}_{ij}^{T} \boldsymbol{\Lambda}_{ijkr}^{(q)} \mathbf{X}_{ij}\right]^{-1} \sum_{i,j} \tilde{z}_{i,j}^{(q)} \tilde{w}_{j\ell}^{(q)} \mathbf{X}_{ijkr}^{T} \boldsymbol{\Lambda}_{ijkr}^{(q)} \\ \textbf{5} \hspace{0.5cm} \sigma_{k\ell r}^{2(q+1)} = \frac{\sum_{i,j} \tilde{z}_{ik}^{(q)} \tilde{w}_{j\ell}^{(q)} \| \sqrt{\boldsymbol{\Lambda}_{ijkr}^{(q)}(\boldsymbol{y}_{ij} - \boldsymbol{X}_{ij} \boldsymbol{\beta}_{kr}^{(q+1)}) \|^{2}}}{\sum_{i,j} \tilde{z}_{ik}^{(q)} \tilde{w}_{j\ell}^{(q)} \operatorname{trace}(\boldsymbol{\Lambda}_{ijkr}^{(q)})} \\ \textbf{6} \hspace{0.5cm} \boldsymbol{w}^{(q)} \text{the curve, } \boldsymbol{\Lambda}_{ijkr}^{(q)} \text{ is the diagonal matrix whose diagonal elements are the posterior segment memberships } \{ \tilde{h}_{ijtr}^{(q)}; t = 1, \ldots, T_{ij} \}. \end{array} \right.$$

 \hookrightarrow It is also possible to use the Classification EM (CEM) approximation of EM [Celeux and Govaert, 1992].

Parameter estimation by an SEM algorithm : SEM-FLBM

- → The SEM algorithm [Celeux and Diebolt, 1985] allows to overcome some drawbacks of the variational-EM algorithm, including its sensitivity to starting values; SEM does not use an approximation.
- Eg. SEM for latent block models for categorical data [Keribin et al., 2012, 2014]
- The formulas of VEM-FLBM and SEM-FLBM are essentially the same, except that we incorporate a stochastic step consisting of sampling binary indicator variables z_{ik} , $w_{j\ell}$ and h_{tr} according to \tilde{z}_{ik} , $\tilde{w}_{j\ell}$ and \tilde{h}_{tr} .

Source codes are/will be made available on github

$\mathsf{Matlab}/\mathsf{R}/\mathsf{Python}$

https://github.com/fchamroukhi

E -> C a GitHub, Inc. [US] https://github.com/fchamrou	khi		☆ 0
and the second s	Overview Repositories Stars: P Folio Popular repositories RHLP_Matlab User-finision and flatble algorithm for time series aggregation model with a Holden Logistic Process (RHLP). Logistic Process (RHLP).	Vers 4 Following (1)	
	● MATLAB ★1	MATLAB	
fchamroukhi Block or report user Professor of Statistics and Data Sciences	HMMR Hidden Markov Model Regression (HMMR) for Times Series Segmentation MATLAB	PWR Placowise regression (PWR) for the optimal segmentation of time series with regime changes MATLAB	
L Caen University			
⑦ Ceen, France ③ https://chamroukhi.com	MixFHMM Clustering and segmentation of time series by mixture of gaussian Holden Markov Models (MxXHMMs) and the EM algorithm • MATLAB	MixFHMMR Clustering and segmentation of time series with regime changes by micture of Hidden Markov Model Regressions (MixFHMMR) and the EM algorithm MATLAB	
	39 contributions in the last year		
	Die Jan Pels Mer Apr Mey Jan Man	Jul Aug Ses Oct Nov Dec	

Data science, Big-Data, Al

The way of the future !

Eg. In France : Interdisciplinary Institutes of Artificial Intelligence (3IA) :



6/11/2018

Interdisciplinary Institutes of Artificial Intelligence (3IA): the four selected projects



The results of the 3IA (Interdisciplinary Institute of Artificial Intelligence) call for expressions of interest were made public on 6 November 2018 by Frédérique Vidal, French Minister of Higher Education, Research and Innovation, and Mounir Mahjoubi, French Secretary of State for Digital Affairs. The projects of the Grenoble (MIAI@Grenoble-Alpes), Nice (3IA Côte d'Azur), Paris (PRAIRIE) and Toulouse (ANITI) sites have been selected. Inria participates in three of the four successful projects.

At the heart of the national AI strategy

Each of the Inria research centres took part in formulating 3IA Institute projects within the framework of the considerable mobilisation of the regional academic and industry ecosystems: they are part of the dynamic stimulated by the national plan on artificial intelligence announced by the French president following the report by Cédric Villani.

Data science models/algorithms

New problems (big data, etc) but ... classical methods?



 Regression, decision trees, and cluster analysis continue to form a triad of core algorithms for most data miners. This has been consistent since the first Data Miner Survey in 2007.



References I

Jeffrey D. Banfield and Adrian E. Raftery. Model-based Gaussian and non-Gaussian clustering. Biometrics, 49(3):803-821, 1993.

- 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.
- H. Bensmail and Jacqueline J. Meulman. Model-based Clustering with Noise : Bayesian Inference and Estimation. Journal of Classification, 20(1):049–076, 2003.
- H. Bensmail, G. Celeux, A. E. Raftery, and C. P. Robert. Inference in model-based cluster analysis. Statistics and Computing, 7 (1):1–10, 1997.

Halima Bensmail. Modèles de régularisation en discrimination et classification bayésienne. PhD thesis, Université Paris 6, 1995.

- Halima Bensmail and Gilles Celeux. Regularized gaussian discriminant analysis through eigenvalue decomposition. Journal of the American statistical Association, 91(436) :1743–1748, 1996.
- D. Blackwell and J. MacQueen. Ferguson Distributions Via Polya Urn Schemes. The Annals of Statistics, 1:353-355, 1973.
- C. Bouveyron and J. Jacques. Model-based clustering of time series in group-specific functional subspaces. Adv. Data Analysis and Classification, 5(4):281–300, 2011.
- C. Bouveyron, L. Bozzi, J. Jacques, and F.-X. Jollois. The functional latent block model for the co-clustering of electricity consumption curves. *Journal of the Royal Statistical Society, Series C*, 2018.
- G. Celeux. Bayesian inference for mixture : the label switching problem. Technical report, INRIA Rhone-Alpes, 1999.
- G. Celeux and J. Diebolt. The SEM algorithm a probabilistic teacher algorithm derived from the EM algorithm for the mixture problem. Computational Statistics Quarterly, 2(1):73–82, 1985.
- G. Celeux and G. Govaert. A classification EM algorithm for clustering and two stochastic versions. Computational Statistics and Data Analysis, 14:315–332, 1992.
- G. Celeux, M. Hurn, and C. P. Robert. Computational and inferential difficulties with mixture posterior distributions. Journal of the American Statistical Association, 95(451):957–970, 2000.

References II

- 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, 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.com/papers/chamroukhi_neucomp_2010.pdf.
- Faicel Chamroukhi and Hien D. Nguyen. Model-based clustering and classification of functional data. Wiley Interdisciplinary Reviews : Data Mining and Knowledge Discovery, To appear, 2018. URL https://chamroukhi.com/papers/MBCC-FDA.pdf. arXiv :1803.00276v2.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. Journal of The Royal Statistical Society, B, 39(1):1–38, 1977.
- Thomas S. Ferguson. A Bayesian Analysis of Some Nonparametric Problems. The Annals of Statistics, 1(2):209-230, 1973.
- F. Ferraty and P. Vieu. Curves discrimination : a nonparametric functional approach. Computational Statistics & Data Analysis, 44(1-2):161–173, 2003.
- 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. Hugueney, Y. Lechevallier, and F. Rossi. Exploratory analysis of functional data via clustering and optimal segmentation. *Neurocomputing*, 73(7-9) :1125–1141, March 2010.
- H. Hotelling. Analysis of a complex of statistical variables into principal components. Journal of Educational Psychology, 24 : 417–441, 1933.

References III

- Julien Jacques and Cristian Preda. Functional data clustering : A survey. Adv. Data Anal. Classif., 8(3) :231–255, September 2014. ISSN 1862-5347. doi: 10.1007/s11634-013-0158-y. URL http://dx.doi.org/10.1007/s11634-013-0158-y.
- G. M. James and T. J. Hastie. Functional linear discriminant analysis for irregularly sampled curves. Journal of the Royal Statistical Society Series B, 63:533–550, 2001.
- G. M. James and C. Sugar. Clustering for sparsely sampled functional data. Journal of the American Statistical Association, 98 (462), 2003.
- 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.
- J. B. MacQueen. Some methods for classification and analysis of multivariate observations. In Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability, pages 281–297, 1967.
- G. J. McLachlan and D. Peel. Finite mixture models. New York : Wiley, 2000.
- K. Pearson. On lines and planes of closest fit to systems of points in space. Philosophical Magazine, 2(6):559-572, 1901.
- J. O. Ramsay and B. W. Silverman. Functional Data Analysis. Springer Series in Statistics. Springer, June 2005.
- S. Roweis. EM algorithms for PCA and SPCA. In M. I. Jordan, M. J. Kearns, and S. A. Solla, editors, Proceedings of the 11th Conference on Advances in Neural Information Processing Systems (NIPS), volume 10. MIT Press, 1998.
- A. Samé, F. Chamroukhi, G. Govaert, and P. Aknin. Model-based clustering and segmentation of time series with changes in regime. Advances in Data Analysis and Classification, pages 1–21, 2011. ISSN 1862-5347.
- C. Spearman. General intelligence, objectively determined and measured. American Journal of psychology, 15 :201-293, 1904.

References IV

- L. L. Thurstone. Multiple Factor Analysis. University of Chicago Press, 1947.
- M. E. Tipping and C. Bishop. Probabilistic principal component analysis. Technical Report NCRG/97/010, Neural Comuting Research Group, Aston University, 1997.
- M. E. Tipping and C. Bishop. Probabilistic principal component analysis. Journal of the Royal Statistical Society, Series B, 61: 611–622, 1999.
- Ka Yee Yeung, Chris Fraley, A. Murua, Adrian E. Raftery, and Walter L. Ruzzo. Model-based clustering and data transformations for gene expression data. *Bioinformatics*, 17(10):977–987, 2001.

Thank you for your attention !