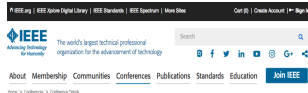


# Learning from Heterogeneous & Non-Stationary Time-Series Data

FAICEL CHAMROUKHI

<https://chamroukhi.com/>

- 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



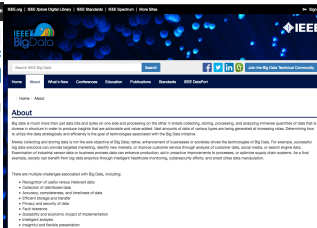
None > Conferences > Conference Details

## 2017 IEEE International Conference on Data Science and Advanced Analytics (DSAA)

IEEE sponsors:

- \* IEEE Computational Intelligence Society

DSAA is a premier forum that brings together researchers, industry practitioners, as well as potential users of data science, big data and advanced analytics, to promote collaborations and exchange of ideas and practices, discuss new opportunities, and investigate the best actionable analytics framework for wide range of applications. DSAA solicits both experimental and theoretical works on data science and advanced analytics along with their application to real life situations. Topics include but not limited to data analytics, machine learning, data mining, knowledge discovery, storage, search, privacy, security, complexity, efficiency, scalability and visualization.



## About

Big data is much more than just data sets and tools for on-line data processing on the other. It is a collection, storing, processing, and analyzing massive quantities of data that is diverse in structure in order to produce insights that are actionable and value added. Real amounts of data of various forms are being generated at increasing rates. Storing them to allow the data integrally and efficiently is the goal of the technologies associated with the Big Data initiative.

Many collecting and storing data is not the sole objective of Big Data; rather, enhancement of businesses or societies drive the technologies of Big Data. For example, successful big data analytics can provide targeted marketing, identify new markets, or improve customer service through analysis of customer data, social media, or search engine data. Identification of fraudulent sensor data in business process data can enhance production, and predictive improvements to processes, or predictive supply chain systems. As a true example, today's car benefit from big data analytics through intelligent features monitoring, cybersecurity efforts, and smart cities data management.

There are multiple challenges associated with Big Data, including:

- Heterogeneity of multi-source relevant data
- Collection of distributed data
- Accesses, completeness, and timeliness of data
- Efficient storage and transfer
- Privacy and security of data
- Fault tolerance
- Scalability and economic impact of implementation
- Intelligent analysis
- Insightful and flexible presentation

## ICLR 2017



## 5th International Conference on Learning Representations

### Overview

The performance of machine learning methods is heavily dependent on the choice of data representation (or features) on which they are applied. The rapidly developing field of representation learning is concerned with questions surrounding how we can best learn meaningful and useful representations of data. In this tutorial, we will take a broad view of the field and include topics such as deep learning and feature learning, matrix learning, convolutional modeling, adversarial models, reinforcement learning, and models regarding large-scale learning and non-linear optimization. The range of examples to which these techniques apply is also very broad, from vision to speech recognition, text understanding, gaming, music, etc.

A non-exhaustive list of relevant topics:

- Unsupervised, semi-supervised, and supervised representation learning
- Representation learning for planning and reinforcement learning
- Matrix learning and kernel learning
- Sparse coding and dimensionality reduction
- Hierarchical models
- Optimization for representation learning
- Learning representations of objects in data
- Representation issues, visualization, software platforms, hardware
- Applications in vision, audio, speech, natural language processing, robotics, neuroscience, or any other field



The University of Michigan (U-M) plans to invest \$100 million over the next five years in a new Data Science Initiative (DSI) that will enhance opportunities for students and faculty researchers across the University to tap into the enormous potential of big data.

The DSI plans to:

- hire 30 new faculty over the next four years and engage existing faculty across campus;
- support interdisciplinary research and foster new technological approaches to big data;
- provide new educational opportunities for students pursuing careers in data science;
- expand its role in research, computing, teaching, and
- strengthen data management, storage, analysis, and training resources.

The Data Science Initiative brings together the newly created Michigan Institute for Data Science (MIDS), Consulting for Statistics, Computing and Analytics Research (CSCAR) and Advanced Research Computing - Technology Services (ARC-TS) to provide a coordinated and comprehensive home for data science as part of Advanced Research Computing (ARC) at the University.



Harvard Business Review



## Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil  
Presented with permission from DSI

[illegible]

ABOUT • PEOPLE
CONTACT

# université **PARIS-SACLAY**

# Paris-Saclay Center for Data Science

---

Tweets by @SaclayCDS

## PARIS-SACLAY Center for Data Science (CDS)

Phase I : Lidex Paris-Saclay (2014 – 2016)

Phase II : IRIS Initiatives de Recherche Stratégiques (2016 – 2019)

**Extracting knowledge from data.**

The project consists of developing methods and tools so as to be capable of analysing gigantic amounts of data and extracting useful information from them for physics, biology, medicine, chemistry, the environment and the human sciences.

This project is multidisciplinary: it requires research on analytical methodologies (statistics, processes of machine learning, extracting knowledge, viewing data), as well as on software design.

More than 250 permanent researchers in 35 laboratories participate in the CDS supporting our datascience projects and events.

### Newes

- Associate/full professor position in the area of Computer Vision 2017/03/20
- Appel à projets, Paris-Saclay Center for Data Science 2017/03/14
- Wikipédia pour la science 2017 2017/03/14
- Permanent position of Professor in Signal Processing at CentraleSupélec @ Université Paris-Saclay 2017/03/13
- Permanent position (Associated Professor) in Machine Learning @ TELECOM-ParisTech 2017/03/13
- Appel à projets émergents 2017 du département STIC 2017/03/01
- One day Workshop and Hackathon on spatio-temporal time series

- 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

CONTRIBUTED ARTICLES

## Data Science and Prediction

By Vasant Dhar

Communications of the ACM, Vol. 56 No. 12, Pages 64-73

10.1145/2500499

Comments (2)

VIEW AS:     SHARE:     



Use of the term “data science” is increasingly common, as is “big data.” But what does it mean? Is there something unique about it? What skills do “data scientists” need to be productive in a world deluged by data? What are the implications for scientific inquiry? Here, I address these questions from the perspective of predictive modeling.

[Back to Top](#)

### Key Insights

- Data science is the study of the generalizable extraction of knowledge from data.
- A common epistemic requirement in assessing whether new knowledge is ascertainable for decision making in its predictive power, not just its ability to explain the past.
- A data scientist requires an integrated skill set spanning mathematics, machine learning, artificial intelligence, statistics, databases, and optimization, along with a deep understanding of the craft of problem formulation to engineer effective solutions.

ASA

Amstat News

ASA Community

The World of Sta

# AMSTATNEWS

The Membership Magazine of the American Statistical Association

HOME ABOUT EDITORIAL CALENDAR PDF ARCHIVES ADVERTISE STATISTICIANS IN HISTORY

Home » Featured

## ASA Statement on the Role of Statistics in Data Science

1 OCTOBER 2015 6,956 VIEWS 13 COMMENTS

### Statement Contributors

David van Dyk, Imperial College (chair)

Montse Fuentes, NCSU

Michael I. Jordan, UC Berkeley

Michael Newton, University of Wisconsin

Bonnie K. Ray, Pegged Software

Duncan Temple Lang, UC Davis

Hadley Wickham, RStudio

The rise of data science, including Big Data and data analytics, has recently attracted enormous attention in the popular press for its spectacular contributions in a wide range of scholarly disciplines and commercial endeavors. These successes are largely the fruit of the innovative and entrepreneurial spirit that characterize this burgeoning field. Nonetheless, its interdisciplinary nature means that a substantial collaborative effort is needed for it to realize its full potential for productivity and innovation. While there is not yet a consensus on what precisely constitutes data science, three professional communities, all within computer science and/or statistics, are emerging as foundational to data science: (i)

Database Management enables transformation, conglomeration, and organization of data resources, (ii) Statistics and Machine Learning convert data into knowledge, and (iii) Distributed and Parallel Systems provide the computational infrastructure to carry out data analysis.

- For a review, see the report of D. Donoho (2015) : “50 years of Data Science”






NOUS CONNAÎTRE **VIE SCIENTIFIQUE** ENSEIGNEMENT DES SCIENCES DIFFUSION DES CONNAISSANCES COLLABORATIONS INTERNATIONALES EXPERTISE ET CONSEIL

14<sup>th</sup> 2014

## La datamasse : directions et enjeux pour les données massives

Publié dans Colloques, conférences et débats



Conférence-débat de l'Académie des sciences

Nous vivons dans une "société de l'information" dont les avancées scientifiques et techniques rapides, associées au développement d'usages nouveaux, conduisent à produire des quantités toujours plus gigantesques de données numériques. Cette situation d'abondance ouvre des perspectives nouvelles tant dans les sciences exactes que dans les sciences humaines. L'utilisation de cette "datamasse" (Big Data en anglais) pose des défis considérables : Comment stocker de telles quantités de données, les manipuler, les analyser, les trier... les valoriser ? Comment concilier leur omniprésence et le respect de la vie privée ? Comment faire qu'elles bénéficient à tous ? Ce sont quelques-uns de ces aspects qui seront mis en avant dans cette rencontre, afin d'en mieux comprendre les possibilités et les limitations, pour en mieux maîtriser les développements.

### Introduction

Serge Abiteboul, directeur de recherche Inria, École normale supérieure de Cachan, membre de l'Académie des sciences et Patrick Flandin, directeur de recherche CNRS, École normale supérieure de Lyon, 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 Malat, professeur à l'École normale supérieure, Paris



### La découverte du cerveau grâce à l'exploration de données massives

Anastasia Ailamaki, professeure à l'École polytechnique fédérale de Lausanne



### 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



### Discussion générale et conclusion



Vidéos réalisées par la cellule Webcast CC-IN2P3 du CNRS  

- There is not yet a consensus on what precisely constitutes Data Science, but
- Data Science can be seen (defined ?) as <sup>a</sup> :
  - ▶ 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

↔ 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 framework) and to establish generalizable 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

# Outline

- 1 Introduction
- 2 Latent data models for temporal data segmentation
- 3 Mixture models for functional data analysis
- 4 Mixture-of-Experts for fitting complex non-normal distributions

# Unsupervised Learning

## 1 Introduction

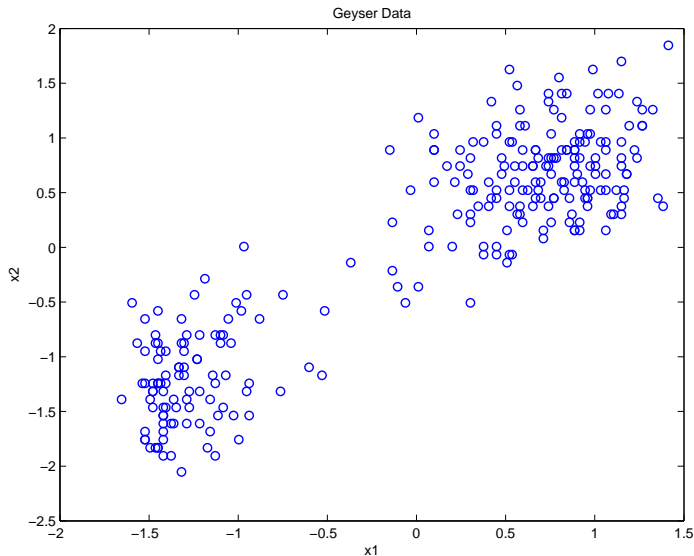
- Statistics and Data Science
- Unsupervised Learning

## 2 Latent data models for temporal data segmentation

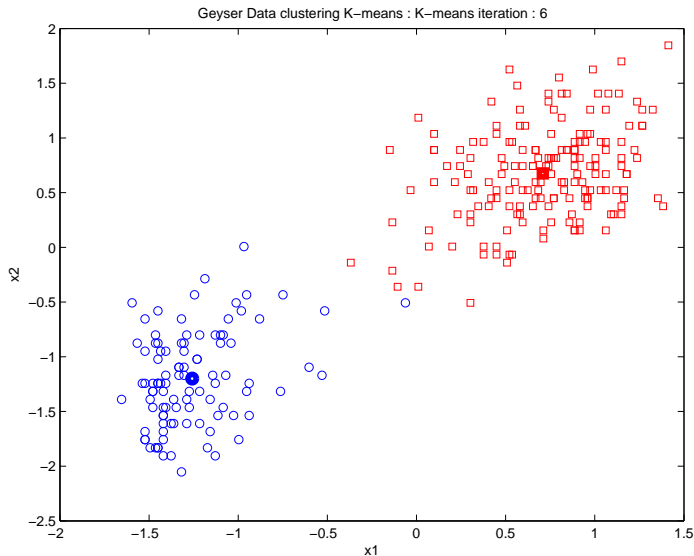
## 3 Mixture models for functional data analysis

## 4 Mixture-of-Experts for fitting complex non-normal distributions

# Clustering of multivariate data



# Clustering of multivariate data



# K-means

- a straightforward and widely used clustering algorithm, is one of the most important algorithms in unsupervised learning.
- Observed data  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$  in  $\mathbb{R}^d$  with unknown cluster labels  $\mathbf{z} = (z_1, \dots, z_n)$  ( $z_i \in 1, \dots, K$ )
- Each of the  $K$  clusters is represented by its mean (cluster centroid)  $\mu_k$  in  $\mathbb{R}^d$ .

## K-means [MacQueen, 1967]

$$(\hat{\mu}_1, \dots, \hat{\mu}_K, \hat{\mathbf{z}}) \in \arg \min_{\mu_1, \dots, \mu_K, \mathbf{z}} \mathcal{J}(\mu_1, \dots, \mu_K, \mathbf{z})$$

$$\text{objective function : } \mathcal{J}(\mu_1, \dots, \mu_K, \mathbf{z}) = \sum_{k=1}^K \sum_{i=1}^n \|\mathbf{x}_i - \mu_{z_i}\|^2$$



# K-means

- a straightforward and widely used clustering algorithm, is one of the most important algorithms in unsupervised learning.
- Observed data  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$  in  $\mathbb{R}^d$  with unknown cluster labels  $\mathbf{z} = (z_1, \dots, z_n)$  ( $z_i \in 1, \dots, K$ )
- Each of the  $K$  clusters is represented by its mean (cluster centroid)  $\mu_k$  in  $\mathbb{R}^d$ .

## K-means [MacQueen, 1967]

$$(\hat{\mu}_1, \dots, \hat{\mu}_K, \hat{\mathbf{z}}) \in \arg \min_{\mu_1, \dots, \mu_K, \mathbf{z}} \mathcal{J}(\mu_1, \dots, \mu_K, \mathbf{z})$$

objective function :  $\mathcal{J}(\mu_1, \dots, \mu_K, \mathbf{z}) = \sum_{k=1}^K \sum_{i=1}^n \|\mathbf{x}_i - \mu_{z_i}\|^2$

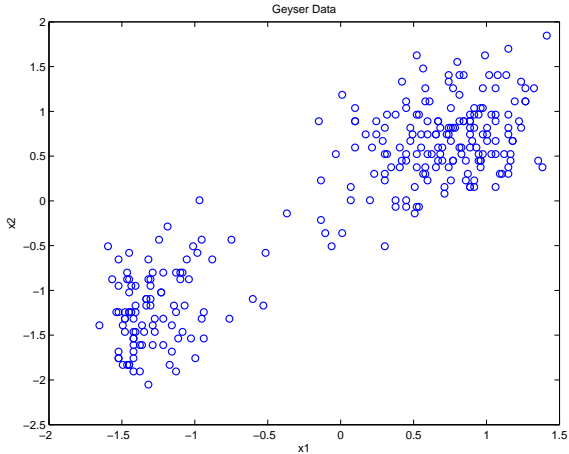
- Initialization :  $(\mu_1^{(0)}, \dots, \mu_K^{(0)})$  (eg, randomly chosen data points)

**1 Assignment step** :  $z_i^{(t)} = \arg \min_{z \in \mathcal{Z}} \|\mathbf{x}_i - \mu_z\|^2$

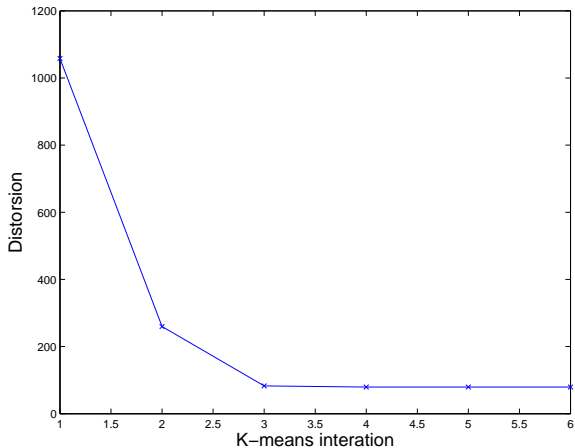
**2 Relocation step** :  $\mu_k^{(t+1)} = \frac{\sum_{i=1}^n z_{ik}^{(t)} \mathbf{x}_i}{\sum_{i=1}^n z_{ik}^{(t)}}$ ,

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

# Example



# Example



# $K$ -means

How to measure uncertainty?

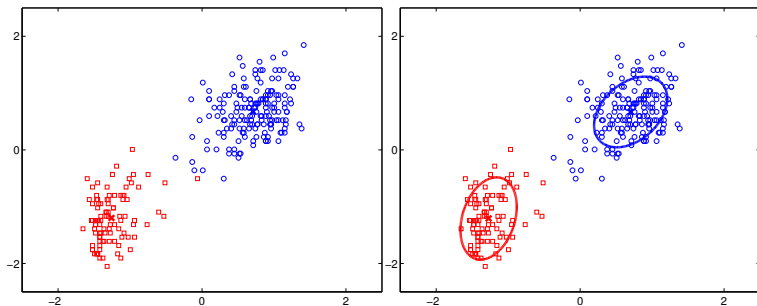


FIGURE –  $K$ -means partition (left) vs GMM-EM partition (right)

## Scientific context

- The data are assumed to represent samples from random variables with unknown probability distributions
- The area of **statistical learning** and **analysis of complex data**.
- **Data** : Complex data  $\hookrightarrow$  *heterogeneous, temporal/dynamical, high-dimensional/functional, incomplete,...*
- **Objective** : Transform the data into knowledge :  
 $\hookrightarrow$  **Reconstruct hidden structure/information, groups/hierarchy of groups, summarizing prototypes, underlying dynamical processes, etc**

## Modeling framework

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

**Generative** formulation :

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

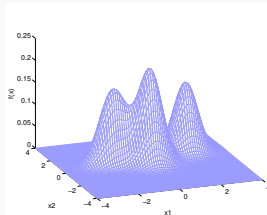
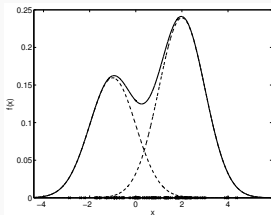
$$x|\mathbf{z} \sim f(x|\mathbf{z}, \boldsymbol{\theta})$$

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

# Mixture models [McLachlan and Peel., 2000]

## Mixture modeling framework

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



- Generative model

$$\begin{aligned} z &\sim \mathcal{M}(1; \pi_1, \dots, \pi_K) \\ x|z &\sim f(x|\theta_z) \end{aligned}$$

→ learn  $\theta$  from the data

# Model-Based Clustering

Clustering based on finite mixture models [McLachlan, 1982, McLachlan and Basford, 1988, Banfield and Raftery, 1993, McLachlan and Peel., 2000]

eg. ; Gaussian mixture models (GMMs) :

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

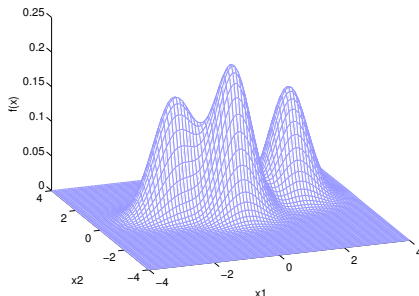


FIGURE – An example of a three-component Gaussian mixture density in  $\mathbb{R}^2$ .

# Mixtures and the EM algorithm (Model-Based Clustering)

Finite Mixture Models [McLachlan and Peel., 2000]

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



# Mixtures and the EM algorithm (Model-Based Clustering)

Finite Mixture Models [McLachlan and Peel., 2000]

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

Maximum-Likelihood Estimation

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

# Mixtures and the EM algorithm (Model-Based Clustering)

## Finite Mixture Models [McLachlan and Peel., 2000]

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

## Maximum-Likelihood Estimation

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

## The EM algorithm [Dempster et al., 1977, McLachlan and Krishnan, 2008]

$$\boldsymbol{\theta}^{new} \in \arg \max_{\boldsymbol{\theta} \in \Omega} \mathbb{E}[\log L_c(\boldsymbol{\theta}) | \mathcal{D}, \boldsymbol{\theta}^{old}]$$

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

# Mixtures and the EM algorithm (Model-Based Clustering)

## Finite Mixture Models [McLachlan and Peel., 2000]

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

## Maximum-Likelihood Estimation

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

## The EM algorithm [Dempster et al., 1977, McLachlan and Krishnan, 2008]

$$\boldsymbol{\theta}^{new} \in \arg \max_{\boldsymbol{\theta} \in \Omega} \mathbb{E}[\log L_c(\boldsymbol{\theta}) | \mathcal{D}, \boldsymbol{\theta}^{old}]$$

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

## Clustering

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

# EM for Gaussian mixture models

**1 E-Step** : calculates the posterior component memberships :

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

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

**2 M-Step** : parameter updates :

$$\pi_k^{(q+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(q)}}{n} = \frac{n_k^{(q)}}{n},$$

$$\boldsymbol{\mu}_k^{(q+1)} = \frac{1}{n_k^{(q)}} \sum_{i=1}^n \tau_{ik}^{(q)} \mathbf{x}_i,$$

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

# Example

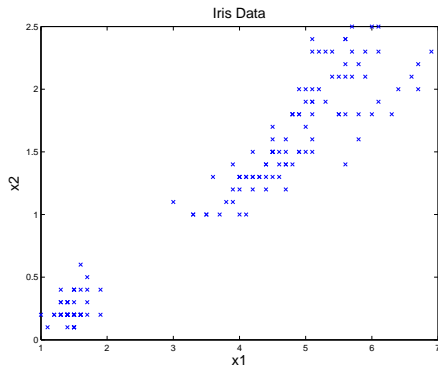
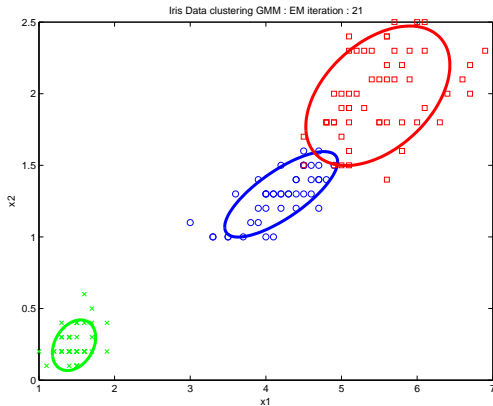
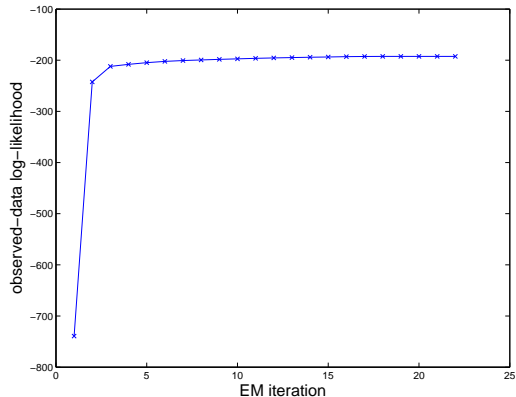


FIGURE — A three-class example of a real data set : Iris data of Fisher.

# Example



# Example



# Mixtures in a high-dimensional setting

- Parsimonious GMMs [Banfield and Raftery, 1993, Celeux and Govaert, 1995] :

- ▶ Eigenvalue decomposition of the covariance matrices :

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

- ▶  $\lambda_k$  the volume of the  $k$ th cluster (the amount of space of the cluster).
- ▶  $\mathbf{D}_k = (\mathbf{v}_{k1}, \dots, \mathbf{v}_{kp})$  orthogonal matrix of eigenvectors  $\mathbf{v}$  of  $\Sigma_k$  : determines the orientation of the cluster.
- ▶  $\mathbf{A}_k = \text{diag}(\lambda_{k1}, \dots, \lambda_{kp}) / |\Sigma_k|^{1/p}$  a normalized diagonal matrix (its determinant is 1) of the eigenvalues of  $\Sigma_k$  arranged in a decreasing order. This matrix is associated with the shape of the cluster.



# Mixtures in a high-dimensional setting

- Parsimonious GMMs [Banfield and Raftery, 1993, Celeux and Govaert, 1995] :

- ▶ Eigenvalue decomposition of the covariance matrices :

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

- ▶  $\lambda_k$  the volume of the  $k$ th cluster (the amount of space of the cluster).
- ▶  $\mathbf{D}_k = (\mathbf{v}_{k1}, \dots, \mathbf{v}_{kp})$  orthogonal matrix of eigenvectors  $\mathbf{v}$  of  $\Sigma_k$  : determines the orientation of the cluster.
- ▶  $\mathbf{A}_k = \text{diag}(\lambda_{k1}, \dots, \lambda_{kp}) / |\Sigma_k|^{1/p}$  a normalized diagonal matrix (its determinant is 1) of the eigenvalues of  $\Sigma_k$  arranged in a decreasing order. This matrix is associated with the shape of the cluster.

for  $p > n$  :

- use regularization (LASSO etc) of the log-likelihood
- Mixtures of Factor Analyzers [McLachlan et al., 2003] (or extensions MCFA, MCUFSA..)

$$\Sigma_k = \mathbf{B}_k \mathbf{B}_k^T + \Lambda_k :$$

$\mathbf{B}_k$  is a  $p \times q$  (with  $q < p$ ) matrix and  $\Lambda_k$  is a diagonal matrix.

$\hookrightarrow (\mathbf{B}_k \mathbf{B}_k^T + \Lambda_k)^{-1}$  and  $|\mathbf{B}_k \mathbf{B}_k^T + \Lambda_k|$  are calculated in a  $q$ -dimensional space !

# How many clusters in the data ?

- 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.

$$\text{score}(\text{model}) = \text{error}(\text{model}) + \text{penalty}(\text{model complexity})$$

which will be minimized.

- Here the complexity of a model  $\mathcal{M}$  is related to the number of its (free) parameters  $\nu$

# Model selection

- Akaike Information Criterion (AIC) [Akaike, 1974] :

$$\text{AIC}(\mathcal{M}_m) = \log L(\hat{\boldsymbol{\theta}}_m) - \nu_m$$

- Bayesian Information Criterion (BIC) [Schwarz, 1978] :

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

- Integrated Classification Likelihood (ICL) [Biernacki et al., 2000] :

$$\text{ICL}(\mathcal{M}_m) = \log L_c(\hat{\boldsymbol{\theta}}_m) - \frac{\nu_m \log(n)}{2}$$

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

$$\nu = \underbrace{(K-1)}_{\pi_k \text{'s}} + \underbrace{K \times d}_{\{\mu_k\}} + \underbrace{K \times \frac{d \times (d+1)}{2}}_{\{\Sigma_k\}} = \frac{K \times (d+1) \times (d+2)}{2} - 1.$$

# Examples

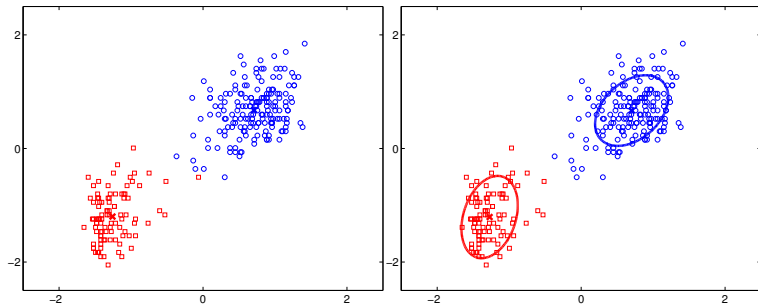


FIGURE – Clustering results obtained with  $K$ -means algorithm (left) with  $K = 2$  and the EM algorithm (right). The cluster centers are shown by the red and blue crosses and the ellipses are the contours of the Gaussian component densities at level 0.4 estimated by EM. The number of clusters for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

# Examples

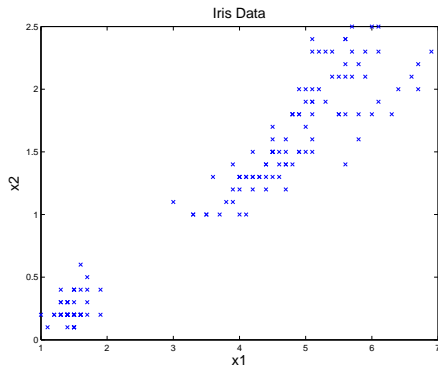


FIGURE — A three-class example of a real data set : Iris data of Fisher.

# Examples

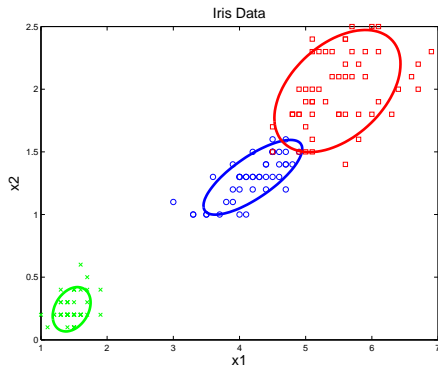
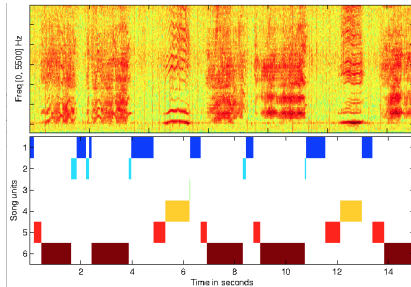
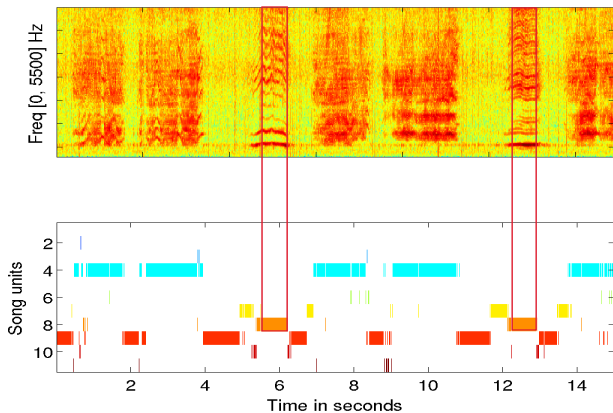


FIGURE — Iris data : Clustering results with EM for a GMM and AIC.

# Unsupervised Sparse Signal Decomposition

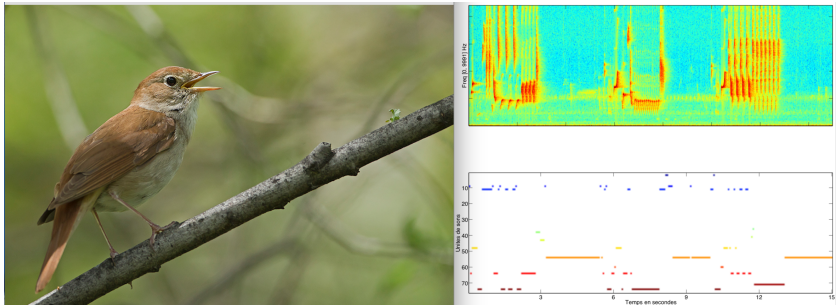


# Unsupervised Sparse Signal Decomposition





# Unsupervised Sparse Signal Decomposition



# Outline

- 1 Introduction
- 2 Latent data models for temporal data segmentation
- 3 Mixture models for functional data analysis
- 4 Mixture-of-Experts for fitting complex non-normal distributions

# Outline

## 1 Introduction

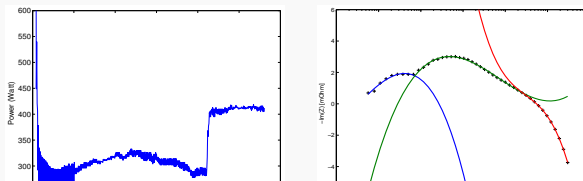
## 2 Latent data models for temporal data segmentation

- Piecewise Regression
- Regression with hidden logistic process
- Multiple hidden process regression
- **SaMURaiS** : Open-Source Software

## 3 Mixture models for functional data analysis

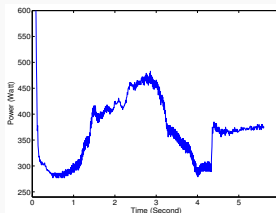
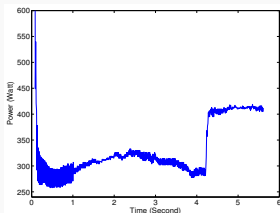
## 4 Mixture-of-Experts for fitting complex non-normal distributions

### Temporal data with regime changes



# 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

# Latent data 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, signal processing, machine learning, finance
- The observed time series is generated by an underlying process  
 $\hookrightarrow$  segmentation  $\equiv$  recovering the parameters the process' states.
- Conventional solutions are subject to limitations in the control of the transitions between these states
- $\hookrightarrow$  Propose latent data modeling for segmentation and approximation
- $\hookrightarrow$  segmentation  $\equiv$  inferring the model parameters and the underlying process

# Piecewise regression [McGee & Carleton 70], Chamroukhi et al. [2009]

- The data :  $((t_1, y_1), \dots, (t_n, y_n))$  where  $y_i$  is the observation et time  $t_i$
- The piecewise polynomial regression model is defined as :

$$\forall i = 1, \dots, n, \quad x_i = \begin{cases} \beta_1^T \mathbf{x}_i + \varepsilon_{i1} & \text{if } i \in I_1 \\ \beta_2^T \mathbf{x}_i + \varepsilon_{i2} & \text{if } i \in I_2 \\ \vdots & \\ \beta_K^T \mathbf{x}_i + \varepsilon_{iK} & \text{if } i \in I_K \end{cases},$$

- ▶  $I_k = ]\gamma_k.. \gamma_{k+1}]$  : indexes of elements of segment  $k$  with  $(\gamma_1=0$  and  $\gamma_{K+1}=n)$ .
- ▶  $\mathbf{x}_i = (1, t_i, \dots, t_i^p)^T$  : time-dependent covariates vector
- ▶  $\beta_k \in \mathbb{R}^{p+1}$  : regression coefficients vector for the  $k^{th}$  segment
- ▶  $\varepsilon_{ik} \sim \mathcal{N}(0, \sigma_k^2)$  : independent additive Gaussian noise on the segment  $k$ .

## The model parameters

$(\psi, \gamma)$  with  $\psi = (\beta_1^T, \dots, \beta_K^T, \sigma_1^2, \dots, \sigma_K^2)$  and  $\gamma = (\gamma_1, \dots, \gamma_{K+1})^T$ .

# Parameter estimation piecewise regression

Maximize the likelihood of  $(\boldsymbol{\psi}, \boldsymbol{\gamma})$  or equivalently minimize, with respect to  $(\boldsymbol{\psi}, \boldsymbol{\gamma})$  :

$$J(\boldsymbol{\psi}, \boldsymbol{\gamma}) = \sum_{k=1}^K \sum_{i \in I_k} \left[ \log \sigma_k^2 + \frac{(y_i - \boldsymbol{\beta}_k^T \mathbf{x}_i)^2}{\sigma_k^2} \right].$$

- Global optimization using by dynamic programming [Bellman 61; Stone 61; Lechevallier 90, C. 2009] since the criterion  $J$  is additive on  $k$

## Time series approximation and segmentation

- $\hat{y}_i = \sum_{k=1}^K \hat{z}_{ik} \hat{\boldsymbol{\beta}}_k^T \mathbf{x}_i \quad ; \quad \forall i = 1, \dots, n$
- $\hat{z}_{ik} = 1$  if  $i \in \hat{I}_k = (\hat{\gamma}_k, \hat{\gamma}_{k+1}]$  ( $y_i$  belongs to the  $k^{th}$  segment) and  $\hat{z}_{ik} = 0$  otherwise

- Using dynamic programming can be computationally expensive
- Provides a hard partition  $\Rightarrow$  adapted for regimes with abrupt changes

# 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]

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

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

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

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

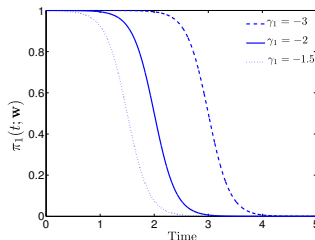
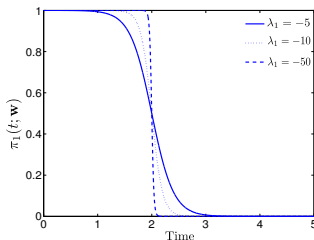
- Both the mixing proportions and the component parameters are time-varying
- Parameter vector of the model :  $\boldsymbol{\theta} = (\mathbf{w}^T, \beta_1^T, \dots, \beta_K^T, \sigma_1^2, \dots, \sigma_K^2)^T$



# Illustration

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

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

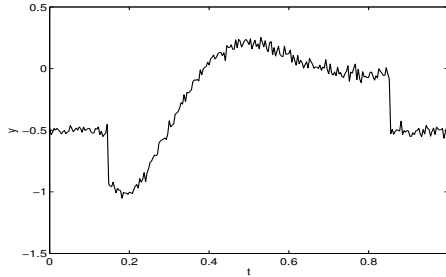


⇒ The parameter  $\lambda_k = w_{k1}$  controls the quality of transitions between regimes

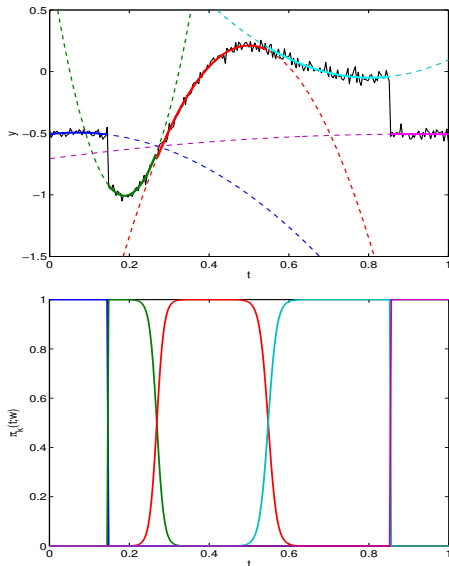
⇒ The parameter  $\gamma_k = w_{k0}/w_{k1}$  is related to the transition time point

- 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

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

- **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)

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

- 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]

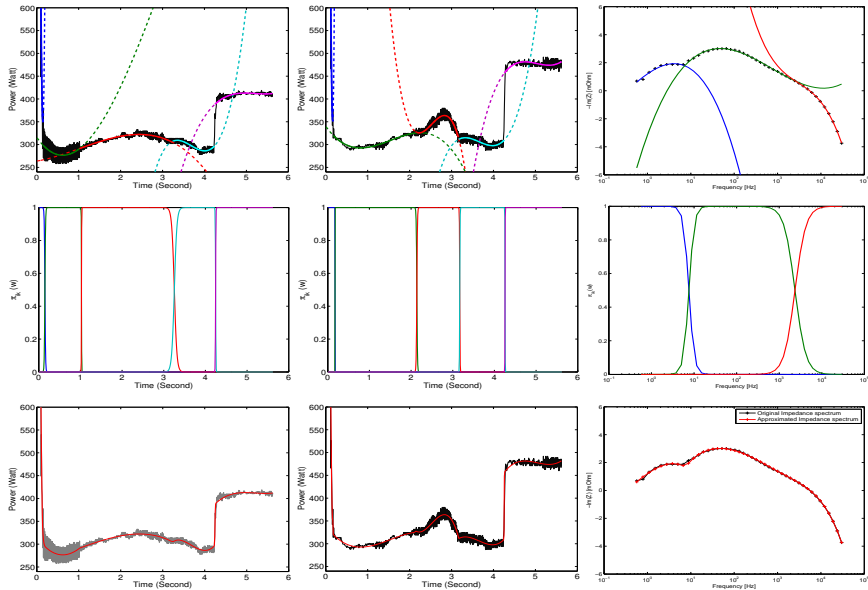
## 2 Curve segmentation :

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

## 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  $\nu_{\boldsymbol{\theta}} = K(p + 4) - 2$ .

# Application to real data





# Joint segmentation of multivariate time series

## Multiple hidden process regression

- Data :  $(\mathbf{y}_1, \dots, \mathbf{y}_n)$  a time series of  $n$  multidimensional observations  $\mathbf{y}_i = (y_i^{(1)}, \dots, y_i^{(d)})^T \in \mathbb{R}^d$  observed at instants  $\mathbf{t} = (t_1, \dots, t_n)$ .
- Model

$$\begin{aligned} y_i^{(1)} &= \beta_{z_i}^{(1)T} \mathbf{x}_i + \sigma_{z_i}^{(1)} \epsilon_i \\ &\vdots \\ y_i^{(d)} &= \beta_{z_i}^{(d)T} \mathbf{x}_i + \sigma_{z_i}^{(d)} \epsilon_i \end{aligned}$$

Vectorial form :  $\mathbf{y}_i = \mathbf{B}_{z_i}^T \mathbf{x}_i + \mathbf{e}_i$  ;  $\mathbf{e}_i \sim \mathcal{N}(\mathbf{0}, \Sigma_{z_i})$ ,  $(i = 1, \dots, n)$

- The latent process  $\mathbf{z} = (z_1, \dots, z_n)$  simultaneously governs the univariate time series components

$\hookrightarrow$  Multiple regression with hidden logistic process : Multiple RHLP [6]

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

# Multiple hidden Markov model regression

- MHMMR : Estimation by the EM algorithm (as for HMMs)

↪ Solve multiple regression problems

## Application to human activity time series

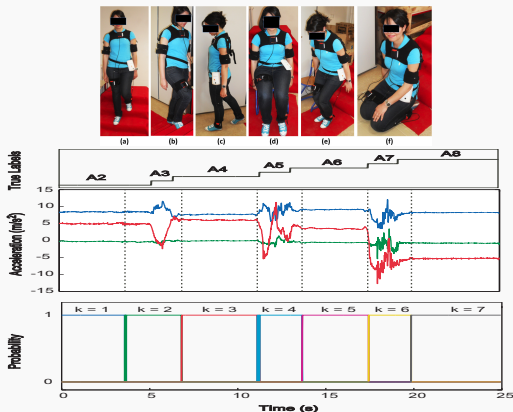


FIGURE – MHMMR Segmentation of acceleration data issued from three body-worn sensors

# Multiple regression with hidden logistic process

- MRHLP : Estimation by the EM algorithm (as for the RHLP)

↔ Solve multiple regression problems

## Application to human activity time series

Problem : Activity recognition from multivariate acceleration time series

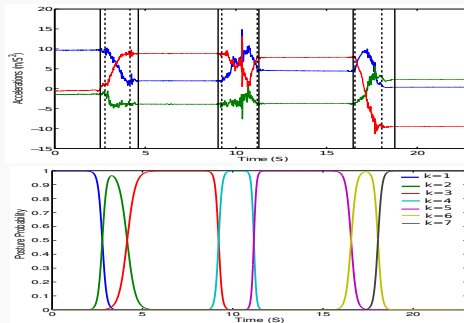


FIGURE – MRHLP segmentation of acceleration data issued from three body-worn sensors

# SaMUraiS : open source software for time-series



**SaMUraiS : StAtistical Models for the Unsupervised segmentAtion of time-Series<sup>1</sup>**

► [R software](#)

► [Matlab software](#)

---

1. credit : pictures above created via dreamscapeapp.com

# SaMUraIS : open source software for time-series

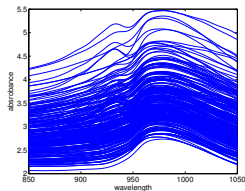
## Available algorithms and Packages

- **RHLP** : Regression with Hidden Logistic Process [▶ R software](#) [▶ Matlab software](#)
- **HMMR** : Hidden Markov Model Regression [▶ R software](#) [▶ Matlab software](#)
- **PWR** : Piece-Wise Regression [▶ R software](#) [▶ Matlab software](#)
- **MRHLP** : Multivariate RHLP [▶ R software](#) [▶ Matlab software](#)
- **MHMMR** : Multivariate HMMR [▶ R software](#) [▶ Matlab software](#)

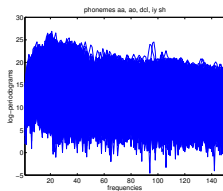
# Outline

- 1 Introduction
- 2 Latent data models for temporal data segmentation
- 3 Mixture models for functional data analysis
  - Mixture of piecewise regressions
  - Mixture of hidden logistic process regressions
  - Mixture of hidden Markov model regression
  - Functional discriminant analysis
  - **FLaMingoS** : Open-Source Software
- 4 Mixture-of-Experts for fitting complex non-normal distributions

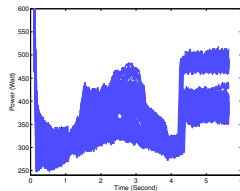
# Functional data are increasingly frequent



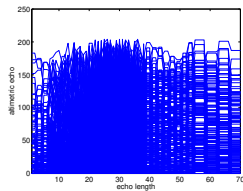
Tecator data



Phonemes curves



Railway switch curves



Satellite waveforms

# Statistical analysis of functional data

A broad literature :

[James and Hastie, 2001, James and Sugar, 2003]

[Ramsay and Silverman, 2005]

[Ferraty and Vieu, 2006]

[Ramsay et al., 2011]

[Bouveyron and Jacques, 2011]

[Samé et al., 2011]

[Delaigle et al., 2012]

[Jacques and Preda, 2014]

[Bouveyron et al., 2018]

[Qiao et al., 2018]

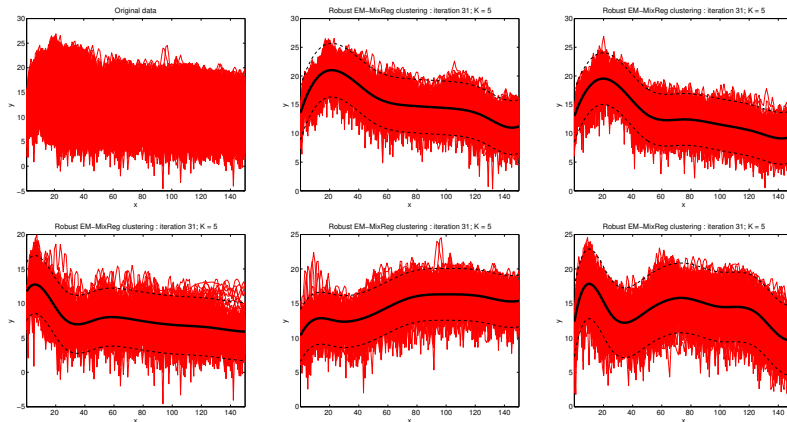
A review can be found in [Chamroukhi and Nguyen, 2018] [▶ pdf available here](#)

- Functional regression
- Functional classification
- Functional clustering, including model-based
- Functional graphical models
- ...



# Clustering of functional data

Phonemes data set<sup>2</sup> :  $n = 1000$  log-periodograms for  $m = 150$  frequencies



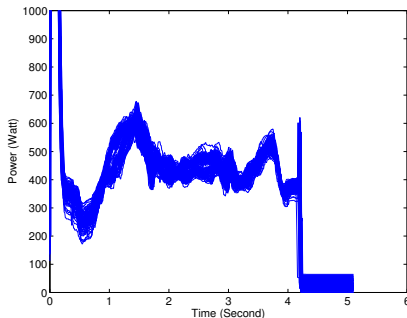
2. Data from <http://www.math.univ-toulouse.fr/staph/npfda/>, used in Ferraty and Vieu [2003]

# Clustering of functional data

Clustering real curves of high-speed railway-switch operations

Data :  $n = 115$  curves of  $m \simeq 510$  observations

$K = 2$  clusters : operating state without/with possible defect

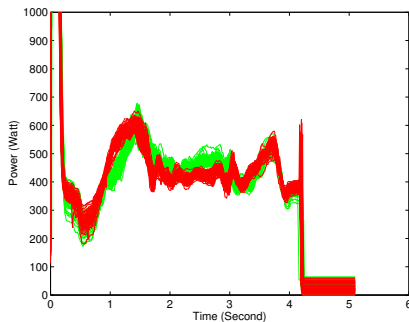


# Clustering switch operations

## Clustering real curves of high-speed railway-switch operations

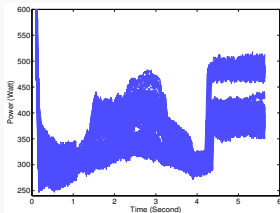
Data :  $n = 115$  curves of  $m \simeq 510$  observations

$K = 2$  clusters : operating state without/with possible defect

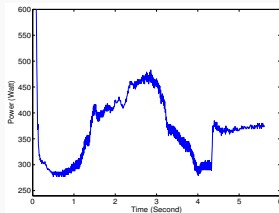


# Functional data analysis context

## Unsupervised analysis of heterogeneous curves with regime changes



Railway switch curves



An individual curve

## Objectives

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

# Outline

- 1 Introduction
- 2 Latent data models for temporal data segmentation
- 3 Mixture models for functional data analysis
  - Mixture of piecewise regressions
  - Mixture of hidden logistic process regressions
  - Mixture of hidden hidden Markov model regression
  - Functional discriminant analysis
  - **FLaMingoS** : Open-Source Software
- 4 Mixture-of-Experts for fitting complex non-normal distributions

# 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)  
 $\Rightarrow$  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)
- 3 Regression model with Markov process (HMMR)

# Piecewise regression mixture model (PWRM) [9]

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

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

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

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



- 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

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

- The conditional expected complete-data log-likelihood

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

# 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 $\Psi^{(q+1)} = \arg \max_{\Psi} Q(\Psi, \Psi^{(q)})$

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

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

$\mathbf{y}_{ir}$  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 Hébrail et al. [2010] :

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

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  $\mathcal{O}(I_{\text{EM}} K R n m^2 p^3)$  : Significant computational load for large

# Simulation results

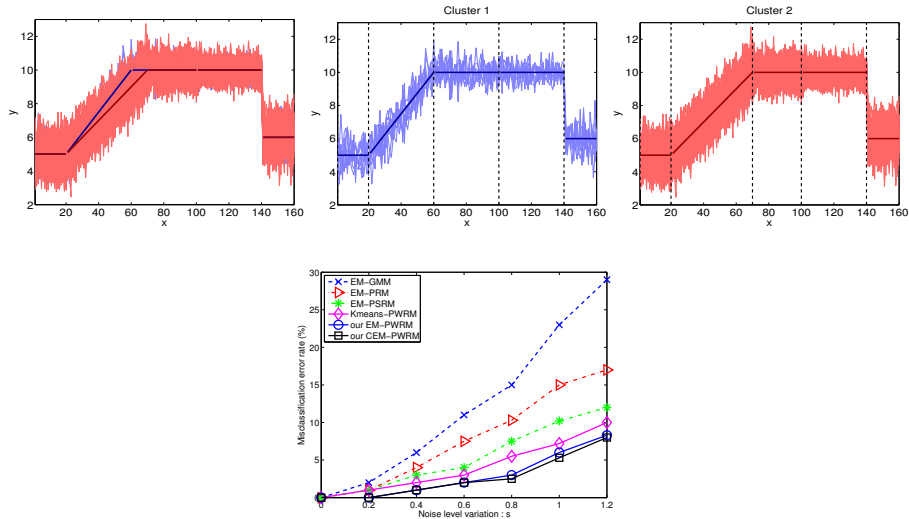
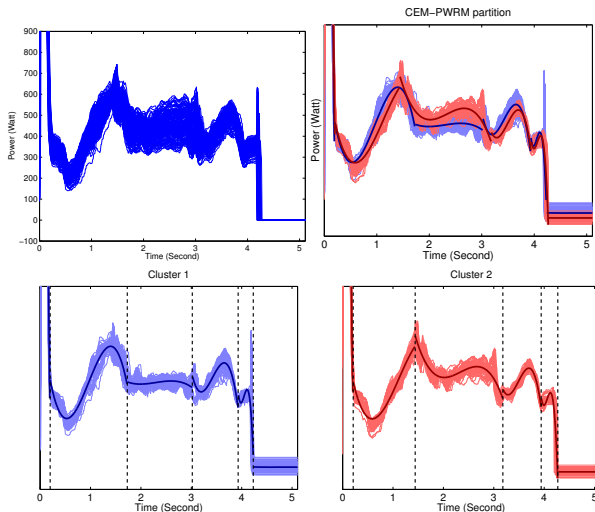


FIGURE – Misclassification error rate versus the noise level variation.

# Application to switch operation curves

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

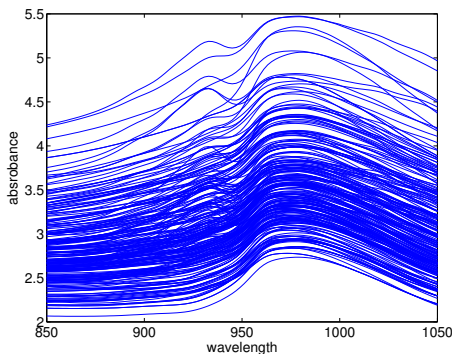
Each curve is composed of  $R = 6$  electromechanical phases (regimes)



# Application to Tecator data

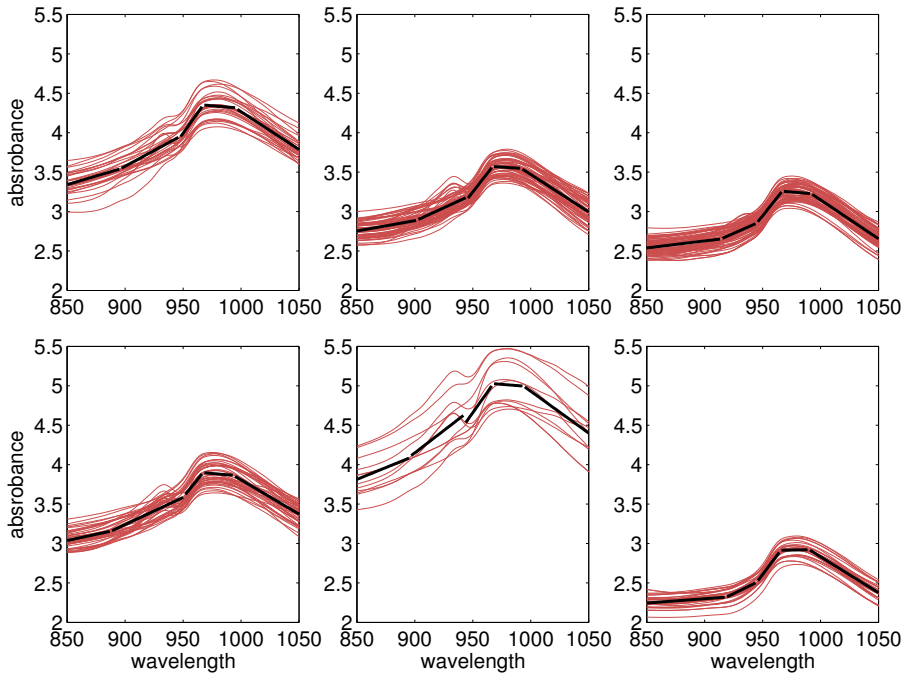
The Tecator data set<sup>3</sup> contains  $n = 240$  spectra with  $m = 100$  observations for each spectrum

Data considered in the same setting as in Hébrail et al. [2010] (six clusters, each cluster is approximated by five linear segments ( $R = 5, p = 1$ ))



---

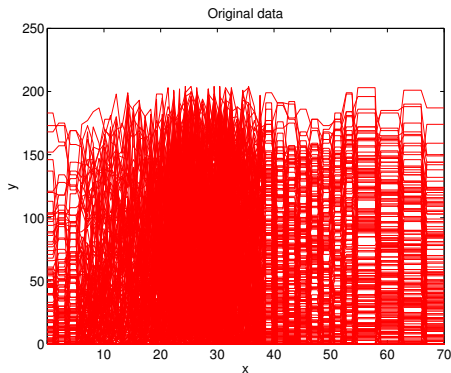
3. Tecator data are available at <http://lib.stat.cmu.edu/datasets/tecator>.



# Topex/Poseidon satellite data

The Topex/Poseidon radar satellite data<sup>4</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 Hébrail et al. [2010].

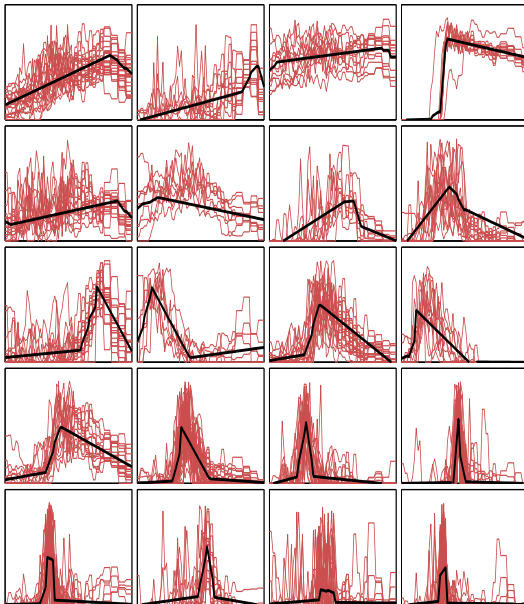


---

4. 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 Hébrail et al. [2010]
- 
- 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 Hébrail et al. [2010].
  - 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 :

$\hookrightarrow$  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

# 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

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

# EM-MixRHLP algorithm

## E-step

- The posterior cluster memberships :

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

- the posterior regime memberships :

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

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

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

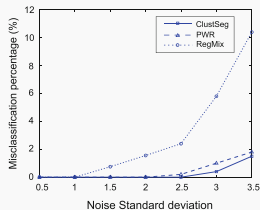
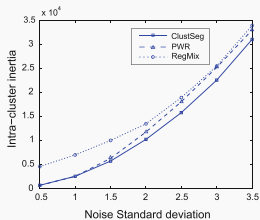
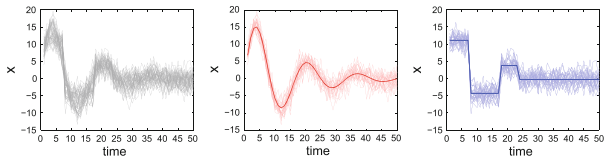
- Regression parameters update : Analytic weighted least-squares problems

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

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  $\mathcal{O}(I_{\text{EM}} I_{\text{IRLS}} K R^3 n m p^3)$  ( $K$ -means like algo. for PWR is in  $\mathcal{O}(I_{\text{KM}} K R n m^2 p^3)$   $\hookrightarrow$  computationally attractive for large  $m$  with moderate value of  $R$ ).

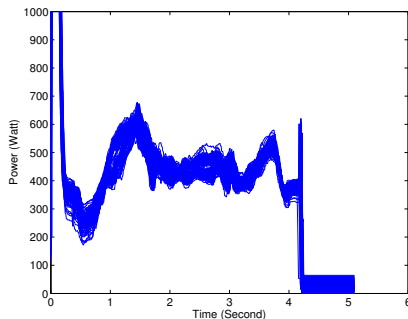
# EM-MixRHLP clustering of simulated data



# 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

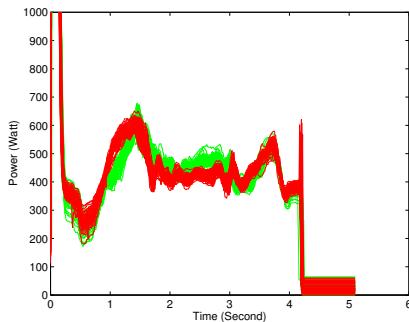




# Clustering switch operations

**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]

# 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

# FLaMingoS : open source softw. for functional data



**FLaMingoS** : Functional Latent data Models for clustering heterogeneous time-Series

► R software

► Matlab software

# FLaMingoS : open source softw. for functional data

## Available algorithms and Packages

- **mixRHLP** : Mixture of Regressions with Hidden Logistic Processes [▶ R software](#)  
[▶ Matlab software](#)
- **mixHMM** : Mixtures of Hidden Markov Models [▶ R software](#) [▶ Matlab software](#)
- **mixHMMR** : Mixtures of Hidden Markov Model Regressions [▶ R software](#)  
[▶ Matlab software](#)
- **PWRM** : Piece-Wise Regression Mixture [▶ R software](#) [▶ Matlab software](#)

## Coming soon : Learning of regression mixtures with unknown number of components

- unsupPRM** [▶ R software](#) [▶ Matlab software](#)
- unsupSRM** [▶ R software](#) [▶ Matlab software](#)
- unsupBSRM** [▶ R software](#) [▶ Matlab software](#)

# Outline

- 1 Introduction
- 2 Latent data models for temporal data segmentation
- 3 Mixture models for functional data analysis
- 4 Mixture-of-Experts for fitting complex non-normal distributions
  - The skew-normal mixture of experts model
  - The  $t$  mixture of experts model
  - The skew  $t$  mixture of experts model
  - **MEteorits** : Open-Source Software

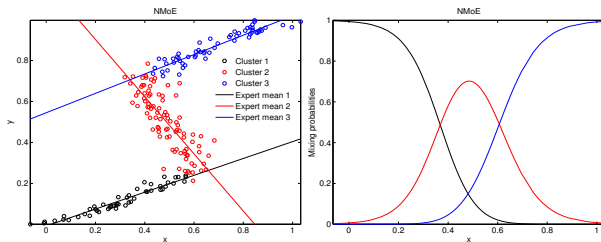
## Mixture-of-Experts (MoE) modeling framework

- Data : an observed i.i.d sample of the pair  $(\mathbf{X}, Y)$  where the response  $Y \in \mathbb{R}$  for the vector of predictors  $\mathbf{X} \in \mathbb{R}^p$  is governed by a hidden categorical variable  $Z$   
 $z_i \in [K]$  is the expert label for  $(\mathbf{X}_i, Y_i)$
- Mixture of experts (MoE) [Jacobs et al., 1991, Jordan and Jacobs, 1994] :

$$f(y|\mathbf{x}; \Psi) = \sum_{k=1}^K \underbrace{\pi_k(\mathbf{x}; \mathbf{w})}_{\text{Gating network}} \underbrace{f_k(y|\mathbf{x}; \Psi_k)}_{\text{Expert Network}}$$

- Gating network (e.g softmax) :  $\pi_k(\mathbf{x}; \mathbf{w}) = \frac{\exp(w_{k0} + \mathbf{w}_k^T \mathbf{x})}{1 + \sum_{\ell=1}^{K-1} \exp(w_{\ell 0} + \mathbf{w}_{\ell}^T \mathbf{x})}$
  - Experts network (e.g Gaussian regressors) :  $f_k(y|\mathbf{x}; \Psi_k) = \phi(y; \mu(\mathbf{x}; \beta_k), \sigma_k^2)$  with parametric (non-)linear regression functions  $\mu(\mathbf{x}; \beta_k)$
  - parameter vector  $\Psi = (\mathbf{w}^T, \Psi_1^T, \dots, \Psi_K^T)^T$
- ↪ For a review, see Nguyen and Chamroukhi [2018] [pdf available here](#)

# Illustration





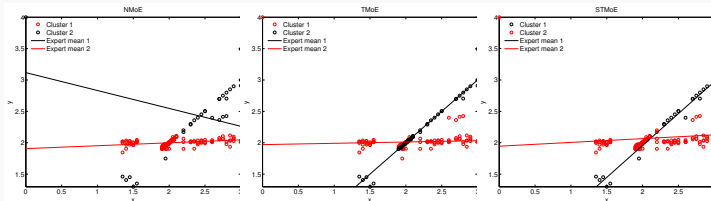


FIGURE – Fitting MoE to the tone data set with ten outliers (0, 4).

For a set of data containing a group or groups of observations with asymmetric behavior, heavy tails or atypical observations, the use of normal experts may be unsuitable and can unduly affect the fit

## Objectives

- Overcome these limitations of MoE modeling with the normal distribution.
- We proposed three non-normal derivations including two robust mixture of experts (MoE) models.  $\leftrightarrow$  suitable to accommodate data which exhibit additional features such as skewness, heavy-tails and which may be affected by atypical data [Chamroukhi, 2017, 2016a,b]

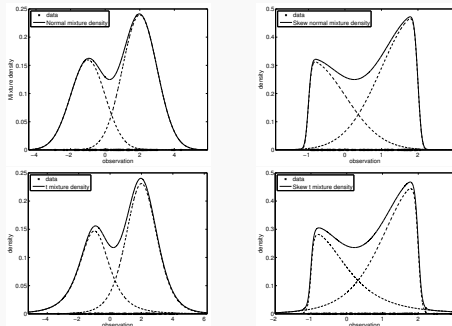
- 1 Introduction
- 2 Latent data models for temporal data segmentation
- 3 Mixture models for functional data analysis
- 4 Mixture-of-Experts for fitting complex non-normal distributions
  - The skew-normal mixture of experts model
  - The  $t$  mixture of experts model
  - The skew  $t$  mixture of experts model
  - **MEteorits** : Open-Source Software

# Non-normal mixtures of experts

## Non-normal mixtures of experts (NNMoE)

- 1 the  $t$  MoE (TMoE) (Robustness, heavy tails) [11]
- 2 the skew-normal MoE (SNMoE) (skewness) [14]
- 3 the skew- $t$  MoE (STMoE) (skewness, robustness, heavy tails) [15]

## Non-normal mixtures



$$\pi_k = [0.4, 0.6], \mu_k = [-1, 2]; \sigma_k = [1, 1]; \nu_k = [3, 7]; \lambda_k = [14, -12];$$

# The skew-normal mixture of experts model

- The SNMoE is defined as

$$f(y|\mathbf{r}, \mathbf{x}; \boldsymbol{\Psi}) = \sum_{k=1}^K \pi_k(\mathbf{r}; \boldsymbol{\alpha}) \text{SN}(y; \mu(\mathbf{x}; \boldsymbol{\beta}_k), \sigma_k^2, \lambda_k)$$

where each expert component  $k$  has indeed a skew-normal distribution, whose density is defined by (1).

- *The skew-normal distribution* [Azzalini, 1985, 1986] with location  $\mu \in \mathbb{R}$ , scale  $\sigma^2 \in (0, \infty)$  and skewness parameter  $\lambda \in \mathbb{R}$  has density

$$\text{SN}(y; \mu, \sigma^2, \lambda) = \frac{2}{\sigma} \phi\left(\frac{y - \mu}{\sigma}\right) \Phi\left(\lambda \left(\frac{y - \mu}{\sigma}\right)\right)$$

$\phi(\cdot)$  and  $\Phi(\cdot)$  denote the pdf and the cdf of the standard normal distribution.

- The parameter vector is  $\boldsymbol{\Psi} = (\boldsymbol{\alpha}_1^T, \dots, \boldsymbol{\alpha}_{K-1}^T, \boldsymbol{\Psi}_1^T, \dots, \boldsymbol{\Psi}_K^T)^T$  with  $\boldsymbol{\Psi}_k = (\boldsymbol{\beta}_k^T, \sigma_k^2, \lambda_k)^T$  the parameter vector for the  $k$ th skewed-normal expert component.
- It is obvious to see that if the skewness parameter  $\lambda_k = 0$  for each  $k$ , the SNMoE model reduces to the NMoE model.

# The skew-normal mixture of experts model

- **Stochastic representation of the SNMoE** : A random variable  $Y_i$  is said to follow the SNMoE model if it has the following representation :

$$Y_i = \mu(\mathbf{x}_i; \boldsymbol{\beta}_{z_i}) + \delta_{z_i} \sigma_{z_i} |U_i| + \sqrt{1 - \delta_{z_i}^2} \sigma_{z_i} E_i.$$

where  $U$  and  $E$  be independent random variables following the standard normal distribution  $\mathcal{N}(0, 1)$  with pdf  $\phi(\cdot)$ ,  $|U|$  denotes the magnitude of  $U$  and  $\delta_{z_i} = \frac{\lambda_{z_i}}{\sqrt{1 + \lambda_{z_i}^2}}$  where  $Z_i \in \{1, \dots, K\}$  is a categorical variable following the multinomial distribution

$$Z_i | \mathbf{r}_i \sim \text{Mult}(1; \pi_1(\mathbf{r}_i; \boldsymbol{\alpha}), \dots, \pi_K(\mathbf{r}_i; \boldsymbol{\alpha}))$$

where each  $\pi_{z_i}(\mathbf{r}_i; \boldsymbol{\alpha}) = \mathbb{P}(Z_i = z_i | \mathbf{r}_i)$  is given by the logistic function.

- **Hierarchical representation of the SNMoE**

$$Y_i | u_i, Z_{ik} = 1, \mathbf{x}_i \sim \mathcal{N}\left(\mu(\mathbf{x}_i; \boldsymbol{\beta}_k) + \delta_k |u_i|, (1 - \delta_k^2) \sigma_k^2\right),$$

$$U_i | Z_{ik} = 1 \sim \mathcal{N}(0, \sigma_k^2),$$

$$\mathbf{Z}_i | \mathbf{r}_i \sim \text{Mult}(1; \pi_1(\mathbf{r}_i; \boldsymbol{\alpha}), \dots, \pi_K(\mathbf{r}_i; \boldsymbol{\alpha}))$$

where  $Z_{ik}$  are the binary latent component-indicators such that  $Z_{ik} = 1$  iff  $Z_i = k$ ,  $\mathbf{Z}_i = (Z_{i1}, \dots, Z_{iK})$  and  $\delta_k = \frac{\lambda_k}{\sqrt{1 + \lambda_k^2}}$

# Maximum likelihood parameter estimation

- Given an observed i.i.d sample of  $n$  observations  $\{(y_i, \mathbf{x}_i, \mathbf{r}_i)\}_{i=1}^n$ , the parameter vector  $\Psi$  of the SNMoE model can be estimated by maximizing the observed-data log-likelihood :

$$\log L(\Psi) = \sum_{i=1}^n \log \sum_{k=1}^K \pi_k(\mathbf{r}_i; \alpha) \text{SN} \left( y; \mu(\mathbf{x}; \beta_k), \sigma_k^2, \lambda_k \right) .$$

- $\Rightarrow$  A dedicated Expectation Conditional Maximization (ECM) algorithm
- The ECM algorithm [Meng and Rubin, 1993] is an EM variant that mainly aims at addressing the optimization problem in the M-step of the EM algorithm. In ECM, the M-step is performed by several conditional maximization (CM) steps by dividing the parameter space into sub-spaces. The parameter vector updates are then performed sequentially, one coordinate block after another in each sub-space.
- This is also the generative process for sampling data according to the SNMoE model

# MLE via the ECM algorithm

- The complete-data log-likelihood of  $\Psi$ , where the complete-data are  $\{y_i, z_i, u_i, \mathbf{x}_i, \mathbf{r}_i\}_{i=1}^n$ , is given by :

$$\log L_c(\Psi) = \log L_c(\alpha) + \sum_{k=1}^K \log L_c(\Psi_k),$$

with

$$\begin{aligned}\log L_c(\alpha) &= \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log \pi_k(\mathbf{r}_i; \alpha), \\ \log L_c(\Psi_k) &= \sum_{i=1}^n Z_{ik} \left[ -\log(2\pi) - \log(\sigma_k^2) - \frac{1}{2} \log(1 - \delta_k^2) \right. \\ &\quad \left. - \frac{d_{ik}^2}{2(1 - \delta_k^2)} + \frac{\delta_k d_{ik} u_i}{(1 - \delta_k^2)\sigma_k} - \frac{u_i^2}{2(1 - \delta_k^2)\sigma_k^2} \right],\end{aligned}$$

where  $d_{ik} = \frac{y_i - \mu(\mathbf{x}_i; \beta_k)}{\sigma_k}$ .

# ECM for the SNMoE : E-Step

**E-Step** calculates the  $Q$ -function

$$Q(\Psi; \Psi^{(m)}) = \mathbb{E}[\log L_c(\Psi) | \{y_i, \mathbf{x}_i, \mathbf{r}_i\}_{i=1}^n; \Psi^{(m)}] = Q_1(\alpha; \Psi^{(m)}) + \sum_{k=1}^K Q_2(\Psi_k; \Psi^{(m)})$$

with

$$\begin{aligned} Q_1(\alpha; \Psi^{(m)}) &= \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}^{(m)} \log \pi_k(\mathbf{r}_i; \alpha), \\ Q_2(\Psi_k; \Psi^{(m)}) &= \sum_{i=1}^n \tau_{ik}^{(m)} \left[ -\log(2\pi) - \log(\sigma_k^2) - \frac{1}{2} \log(1 - \delta_k^2) \right. \\ &\quad \left. + \frac{\delta_k d_{ik} e_{1,ik}^{(m)}}{(1 - \delta_k^2) \sigma_k} - \frac{e_{2,ik}^{(m)}}{2(1 - \delta_k^2) \sigma_k^2} - \frac{d_{ik}^2}{2(1 - \delta_k^2)} \right] \end{aligned}$$

where the required conditional expectations (analytic) are given by :

$$\begin{aligned} \tau_{ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [Z_{ik} | y_i, \mathbf{x}_i, \mathbf{r}_i], \\ e_{1,ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [U_i | Z_{ik} = 1, y_i, \mathbf{x}_i, \mathbf{r}_i], \\ e_{2,ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [U_i^2 | Z_{ik} = 1, y_i, \mathbf{x}_i, \mathbf{r}_i]. \end{aligned}$$



**CM-Step 1** Calculate  $\alpha^{(m+1)} = \arg \max_{\alpha} Q_1(\alpha; \Psi^{(m)})$ . does not exist in closed form (Unlike in skew-normal (regression) mixtures)

**The Iteratively Reweighted Least Squares (IRLS) algorithm :**

$$\alpha^{(l+1)} = \alpha^{(l)} - \left[ \frac{\partial^2 Q_1(\alpha, \Psi^{(m)})}{\partial \alpha \partial \alpha^T} \right]_{\alpha=\alpha^{(l)}}^{-1} \frac{\partial Q_1(\alpha, \Psi^{(m)})}{\partial \alpha} \Big|_{\alpha=\alpha^{(l)}}$$

Then, for  $k = 1 \dots, K$ ,

**CM-Step 2** Calculate  $\beta_k^{(m+1)}$  by maximizing  $Q_2(\Psi_k; \Psi^{(m)})$

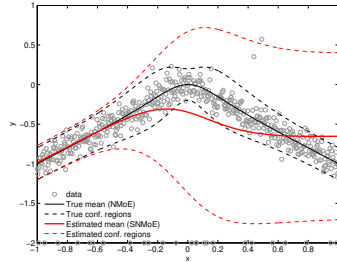
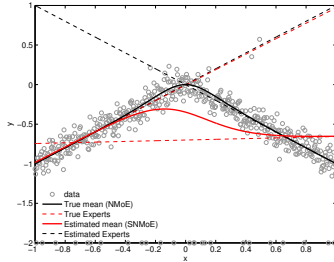
$$\beta_k^{(m+1)} = \left[ \sum_{i=1}^n \tau_{ik}^{(m)} \mathbf{x}_i \mathbf{x}_i^T \right]^{-1} \sum_{i=1}^n \tau_{ik}^{(q)} \left( y_i - \delta_k^{(m)} e_{1,ik}^{(m)} \right) \mathbf{x}_i.$$

**CM-Step 3 :** Calculate  $\sigma_k^{2(m+1)}$  by maximizing  $Q_2(\Psi_k; \Psi^{(m)})$

$$\sigma_k^{2(m+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(m)} \left[ \left( y_i - \beta_k^{T(m+1)} \mathbf{x}_i \right)^2 - 2\delta_k^{(m+1)} e_{1,ik}^{(m)} (y_i - \beta_k^{T(m+1)} \mathbf{x}_i) + e_{2,ik}^{(m)} \right]}{2 \left( 1 - \delta_k^{2(m)} \right) \sum_{i=1}^n \tau_{ik}^{(m)}}.$$

**CM-Step 4** Calculate  $\lambda_k^{(m+1)}$  by maximizing  $Q_2(\Psi_k; \Psi^{(m)})$  : Solution of :

$$\sigma_k^{2(m+1)} \delta_k (1 - \delta_k^2) \sum_{i=1}^n \tau_{ik}^{(m)} + (1 + \delta_k^2) \sum_{i=1}^n \tau_{ik}^{(m)} (y_i - \beta_k^{T(m+1)} \mathbf{x}_i) e_{1,ik}^{(m)} - \delta_k \sum_{i=1}^n \tau_{ik}^{(m)} \left[ e_{2,ik}^{(m)} + \left( y_i - \beta_k^{T(m+1)} \mathbf{x}_i \right)^2 \right] = 0. \text{ root finding (Brent's method [Brent, 1973])}.$$



- SNMoE model is tailored to model the skewness in the data, it may be not adapted to handle data containing groups or a group with heavy-tailed distribution.
- The SNMoE as NMoE may thus be affected by outliers.
- $\Rightarrow$  Handle the problem of sensitivity of normal mixture of experts to outliers and heavy tails.  
 $\hookrightarrow$  robust mixture of experts modeling using the  $t$  distribution.

# The $t$ mixture of experts model

- The proposed  $t$  mixture of experts model extends the  $t$  mixture model, first proposed by Mclachlan and Peel [1998], Peel and Mclachlan [2000] for multivariate data, as well as the regression mixture model using the  $t$ -distribution as in [Bai et al., 2012, Wei, 2012, Ingrassia et al., 2012] to the MoE framework.
- A  $K$ -component TMoE model is defined by :

$$f(y|\mathbf{r}, \mathbf{x}; \Psi) = \sum_{k=1}^K \pi_k(\mathbf{r}; \boldsymbol{\alpha}) t_{\nu_k}(y; \mu(\mathbf{x}; \boldsymbol{\beta}_k), \sigma_k^2, \nu_k).$$

The  $t$ -distribution with location  $\mu \in \mathbb{R}$ , scale  $\sigma^2 \in (0, \infty)$  and degrees of freedom  $\nu \in (0, \infty)$  has the probability density function

$$t_{\nu_k}(y; \mu, \sigma^2, \nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi} \Gamma(\frac{\nu}{2})} \left(1 + \frac{d_y^2}{\nu}\right)^{-\frac{\nu+1}{2}},$$

where  $d_y^2 = \left(\frac{y-\mu}{\sigma}\right)^2$  denotes the squared Mahalanobis distance

- The parameter vector of the TMoE model is given by  $\Psi = (\boldsymbol{\alpha}_1^T, \dots, \boldsymbol{\alpha}_{K-1}^T, \boldsymbol{\Psi}_1^T, \dots, \boldsymbol{\Psi}_K^T)^T$  where  $\boldsymbol{\Psi}_k = (\boldsymbol{\beta}_k^T, \sigma_k^2, \nu_k)^T$
- When the robustness parameter  $\nu_k \rightarrow \infty$  for each experts  $k$ , the TMoE model approaches the NMoE model

# The $t$ mixture of experts model

- **Stochastic representation for the TMoE** Let  $E \sim \phi(\cdot)$ . Suppose that, conditional on the hidden variable  $Z_i = z_i$ , a random variable  $W_i$  is distributed as  $\text{Gamma}(\frac{\nu_{z_i}}{2}, \frac{\nu_{z_i}}{2})$ . Then, given the covariates  $(\mathbf{x}_i, \mathbf{r}_i)$ , a random variable  $Y_i$  is said to follow the TMoE model if

$$Y_i = \mu(\mathbf{x}_i; \boldsymbol{\beta}_{z_i}) + \sigma_{z_i} \frac{E_i}{\sqrt{W_{z_i}}},$$

where the categorical variable  $Z_i | \mathbf{r}_i$  is multinomial

- **Hierarchical representation of the TMoE model**

$$\begin{aligned} Y_i | w_i, Z_{ik} = 1, \mathbf{x}_i &\sim \mathcal{N}\left(\mu(\mathbf{x}_i; \boldsymbol{\beta}_k), \frac{\sigma_k^2}{w_i}\right), \\ W_i | Z_{ik} = 1 &\sim \text{Gamma}\left(\frac{\nu_k}{2}, \frac{\nu_k}{2}\right) \\ \mathbf{Z}_i | \mathbf{r}_i &\sim \text{Mult}(1; \pi_1(\mathbf{r}_i; \boldsymbol{\alpha}), \dots, \pi_K(\mathbf{r}_i; \boldsymbol{\alpha})). \end{aligned}$$

- This hierarchical representation involves the hidden variables  $Z_i$  and  $W_i$  facilitates the ML inference of model parameters  $\boldsymbol{\Psi}$  via E(C)M.

# MLE of the TMoE model

- Given an i.i.d sample of  $n$  observations,  $\Psi$  can be estimated by maximizing the observed-data log-likelihood :

$$\log L(\Psi) = \sum_{i=1}^n \log \sum_{k=1}^K \pi_k(\mathbf{r}_i; \boldsymbol{\alpha}) t\nu_k(y; \mu(\mathbf{x}; \boldsymbol{\beta}_k), \sigma_k^2, \nu_k) .$$

- $\Rightarrow$  EM algorithm and then describe an ECM extension
- The complete data consist of the responses  $(y_1, \dots, y_n)$  and their corresponding predictors  $(\mathbf{x}_1, \dots, \mathbf{x}_n)$  and  $(\mathbf{r}_1, \dots, \mathbf{r}_n)$ , as well as the latent variables  $(w_1, \dots, w_n)$  (in the hierarchical representation) and the latent labels  $(z_1, \dots, z_n)$ .

# MLE of the TMoE model

- $\Rightarrow$  The complete-data log-likelihood of  $\Psi$  is given by :

$$\log L_c(\Psi) = \log L_{1c}(\alpha) + \sum_{k=1}^K [\log L_{2c}(\Psi_k) + \log L_{3c}(\nu_k)],$$

where

$$\log L_{1c}(\alpha) = \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log \pi_k(\mathbf{r}_i; \alpha),$$

$$\log L_{2c}(\Psi_k) = \sum_{i=1}^n Z_{ik} \left[ -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma_k^2) - \frac{1}{2} w_i d_{ik}^2 \right],$$

$$\log L_{3c}(\nu_k) = \sum_{i=1}^n Z_{ik} \left[ -\log \Gamma\left(\frac{\nu_k}{2}\right) + \left(\frac{\nu_k}{2}\right) \log\left(\frac{\nu_k}{2}\right) + \left(\frac{\nu_k}{2} - 1\right) \log(w_i) - \left(\frac{\nu_k}{2}\right) w_i \right].$$

# MLE of the TMoE model : E-Step

**E-Step** Calculate the  $Q$ -function :

$$Q(\Psi; \Psi^{(m)}) = Q_1(\alpha; \Psi^{(m)}) + \sum_{k=1}^K \left[ Q_2(\theta_k, \Psi^{(m)}) + Q_3(\nu_k, \Psi^{(m)}) \right],$$

where  $\theta_k = (\beta_k^T, \sigma_k^2)^T$  and

$$Q_1(\alpha; \Psi^{(m)}) = \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}^{(m)} \log \pi_k(\mathbf{r}_i; \alpha),$$

$$Q_2(\theta_k; \Psi^{(m)}) = \sum_{i=1}^n \tau_{ik}^{(m)} \left[ -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma_k^2) - \frac{1}{2} w_{ik}^{(m)} d_{ik}^2 \right].$$

$$Q_3(\nu_k; \Psi^{(m)}) = \sum_{i=1}^n \tau_{ik}^{(m)} \left[ -\log \Gamma\left(\frac{\nu_k}{2}\right) + \left(\frac{\nu_k}{2}\right) \log\left(\frac{\nu_k}{2}\right) - \left(\frac{\nu_k}{2}\right) w_{ik}^{(m)} + \left(\frac{\nu_k}{2} - 1\right) e_{1,ik}^{(m)} \right]$$

→ requires the following conditional expectations (analytic) :

$$\begin{aligned} \tau_{ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [Z_{ik} | y_i, \mathbf{x}_i, \mathbf{r}_i], \\ w_{ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [W_i | y_i, Z_{ik} = 1, \mathbf{x}_i, \mathbf{r}_i], \\ e_{1,ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [\log(W_i) | y_i, Z_{ik} = 1, \mathbf{x}_i, \mathbf{r}_i]. \end{aligned}$$

# MLE of the TMoE model : M-Step

**M-Step 1** Calculate  $\alpha^{(m+1)}$  by maximizing  $Q_1(\alpha; \Psi^{(m)})$  w.r.t  $\alpha$ .  $\Rightarrow$  Iteratively via IRLS (92) as for the mixture of SNMoE.

**M-Step 2** Calculate  $\theta_k^{(m+1)}$  by maximizing  $Q_2(\theta_k; \Psi^{(m)})$  w.r.t  $\theta_k$

$$\beta_k^{(m+1)} = \left[ \sum_{i=1}^n \tau_{ik}^{(m)} w_{ik}^{(m)} \mathbf{x}_i \mathbf{x}_i^T \right]^{-1} \sum_{i=1}^n \tau_{ik}^{(m)} w_{ik}^{(m)} y_i \mathbf{x}_i,$$
$$\sigma_k^{2(m+1)} = \frac{1}{\sum_{i=1}^n \tau_{ik}^{(m)}} \sum_{i=1}^n \tau_{ik}^{(m)} w_{ik}^{(m)} \left( y_i - \beta_k^{T(m+1)} \mathbf{x}_i \right)^2.$$

**M-Step 3** Calculate  $\nu_k^{(m+1)}$  by maximizing  $Q_3(\nu_k; \Psi^{(m)})$  w.r.t  $\nu_k$

$\Rightarrow$  iteratively solve the following equation in  $\nu_k$  :

$$-\psi\left(\frac{\nu_k}{2}\right) + \log\left(\frac{\nu_k}{2}\right) + 1 + \frac{\sum_{i=1}^n \tau_{ik}^{(m)} (\log(w_{ik}^{(m)}) - w_{ik}^{(m)})}{\sum_{i=1}^n \tau_{ik}^{(m)}} + \psi\left(\frac{\nu_k^{(m)} + 1}{2}\right) - \log\left(\frac{\nu_k^{(m)} + 1}{2}\right) = 0.$$

This scalar non-linear equation can be solved with a root finding algorithm, such as Brent's method [Brent, 1973].



# The skew $t$ mixture of experts (STMoe) model

- A  $K$ -component mixture of skew  $t$  experts (STMoe) is defined by :

$$f(y|\mathbf{r}, \mathbf{x}; \Psi) = \sum_{k=1}^K \pi_k(\mathbf{r}; \alpha) \text{ST}(y; \mu(\mathbf{x}; \beta_k), \sigma_k^2, \lambda_k, \nu_k)$$

- $k$ th expert : has a skew  $t$  distribution [Azzalini and Capitanio, 2003] :

$$f(y|\mathbf{x}; \mu(\mathbf{x}; \beta_k), \sigma^2, \lambda, \nu) = \frac{2}{\sigma} t_{\nu}(d_y(\mathbf{x})) T_{\nu+1} \left( \lambda d_y(\mathbf{x}) \sqrt{\frac{\nu+1}{\nu + d_y^2(\mathbf{x})}} \right)$$

The skew  $t$  mixture of experts (STMoe) model extends the univariate skew  $t$  mixture model Lin et al. [2007], to the MoE framework.

## Model characteristics

$\hookrightarrow$  For  $\{\nu_k\} \rightarrow \infty$ , the STMoe reduces to the SNMoe

$\hookrightarrow$  For  $\{\lambda_k\} \rightarrow 0$ , the STMoe reduces to the TMoe.

$\hookrightarrow$  For  $\{\nu_k\} \rightarrow \infty$  and  $\{\lambda_k\} \rightarrow 0$ , it approaches the NMoe.

$\hookrightarrow$  The STMoe is flexible as it generalizes the previously described models

# Representation of the STMoE model

- **Stochastic representation** Suppose that conditional on a Multinomial categorical variable  $Z_i$ ,  $E_i$  and  $W_i$  are independent univariate random variables such that  $E_i \sim \text{SN}(\lambda_{z_i})$  and  $W_i \sim \text{Gamma}(\frac{\nu_{z_i}}{2}, \frac{\nu_{z_i}}{2})$ , and  $\mathbf{x}_i$  and  $\mathbf{r}_i$  are given covariates. A variable  $Y_i$  having the following representation :

$$Y_i = \mu(\mathbf{x}_i; \boldsymbol{\beta}_{z_i}) + \sigma_{z_i} \frac{E_i}{\sqrt{W_i}}$$

is said to follow the STMoE distribution

- **Hierarchical representation**

$$Y_i | u_i, w_i, Z_{ik} = 1, \mathbf{x}_i \sim \text{N}\left(\mu(\mathbf{x}_i; \boldsymbol{\beta}_k) + \delta_k |u_i|, \frac{1 - \delta_k^2}{w_i} \sigma_k^2\right),$$

$$U_i | w_i, Z_{ik} = 1 \sim \text{N}\left(0, \frac{\sigma_k^2}{w_i}\right),$$

$$W_i | Z_{ik} = 1 \sim \text{Gamma}\left(\frac{\nu_k}{2}, \frac{\nu_k}{2}\right)$$

$$\mathbf{Z}_i | \mathbf{r}_i \sim \text{Mult}(1; \pi_1(\mathbf{r}_i; \boldsymbol{\alpha}), \dots, \pi_K(\mathbf{r}_i; \boldsymbol{\alpha})).$$

The variables  $U_i$  and  $W_i$  are hidden in this hierarchical representation

# MLE via the ECM algorithm

- Maximize the observed-data log-likelihood :

$$\log L(\Psi) = \sum_{i=1}^n \log \sum_{k=1}^K \pi_k(\mathbf{r}_i; \alpha) \text{ST}(y; \mu(\mathbf{x}_i; \beta_k), \sigma_k^2, \lambda_k, \nu_k) .$$

- $\Rightarrow$  This is performed iteratively by a dedicated ECM algorithm.
- The complete-data log-likelihood :

$$\log L_c(\Psi) = \log L_{1c}(\alpha) + \sum_{k=1}^K [\log L_{2c}(\theta_k) + \log L_{3c}(\nu_k)]; \quad \theta_k = (\beta_k^T, \sigma_k^2, \lambda_k)^T$$

$$\log L_{1c}(\alpha) = \sum_{i=1}^n \sum_{k=1}^K Z_{ik} \log \pi_k(\mathbf{r}_i; \alpha),$$

$$\log L_{2c}(\theta_k) = \sum_{i=1}^n Z_{ik} \left[ -\log(2\pi) - \log(\sigma_k^2) - \frac{1}{2} \log(1 - \delta_k^2) - \frac{w_i d_{ik}^2}{2(1 - \delta_k^2)} + \frac{w_i u_i \delta_k d_{ik}}{(1 - \delta_k^2)\sigma_k} - \frac{w_i u_i^2}{2(1 - \delta_k^2)\sigma_k^2} \right]$$

$$\log L_{3c}(\nu_k) = \sum_{i=1}^n Z_{ik} \left[ -\log \Gamma\left(\frac{\nu_k}{2}\right) + \left(\frac{\nu_k}{2}\right) \log\left(\frac{\nu_k}{2}\right) + \left(\frac{\nu_k}{2}\right) \log(w_i) - \left(\frac{\nu_k}{2}\right) w_i \right] .$$

# MLE via the ECM algorithm : E-Step

- **E-Step** Calculates the  $Q$ -function, that is the conditional expectation of the complete-data log-likelihood , given the observed data  $\{y_i, \mathbf{x}_i, \mathbf{r}_i\}_{i=1}^n$  and a current parameter estimation  $\Psi^{(m)}$  given by :

$$Q(\Psi; \Psi^{(m)}) = Q_1(\alpha; \Psi^{(m)}) + \sum_{k=1}^K \left[ Q_2(\theta_k, \Psi^{(m)}) + Q_3(\nu_k, \Psi^{(m)}) \right],$$

where

$$Q_1(\alpha; \Psi^{(m)}) = \sum_{i=1}^n \sum_{k=1}^K \tau_{ik}^{(m)} \log \pi_k(\mathbf{r}_i; \alpha),$$

$$Q_2(\theta_k; \Psi^{(m)}) = \sum_{i=1}^n \tau_{ik}^{(m)} \left[ -\log(2\pi\sigma_k^2) - \frac{1}{2} \log(1 - \delta_k^2) - \frac{w_{ik}^{(m)} d_{ik}^2}{2(1 - \delta_k^2)} + \frac{\delta_k d_{ik} e_{1,ik}^{(m)}}{(1 - \delta_k^2)\sigma_k} - \frac{e_{2,ik}^{(m)}}{2(1 - \delta_k^2)\sigma_k^2} \right],$$

$$Q_3(\nu_k; \Psi^{(m)}) = \sum_{i=1}^n \tau_{ik}^{(m)} \left[ -\log \Gamma\left(\frac{\nu_k}{2}\right) + \left(\frac{\nu_k}{2}\right) \log\left(\frac{\nu_k}{2}\right) - \left(\frac{\nu_k}{2}\right) w_{ik}^{(m)} + \left(\frac{\nu_k}{2}\right) e_{3,ik}^{(m)} \right].$$

# MLE via the ECM algorithm : E-Step

- $\Rightarrow$  The E-Step requires the following conditional expectations :

$$\begin{aligned}\tau_{ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [Z_{ik} | y_i, \mathbf{x}_i, \mathbf{r}_i], \\ w_{ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [W_i | y_i, Z_{ik} = 1, \mathbf{x}_i, \mathbf{r}_i], \\ e_{1,ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [W_i U_i | y_i, Z_{ik} = 1, \mathbf{x}_i, \mathbf{r}_i], \\ e_{2,ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [W_i U_i^2 | y_i, Z_{ik} = 1, \mathbf{x}_i, \mathbf{r}_i], \\ e_{3,ik}^{(m)} &= \mathbb{E}_{\Psi^{(m)}} [\log(W_i) | y_i, Z_{ik} = 1, \mathbf{x}_i, \mathbf{r}_i].\end{aligned}$$

- These conditional expectations are calculated analytically except  $e_{3,ik}^{(m)}$  for which I adopted a one-step-late (OSL) approach as in Lee and McLachlan [2014], rather than using a Monte Carlo approximation as in Lin et al. [2007].
- I also mention that, for the multivariate skew  $t$  mixture models, recently Lee and McLachlan [2015] presented a series-based truncation approach, which exploits an exact representation of this conditional expectation and which can also be used here.

# MLE via the ECM algorithm : M-Step

- **CM-Step 1** update the mixing parameters  $\alpha^{(m+1)}$  by maximizing the function  $Q_1(\alpha; \Psi^{(m)})$  by using IRLS. Then, for  $k = 1 \dots, K$ ,
- **CM-Step 2** Update the regression params  $(\beta_k^{T(m+1)}, \sigma_k^{2(m+1)})$  :

$$\beta_k^{(m+1)} = \left[ \sum_{i=1}^n \tau_{ik}^{(q)} w_{ik}^{(m)} \mathbf{x}_i \mathbf{x}_i^T \right]^{-1} \sum_{i=1}^n \tau_{ik}^{(q)} \left( w_{ik}^{(m)} y_i - e_{1,ik}^{(m)} \delta_k^{(m+1)} \right) \mathbf{x}_i,$$

$$\sigma_k^{2(m+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(m)} \left[ w_{ik}^{(m)} \left( y_i - \beta_k^{T(m+1)} \mathbf{x}_i \right)^2 - 2 \delta_k^{(m+1)} e_{1,ik}^{(m)} (y_i - \beta_k^{T(m+1)} \mathbf{x}_i) + e_{2,ik}^{(m)} \right]}{2 \left( 1 - \delta_k^{2(m)} \right) \sum_{i=1}^n \tau_{ik}^{(m)}}$$

- **CM-Step 3** Update the skewness parameters  $\lambda_k$  by solving the following equation :

$$\delta_k (1 - \delta_k^2) \sum_{i=1}^n \tau_{ik}^{(m)} + (1 + \delta_k^2) \sum_{i=1}^n \tau_{ik}^{(m)} \frac{d_{ik}^{(m+1)} e_{1,ik}^{(m)}}{\sigma_k^{2(m+1)}} - \delta_k \sum_{i=1}^n \tau_{ik}^{(m)} \left[ w_{ik}^{(m)} d_{ik}^{2(m+1)} + \frac{e_{2,ik}^{(m)}}{\sigma_k^{2(m+1)}} \right] = 0.$$

- **CM-Step 4** Update the degree of freedom  $\nu_k$  by solving of the following equation :

$$-\psi \left( \frac{\nu_k}{2} \right) + \log \left( \frac{\nu_k}{2} \right) + 1 + \frac{\sum_{i=1}^n \tau_{ik}^{(m)} \left( e_{3,ik}^{(m)} - w_{ik}^{(m)} \right)}{\sum_{i=1}^n \tau_{ik}^{(m)}} = 0.$$

# Prediction, clustering and model selection

- **Prediction** Predicted response :  $\hat{y} = \mathbb{E}_{\hat{\Psi}}(Y|\mathbf{r}, \mathbf{x})$  with

$$\mathbb{E}_{\hat{\Psi}}(Y|\mathbf{r}, \mathbf{x}) = \sum_{k=1}^K \pi_k(\mathbf{r}; \hat{\alpha}_n) \mathbb{E}_{\hat{\Psi}}(Y|Z = k, \mathbf{x}),$$

$$\mathbb{V}_{\hat{\Psi}}(Y|\mathbf{r}, \mathbf{x}) = \sum_{k=1}^K \pi_k(\mathbf{r}; \hat{\alpha}_n) [(\mathbb{E}_{\hat{\Psi}}(Y|Z = k, \mathbf{x}))^2 + \mathbb{V}_{\hat{\Psi}}(Y|Z = k, \mathbf{x})] - [\mathbb{E}_{\hat{\Psi}}(Y|\mathbf{r}, \mathbf{x})]^2$$

where  $\mathbb{E}_{\hat{\Psi}}(Y|Z = k, \mathbf{x})$  and  $\mathbb{V}_{\hat{\Psi}}(Y|Z = k, \mathbf{x})$  are respectively the component-specific (expert) means and variances.

- **Clustering of regression data** Calculate the cluster label as

$$\hat{z}_i = \arg \max_{k=1}^K \mathbb{E}[Z_i|\mathbf{r}_i, \mathbf{x}_i; \hat{\Psi}] = \arg \max_{k=1}^K \frac{\pi_k(\mathbf{r}; \hat{\Psi}) f_k(y_i|\mathbf{r}_i, \mathbf{x}_i; \hat{\Psi})}{\sum_{k'=1}^K \pi_{k'}(\mathbf{r}; \hat{\alpha}) f_{k'}(y_i|\mathbf{r}_i, \mathbf{x}_i; \hat{\Psi}_{k'})}$$

- **Model selection** The value of  $(K, p)$  can be computed by using BIC, ICL

Number of free parameters :

$\eta_{\Psi} = K(p + 4) - 2$  for the NMoE model,

$\eta_{\Psi} = K(p + 5) - 2$  for both the SNMoE and the TMoE models,

$\eta_{\Psi} = K(p + 6) - 2$  for the STMoE model.

# Illustration on Bishop's data set

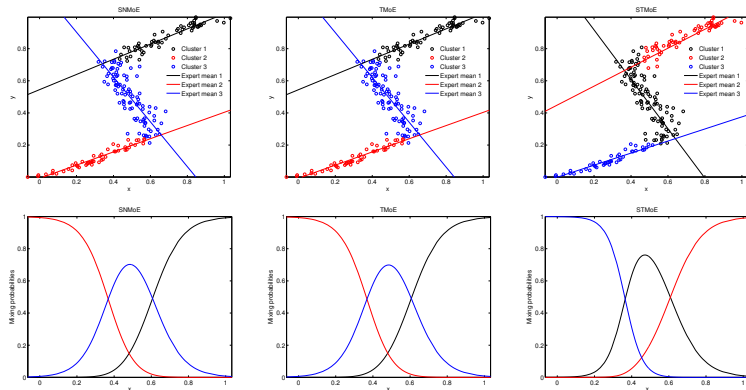


FIGURE – Fitting the the non-normal mixture of experts models (SNMoE, TNMoE, STMoE) to the toy data set analyzed in Bishop and Svensén [2003] :  $n = 250$  values of input variables  $x_i$  generated uniformly in  $(0, 1)$  and output variables  $y_i$  generated as  $y_i = x_i + 0.3 \sin(2\pi x_i) + \epsilon_i$ , with  $\epsilon_i$  drawn from a zero mean Normal distribution with standard deviation 0.05.



# Robustness of the NNMoe

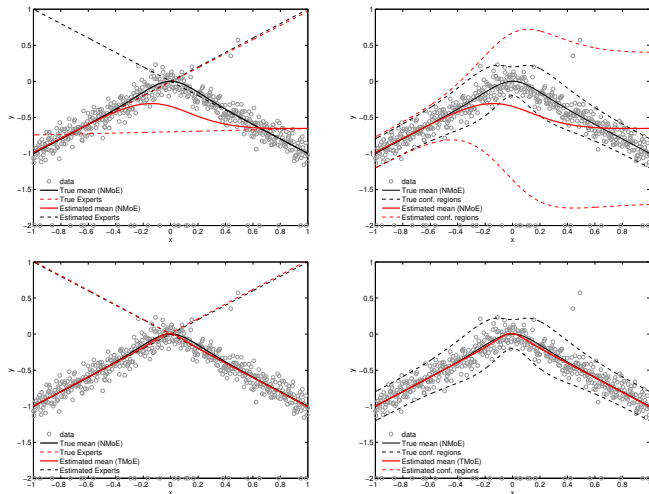


FIGURE – Fitted MoE to  $n = 500$  observations generated according to the NMoE with 5% of outliers ( $x; y = -2$ ) : NMoE fit (top), TMoE fit (bottom).

# Robustness of the NNMoe

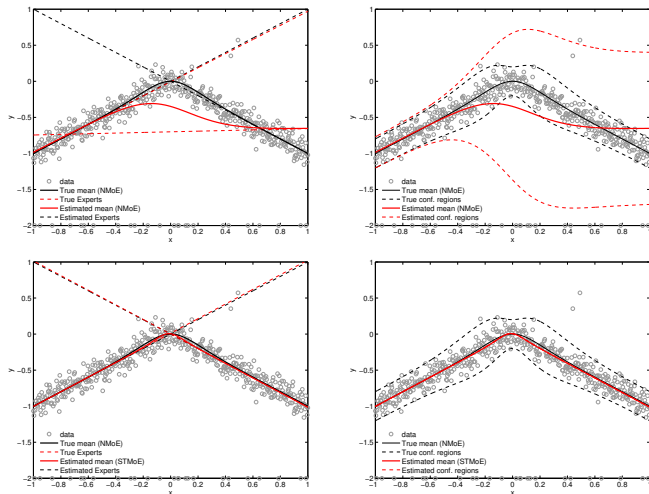


FIGURE – Fitted MoE to  $n = 500$  observations generated according to the NMoE with 5% of outliers ( $x; y = -2$ ) : NMoE fit (top), STMoE fit (bottom).

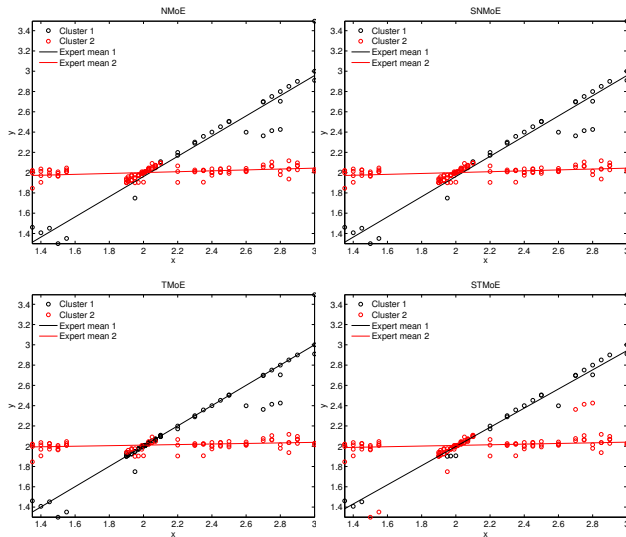


FIGURE – Fitting the MoE models to the tone data set studied by Bai et al. [2012] and Song et al. [2014] by using robust regression mixture models based on, respectively, the  $t$  distribution and the Laplace distribution :  $n = 150$  pairs of “tuned” predictors ( $x$ ), and their corresponding “strech ratio” responses ( $y$ ).

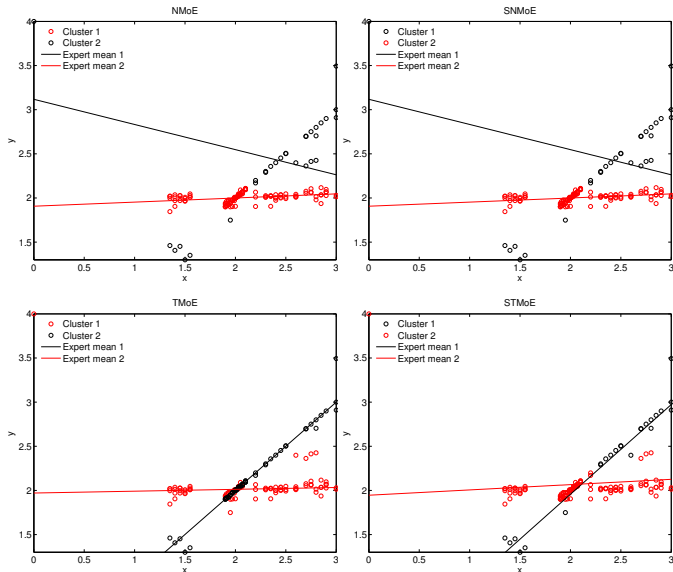


FIGURE – Fitting MoLE to the tone data set with ten added outliers  $(0, 4)$ .

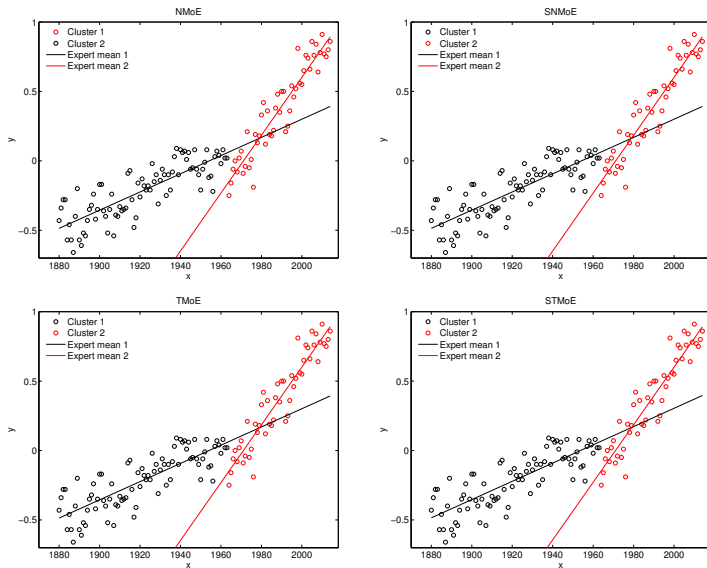
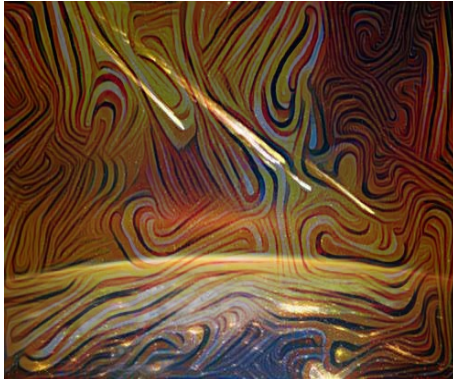


FIGURE – Fitting the MoLE models to the temperature anomalies data set.

# MEteorits : open-source soft. for mixtures-of-experts



**MEteorits** : Mixtures-of-ExperTs modEling for cOMplex and non-noRmal dIsTributionS

► R software

► Matlab software

# MEteorits : open-source soft. for mixtures-of-experts

## Available algorithms and Packages

- **NMoE** : Normal Mixture-of-Experts [▶ R software](#) [▶ Matlab software](#)
- **SNMoE** : Skew-Normal Mixture-of-Experts [▶ R software](#) [▶ Matlab software](#)
- **tMoE** : Robust modeling of mixture-of-experts using the  $t$ -distribution [▶ R software](#)  
[▶ Matlab software](#)
- **StMoE** : Skew- $t$  Mixture-of-Experts [▶ R software](#) [▶ Matlab software](#)

## High-dimensional Mixtures-of-Experts [Huynh and Chamroukhi, 2019]

Estimation and Feature Selection in Mixtures of Generalized Linear Experts Models

- **prEMME** : proximal Newton EM for estimation and feature selection in high-dimensional Mixtures-of-Experts [Huynh and Chamroukhi, 2019] [▶ R software](#)
- Expert models : Poisson [▶ R software](#) Logistic [▶ R software](#) Gaussian [▶ R software](#)

## Coming soon : MoE for functional data (functional predictors)

- **FunME** [▶ Matlab software](#)

# MoE models in high-dimension

## Maximum Likelihood Estimation via EM [Dempster et al., 1977, Jacobs et al., 1991]

- MLE :  $\Psi$  is commonly estimated by maximizing the observed-data log-likelihood :

$$\hat{\Psi}_n \in \arg \max_{\Psi \in \Theta} L(\Psi) \text{ with } L(\Psi) = \sum_{i=1}^n \log \sum_{k=1}^K \pi_k(\mathbf{x}_i; \mathbf{w}) f(\mathbf{y}_i | \mathbf{x}_i; \Psi_k)$$

## Regularized MLE of the MoE [Khalili, 2010] [Huyhn and Chamroukhi, 2019]

$\Psi$  is estimated by maximizing a penalized observed-data log-likelihood :

$$\hat{\Psi}_n \in \arg \max_{\Psi \in \Theta} L(\Psi) - \text{Pen}_{\lambda}(\Psi)$$

- $\hookrightarrow \text{Pen}_{\lambda}(\Psi)$  LASSO penalties for experts and the gating network
- encourages sparse solutions
- parameter estimation and selection problem

## High-dimensional Mixtures-of-Experts

Estimation and Feature Selection in Mixtures of Generalized Linear Experts Models

- **prEMME** : proximal Newton EM for estimation and feature selection in high-dimensional Mixtures-of-Experts [▶ R software](#)
- Poisson [▶ R software](#) Logistic [▶ R software](#) Gaussian [▶ R software](#)



# FunME : Functional Mixtures-of-Experts

Ongoing work : MoE with functional predictors/responses

- Let  $\{X_i(\cdot), Y_i\}_{i=1}^n$ , be a random i.i.d sample where  $Y_i \in \mathbb{R}$  is the response and  $X_i(t); t \in \mathcal{T} \subset \mathbb{R}$  is a functional predictor, for example the time in time series.
- The input  $X(\cdot)$  is a function (eg. data continuously recorded for some time period)  
eg.  $\mathbf{X}(\cdot)$  are data continuously recorded from multiple subject' sensors

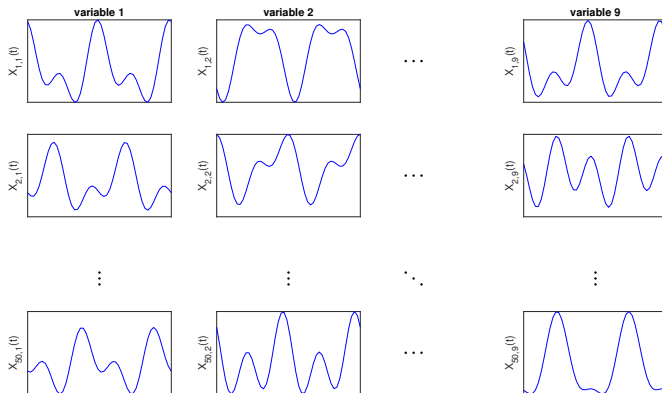


FIGURE – Functional predictors  $X_{ij}(t)$   $t \in \mathcal{T}$ ,  $i = 1, \dots, n$  and  $j = 1, \dots, p$ .

# References I

- H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6) :716–723, 1974.
- A. Azzalini. A class of distributions which includes the normal ones. *Scandinavian Journal of Statistics*, pages 171–178, 1985.
- A. Azzalini. Further results on a class of distributions which includes the normal ones. *Scandinavian Journal of Statistics*, pages 199–208, 1986.
- A. Azzalini and A. Capitanio. Distributions generated by perturbation of symmetry with emphasis on a multivariate skew  $t$  distribution. *Journal of the Royal Statistical Society, Series B*, 65 :367–389, 2003.
- Xiuqin Bai, Weixin Yao, and John E. Boyer. Robust fitting of mixture regression models. *Computational Statistics & Data Analysis*, 56(7) :2347 – 2359, 2012.
- Jeffrey D. Banfield and Adrian E. Raftery. Model-based Gaussian and non-Gaussian clustering. *Biometrics*, 49(3) :803–821, 1993.
- C. Biernacki, G. Celeux, and G. Govaert. Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(7) :719–725, 2000.
- C. Bishop and M. Svensén. Bayesian hierarchical mixtures of experts. In *In Uncertainty in Artificial Intelligence*, 2003.
- 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.
- Richard P. Brent. *Algorithms for minimization without derivatives*. Prentice-Hall series in automatic computation. Englewood Cliffs, N.J. Prentice-Hall, 1973. ISBN 0-13-022335-2.
- G. Celeux and G. Govaert. Gaussian Parsimonious Clustering Models. *Pattern Recognition*, 28(5) :781–793, 1995.
- F. Chamroukhi. Skew-normal mixture of experts. 2016a. URL <https://chamroukhi.com/papers/Chamroukhi-SNMoe.pdf>.
- F. Chamroukhi. Robust mixture of experts modeling using the  $t$ -distribution. *Neural Networks - Elsevier*, 79 :20–36, 2016b. URL <https://chamroukhi.com/papers/TMoE.pdf>.

# References II

- F. Chamroukhi. Skew  $t$  mixture of experts. *Neurocomputing - Elsevier*, 266 :390–408, 2017. URL <https://chamroukhi.com/papers/STMoE.pdf>.
- F. Chamroukhi, A. Samé, G. Govaert, and P. Akin. A regression model with a hidden logistic process for feature extraction from time series. In *International Joint Conference on Neural Networks (IJCNN)*, 2009.
- Faïcel Chamroukhi and Hien D. Nguyen. Model-based clustering and classification of functional data. *Wiley Interdisciplinary Reviews : Data Mining and Knowledge Discovery*, Dec 2018. URL <https://chamroukhi.com/papers/MBCC-FDA.pdf>. DOI : 10.1002/widm.1298.
- A. Delaigle, P. Hall, and N. Bathia. Componentwise classification and clustering of functional data. *Biometrika*, 99(2) :299–313, 2012.
- 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.
- F. Ferraty and P. Vieu. Curves discrimination : a nonparametric functional approach. *Computational Statistics & Data Analysis*, 44(1-2) :161–173, 2003.
- F. Ferraty and P. Vieu. *Nonparametric functional data analysis : theory and practice*. Springer series in statistics, 2006. ISBN 0-387-30369-3.
- G. Hébrail, B. Huguency, Y. Lechevallier, and F. Rossi. Exploratory analysis of functional data via clustering and optimal segmentation. *Neurocomputing*, 73(7-9) :1125–1141, March 2010.
- T. Huynh and F. Chamroukhi. Estimation and feature selection in mixtures of generalized linear experts models. *arXiv :1810.12161*, July 2019. URL <https://chamroukhi.com/papers/preMME.pdf>.
- Salvatore Ingrassia, Simona Minotti, and Giorgio Vittadini. Local statistical modeling via a cluster-weighted approach with elliptical distributions. *Journal of Classification*, 29(3) :363–401, 2012.
- R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton. Adaptive mixtures of local experts. *Neural Computation*, 3(1) : 79–87, 1991.

# 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.
- M. I. Jordan and R. A. Jacobs. Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, 6 :181–214, 1994.
- A. Khalili. New estimation and feature selection methods in mixture-of-experts models. *Canadian Journal of Statistics*, 38(4) : 519–539, 2010.
- Sharon X. Lee and Geoffrey J. McLachlan. Finite mixtures of multivariate skew  $t$ -distributions : some recent and new results. *Statistics and Computing*, 24(2) :181–202, 2014.
- Sharon X. Lee and Geoffrey J. McLachlan. Finite mixtures of canonical fundamental skew  $t$ -distributions. *Statistics and Computing (To appear)*, 2015. doi: 10.1007/s11222-015-9545-x.
- Tsung I. Lin, Jack C. Lee, and Wan J. Hsieh. Robust mixture modeling using the skew  $t$  distribution. *Statistics and Computing*, 17(2) :81–92, 2007.
- 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. The classification and mixture maximum likelihood approaches to cluster analysis. In P.R. Krishnaiah and L. Kanal, editors, *In Handbook of Statistics, Vol. 2*, pages 199–208. Amsterdam : North-Holland, 1982.
- G. J. McLachlan and T. Krishnan. *The EM algorithm and extensions*. New York : Wiley, second edition, 2008.
- G. J. McLachlan and D. Peel. *Finite mixture models*. New York : Wiley, 2000.
- G. J. McLachlan, D. Peel, and R. W. Bean. Modelling high-dimensional data by mixtures of factor analyzers. *Computational Statistics & Data Analysis*, 41(3-4) :379–388, January 2003. URL <http://www.sciencedirect.com/science/article/B6V8V-472JRC1-12/1/4d40244841bb6f7c8c454ca92e6cc347>.

# References IV

- Geoffrey J. McLachlan and David Peel. Robust cluster analysis via mixtures of multivariate t-distributions. In *Lecture Notes in Computer Science*, pages 658–666. Springer-Verlag, 1998.
- G.J. McLachlan and K.E. Basford. *Mixture Models : Inference and Applications to Clustering*. Marcel Dekker, New York, 1988.
- X. L. Meng and D. B. Rubin. Maximum likelihood estimation via the ECM algorithm : A general framework. *Biometrika*, 80(2) : 267–278, 1993.
- Hien D. Nguyen and Faïcel Chamroukhi. Practical and theoretical aspects of mixture-of-experts modeling : An overview. *Wiley Interdisciplinary Reviews : Data Mining and Knowledge Discovery*, pages e1246–n/a, Feb 2018. ISSN 1942-4795. doi: 10.1002/widm.1246. URL <http://dx.doi.org/10.1002/widm.1246>.
- D. Peel and G. J. McLachlan. Robust mixture modelling using the t distribution. *Statistics and Computing*, 10 :339–348, 2000.
- Xinghao Qiao, Shaojun Guo, and Gareth M. James. Functional graphical models. *Journal of the American Statistical Association*, 0(0) :1–12, 2018. doi: 10.1080/01621459.2017.1390466.
- J. O. Ramsay and B. W. Silverman. *Functional Data Analysis*. Springer Series in Statistics. Springer, June 2005.
- J.O. Ramsay, T.O. Ramsay, and L.M. Sangalli. Spatial functional data analysis. In F. Ferraty, editor, *Recent Advances in Functional Data Analysis and Related Topics*, pages 269–275. Springer, 2011.
- 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.
- G. Schwarz. Estimating the dimension of a model. *Annals of Statistics*, 6 :461–464, 1978.
- Weixing Song, Weixin Yao, and Yanru Xing. Robust mixture regression model fitting by laplace distribution. *Computational Statistics & Data Analysis*, 71(0) :128 – 137, 2014.
- Y. Wei. Robust mixture regression models using t-distribution. Technical report, Master Report, Department of Statistics, Kansas State University, 2012.

Thank you for your attention !