

Learning mixtures-of-experts from heterogenous and high-dimensional data

FAÏCEL CHAMROUKHI



CMStatistics, December 17, 2023

16th International Conference of the ERCIM WG on
Computational and Methodological Statistics (CMStatistics 2023)

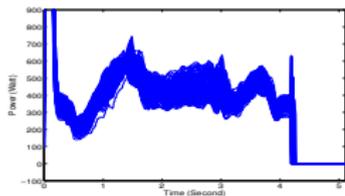
16-18 December 2023, HTW Berlin, University of Applied Sciences, Berlin, Germany



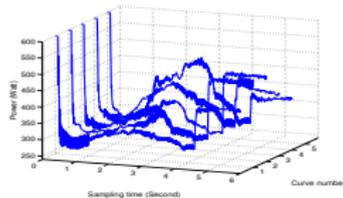
Real-world data are complex

- Heterogenous, Multimodal, High-Dimensional, Unlabeled, Possibly Massive ...
- Need for adapted analysis tools

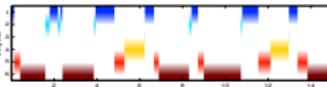
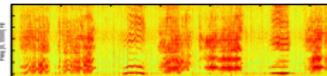
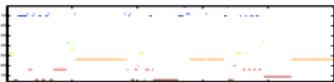
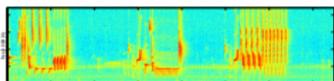
Transport : Railway switch curves diagnostic



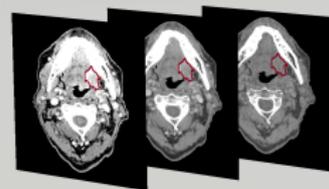
Predictive Maintenance



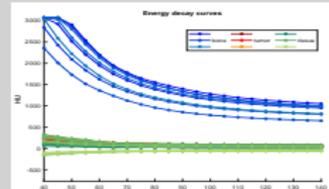
Acoustics : scene listening (marine, terrestrial)



Health : Medical images



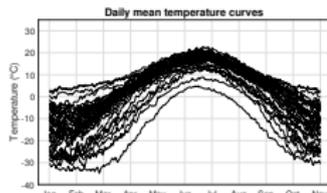
Dual-energy computed tomography



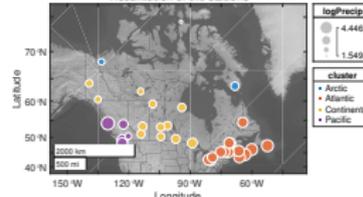
Health & Well Being : Activity recog.



Climate/Environment : meteorological data



Visualization of the stations



- Establish well-principled (with statistical guarantees) predictions in heterogeneous and high-dimensional situations,
- Construct efficient algorithms that operate in unsupervised way and provide interpretable solutions with computational guarantees.

Scientific framework

↪ **Latent variable models** : $f(x|\theta) = \int_{\mathcal{Z}} f(x, z|\theta) dz$

↪ **Inference, unsupervised representation and Selection** in high-dimension

1 Scientific Challenges

2 Latent Variable Models

- Mixture models
- Mixtures of Experts Models

3 High-Dimensional Learning

- Learning with high-dimensional predictors
- Learning with functional predictors

Density approximation in Unsupervised Learning

- **Data** : observations $\{\mathbf{x}_i\}$ from $\mathbf{X} \in \mathbb{X} \subset \mathbb{R}^d$ of density (multimodal) $f \in \mathcal{F}$
- **Objective** : approximate the density f (and represent the data, e.g. *clustering*)
- **Solution** : Approximate f within the class $\mathcal{H}^\varphi = \bigcup_{K \in \mathbb{N}^*} \mathcal{H}_K^\varphi$ of finite location-scale mixture h_K^φ (of K -components) of density φ (e.g., Gaussian), where

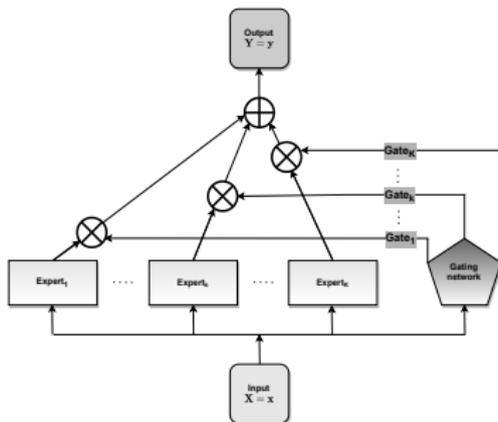
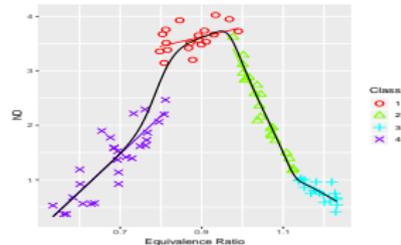
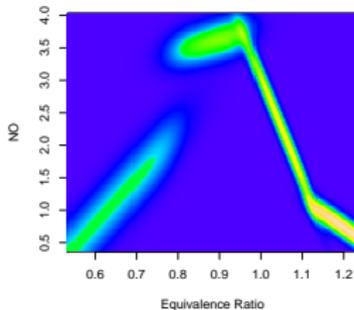
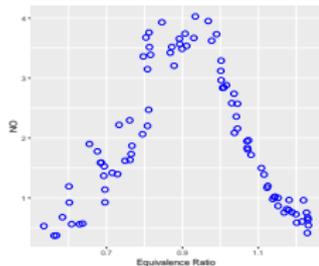
$$\mathcal{H}_K^\varphi = \left\{ h_K^\varphi(\mathbf{x}) := \sum_{k=1}^K \pi_k \frac{1}{\sigma_k^d} \varphi\left(\frac{\mathbf{x} - \boldsymbol{\mu}_k}{\sigma_k}\right), \boldsymbol{\mu}_k \in \mathbb{R}^d, \sigma_k \in \mathbb{R}_+, \pi_k > 0 \forall k \in [K], \sum_{k=1}^K \pi_k = 1 \right\}$$

Theorem : Universal approximation of finite location-scale mixtures

- Given any p.d.f $f, \varphi \in \mathcal{C}$ and a compact set $\mathbb{X} \subset \mathbb{R}^d$, there exists a sequence $(h_K^\varphi) \subset \mathcal{H}^\varphi$, such that $\lim_{K \rightarrow \infty} \sup_{\mathbf{x} \in \mathbb{X}} |f(\mathbf{x}) - h_K^\varphi(\mathbf{x})| = 0$.
- For $p \in [1, \infty)$, if $f \in \mathcal{L}_p$ (Lebesgue p.d.f) and $\varphi \in \mathcal{L}_\infty$ (essentially bounded p.d.f), there exists a sequence $(h_K^\varphi) \subset \mathcal{H}^\varphi$, such that $\lim_{K \rightarrow \infty} \|f - h_K^\varphi\|_{\mathcal{L}_p} = 0$.

Heterogeneous regression-type data

Mixtures-of-Experts as good candidates to model a response Y given covariates X when governed by a hidden structure accounting for heterogeneity



- **Context** : n observations $\{\mathbf{x}_i, \mathbf{y}_i\}$ from a pair $(\mathbf{X}, \mathbf{Y}) \in \mathbb{X} \times \mathbb{Y}$ with unknown conditional p.d.f $f \in \mathcal{F} = \{f : \mathbb{X} \times \mathbb{Y} \rightarrow \mathbb{R}_+ \mid \int_{\mathbb{Y}} f(\mathbf{y}|\mathbf{x}) d\lambda(\mathbf{y}) = 1, \forall \mathbf{x} \in \mathbb{X}\}$
- **High-dimensional setting** : $\mathbb{X} \subseteq \mathbb{R}^d, \mathbb{Y} \subseteq \mathbb{R}^q$, with $d, q \gg n$ and **heterogeneous**.
- **Objectifs** : Regression ; Clustering ; Model selection
- **Solution** : Approximate f within the class of **mixtures-of-experts** :

Let φ be a p.d.f (compactly supported on $\mathbb{Y} \subseteq \mathbb{R}^q$), we define the functional classes :

- Location-scale family : $\mathcal{E}_\varphi = \left\{ \phi_q(\mathbf{y}; \boldsymbol{\mu}, \sigma) := \frac{1}{\sigma^q} \varphi\left(\frac{\mathbf{y}-\boldsymbol{\mu}}{\sigma}\right); \boldsymbol{\mu} \in \mathbb{Y}, \sigma \in \mathbb{R}_+ \right\}$.
- Mixture of location-scale experts with softmax activation network : SGaME :

$$\mathcal{H}_S^\varphi = \left\{ h_K^\varphi(\mathbf{y}|\mathbf{x}) := \sum_{k=1}^K g_k(\mathbf{x}; \boldsymbol{\gamma}) \phi_q(\mathbf{y}; \boldsymbol{\mu}_k, \sigma_k) ; \phi_q \in \mathcal{E}_\varphi \cap \mathcal{L}_\infty, g_k(\cdot; \boldsymbol{\gamma}) \in \{\text{softmax}\} \right\}$$

Theorem : Approximation capabilities of isotropic mixtures-of-experts SGaME

- For $p \in [1, \infty)$, $f \in \mathcal{F}_p \cap \mathcal{C}$, $\varphi \in \mathcal{F} \cap \mathcal{C}$, $\mathbb{X} = [0, 1]^d$, there exists a sequence $(h_K^\varphi) \subset \mathcal{H}_S^\varphi$ such that $\lim_{K \rightarrow \infty} \|f - h_K^\varphi\|_{\mathcal{L}_p} = 0$.
- For $f \in \mathcal{F} \cap \mathcal{C}$, if $\varphi \in \mathcal{F} \cap \mathcal{C}$, $d = 1$, there exists a sequence $(h_K^\varphi) \subset \mathcal{H}_S^\varphi$ such that $\lim_{K \rightarrow \infty} h_K^\varphi = f$ almost uniformly.

Learning via the EM algorithm

SaMuraiS : open source software for statistical time-series analysis



SaMuraiS : StAtistical Models for the UnsupERvised segmentAtion of time-Series

▶ [Github](#)

▶ [CRAN](#)

▶ [Matlab software](#)

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](#)

MPWR : Multivariate PWR

▶ [R software](#)

▶ [Matlab software](#)

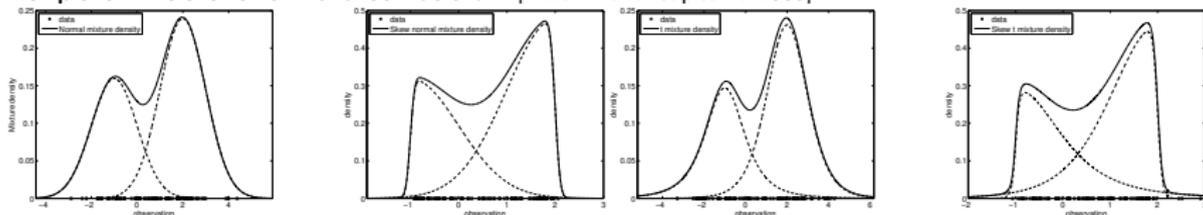
Include estimation, segmentation, approximation, model selection, and sampling

Principled robustness in regression and clustering

- **Questionings** : Prediction (non-linear regr., classification) & clustering in presence of **Outliers**, with potentially **skewed**, **heavy-tailed** distributions
- **Answering** : Robust MoE that accommodate asymmetry, heavy tails, and outliers

$$m(y|\mathbf{r}, \mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^K \underbrace{g_k(\mathbf{r}; \boldsymbol{\alpha})}_{\text{Softmax Gating Network}} \underbrace{ST(y; \mu(\mathbf{x}; \boldsymbol{\beta}_k), \sigma_k, \lambda_k, \nu_k)}_{\text{Skew-}t \text{ Expert Network}}$$

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

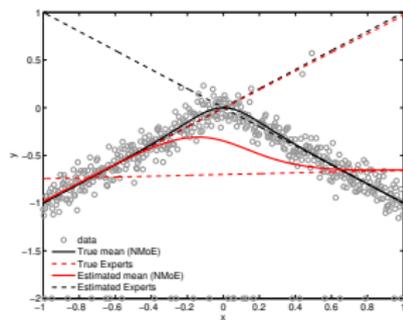


$$\pi_k = [0.4, 0.6], \mu_k = [-1, 2]; \sigma_k = [1, 1]; \nu_k = [3, 7]; \lambda_k = [14, -12];$$

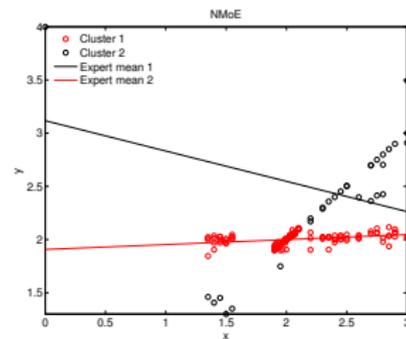
Flexible and robust generalization of the standard MoE models

For $\{\nu_k\} \rightarrow \infty$, STMoE reduces to SNMoE; For $\{\lambda_k\} \rightarrow 0$, STMoE reduces to TMoE.
 For $\{\nu_k\} \rightarrow \infty$ and $\{\lambda_k\} \rightarrow 0$, StMoE approaches the NMoE.

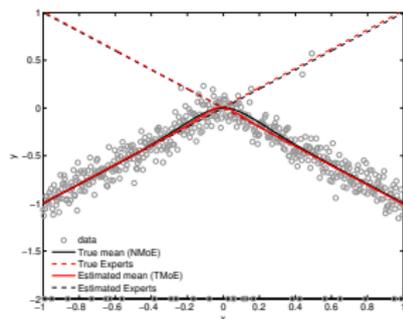
Robust learning with mixtures-of-experts models



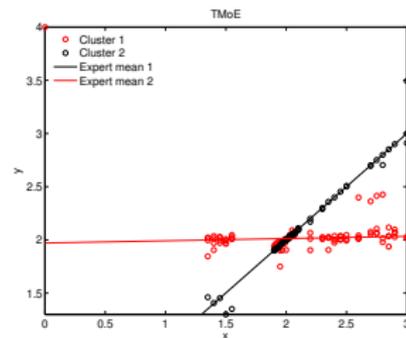
$n = 500$ observations with 5% of outliers ($x; y = -2$) : Normal fit



Tone data with 10 outliers (0, 4) : Normal fit



$n = 500$ observations with 5% of outliers ($x; y = -2$) : Robust fit



Tone data with 10 outliers (0, 4) : Robust fit

MEteorits : open-source soft. Robust learning with mixtures-of-experts models



MEteorits : **M**ixtures-of-**E**xper**T**s **m**o**D**ELing for **c**o**M**plex and non-**n**o**R**mal **d**IS**T**ribution**S**

[▶ Github](#) [▶ CRAN](#) [▶ Matlab software](#)

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 MoE using the t -distribution

[▶ R software](#)

[▶ Matlab software](#)

StMoE : Skew- t Mixture-of-Experts

[▶ R software](#)

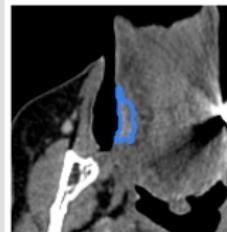
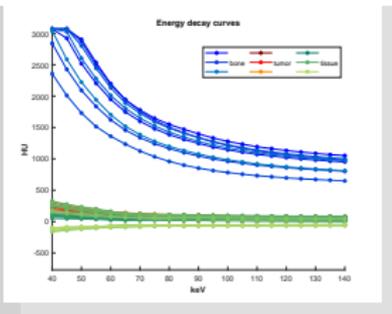
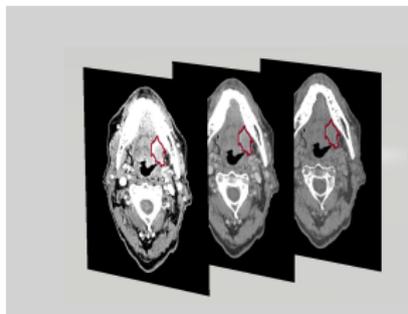
[▶ Matlab software](#)

- Meteorits include sampling, fitting, prediction, clustering with each MoE model
- Non-normal mixtures (and MoE) is a very recent topic in the field

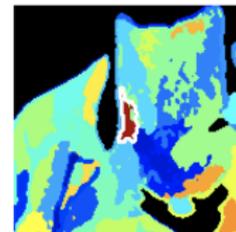
- Learning from Multimodal information in Healthcare/Radiology
- Cancer detection in Radiology : DECT clustering [Diagnostics (AI in medicine), 2022]

Spatial mixture of functional regressions for dual-energy CT images

$$m(\mathbf{y}|\mathbf{x}, \mathbf{v}; \boldsymbol{\theta}) = \sum_{k=1}^K \alpha_k(\mathbf{v}; \boldsymbol{\alpha}) f_k(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta}_k) \text{ where } \alpha_k(\mathbf{v}; \boldsymbol{\alpha}) = \frac{w_k \phi_3(\mathbf{v}; \boldsymbol{\mu}_k, \mathbf{R}_k)}{\sum_{\ell=1}^K w_\ell \phi_3(\mathbf{v}; \boldsymbol{\mu}_\ell, \mathbf{R}_\ell)}$$

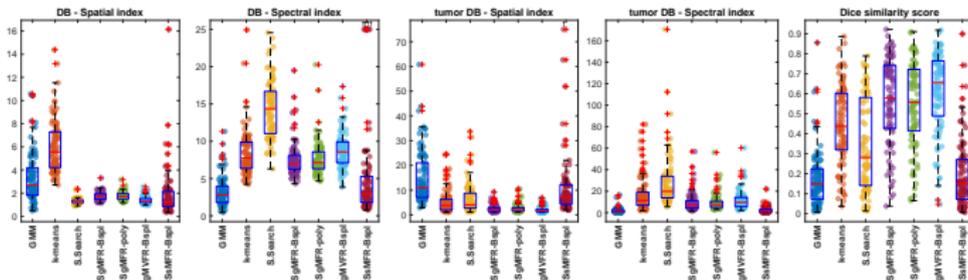


(a) Original slice



(b) Dice=0.84, DB=1.64/6.92

DECT multimodal Data : 3D voxels & energy levels Expert Annotation Automatic Annotation



Questioning : Prediction (non-linear regr., classification) & clustering in presence of

[1.] **High-dimensional** predictors : $\mathbf{X}_i \in \mathbb{R}^p$ with $p \gg n$

[2.] **Functional** predictors : $X_i(t), t \in \mathcal{T} \subseteq \mathbb{R}$ {eg. continuously recorded variables}

↔ Look for parsimonious and interpretable methods

[1.] HDME : High-Dimensional Mixtures-of-Experts

■ Learning : PMLE $\hat{\theta}_n \in \arg \max_{\theta} \sum_{i=1}^n \log h_K^{\varphi}(\mathbf{y}_i | \mathbf{x}_i; \theta) - \text{pen}(\theta)$

■ ↔ LASSO penalty : $\text{Pen}_{\lambda}(\theta) = \underbrace{\sum_{k=1}^K \lambda_k \|\beta_k\|_1}_{\text{Experts Net.}} + \underbrace{\sum_{k=1}^{K-1} \gamma_k \|\mathbf{w}_k\|_1}_{\text{Gating Net.}}$

↔ encourages sparse solutions & performs estimation and feature selection

↔ computationally attractive (Avoid matrix inversion ; univariate updates)

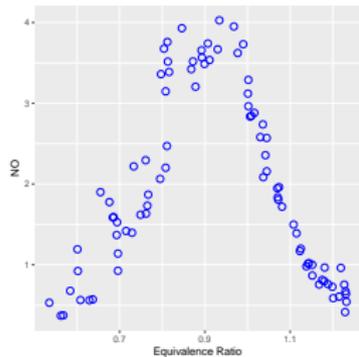
▶ [Software Toolbox HDME on Github](#) (GaussRMOE, LogisticRMOE, PoissonRMOE)

[PhD] Bao Tuyen Huynh. *Estimation and Feature Selection in High-Dimensional Mixtures-of-Experts Models*. PhD Thesis, Normandie Université, 2019.

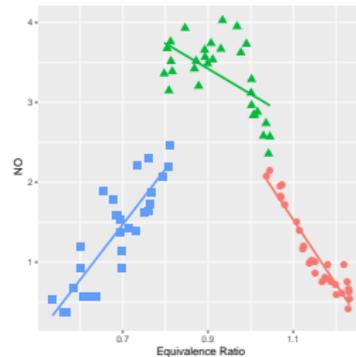
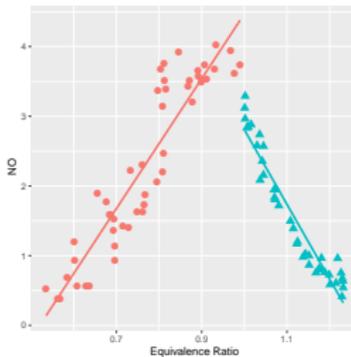
[J] Chamroukhi & Huynh. Regularized Maximum Likelihood Estimation and Feature Selection in Mixtures-of-Experts Models. *Journal de la Société Française de Statistique*, Vol. 160(1), pp :57–85, 2019

[J] Huynh & C. Estimation and Feature Selection in Mixtures of Generalized Linear Experts Models. arXiv :1810.12161, 2019

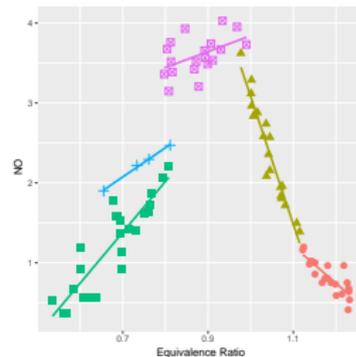
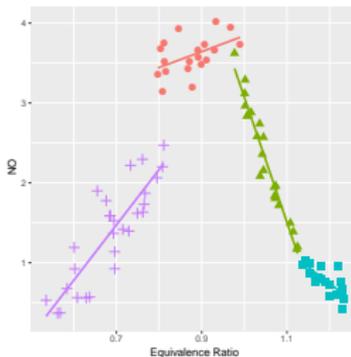
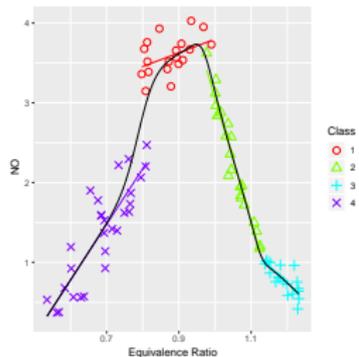
(a) Raw Ethanol data set



Collection of MoE models with linear mean functions characterized by 2-5 clusters



(b) Our best data-driven MoE model



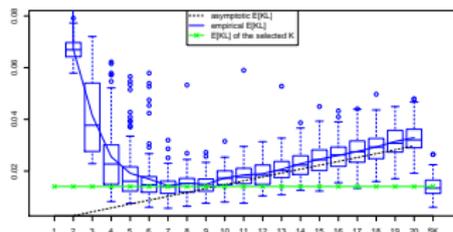
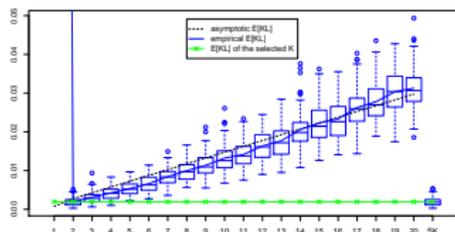
Questioning : Prediction (non-linear regr., classification) & clustering in presence of **High-dimensional** predictors : Data $\mathcal{D}_n = (\mathbf{X}_i, Y_i)_{i=1}^n$ where $\mathbf{X}_i \in \mathbb{R}^p$ with $p \gg n$
 HDME : High-Dimensional MoE : PMLE $\hat{\theta}_n \in \arg \max_{\theta} \sum_{i=1}^n \log h_K^{\varphi}(\mathbf{y}_i | \mathbf{x}_i; \theta) - \text{pen}(\theta)$

Theorem : Non-asymptotic oracle inequality for collection of MoE models

Result : \exists constants C et $\kappa(\rho, C_1) > 0$ ($C_1 > 1$) s.that whenever for $\mathbf{m} \in \mathcal{M}$, $\text{pen}(\mathbf{m}) \geq \kappa(\rho, C_1) [(C + \ln n) \dim(\mathcal{H}_{\mathbf{m}}) + z_{\mathbf{m}}]$, the estimator PMLE $\hat{h}_{\widehat{\mathbf{m}}}$ satisfies

$$\mathbb{E} \left[\text{JKL}_{\rho}^{\otimes n} \left(f, \hat{h}_{\widehat{\mathbf{m}}} \right) \right] \leq C_1 \inf_{\mathbf{m} \in \mathcal{M}} \left(\inf_{h_{\mathbf{m}} \in \mathcal{H}_{\mathbf{m}}} \text{KL}^{\otimes n} (f, h_{\mathbf{m}}) + \frac{\text{pen}(\mathbf{m})}{n} \right) + \frac{\kappa(\rho, C_1) C_1 \xi}{n} + \frac{\eta + \eta'}{n}$$

- A non-asymptotic result. If $\text{pen}(\mathbf{m})$ is well chosen, then our PMLE behaves in a comparable manner compared to **the best (oracle) model** $\mathcal{H}_{\mathbf{m}^*}$ in the collection, minimizing the risk : $\inf_{\mathbf{m} \in \mathcal{M}} \left(\inf_{h_{\mathbf{m}} \in \mathcal{H}_{\mathbf{m}}} \text{KL}^{\otimes n} (f, h_{\mathbf{m}}) + \frac{\text{pen}(\mathbf{m})}{n} \right)$ (f is unknown).



FLaMingoS : open source software for learning from functions



FLaMingoS : Functional Latent datA Models for clusterING heterogeneous time-Series

▶ [Github](#)

▶ [CRAN](#)

▶ [Matlab software](#)

Available algorithms and Packages

mixRHLP : Mixture of Regressions with HLPs

▶ [R](#)

▶ [Matlab](#)

mixHMM : Mixture of Hidden Markov Models (HMMs)

▶ [R](#)

▶ [Matlab](#)

mixHMMR : Mixture of HMM Regressions

▶ [R](#)

▶ [Matlab](#)

PWRM : Piece-Wise Regression Mixture

▶ [R](#)

▶ [Matlab](#)

uReMix : Unsupervised Regression Mixtures

▶ [R](#)

▶ [Matlab](#)

↔ A flexible full generative modeling for FDA

↔ Could be extended to the multivariate case without a major effort

[2.] Learning with functional predictors

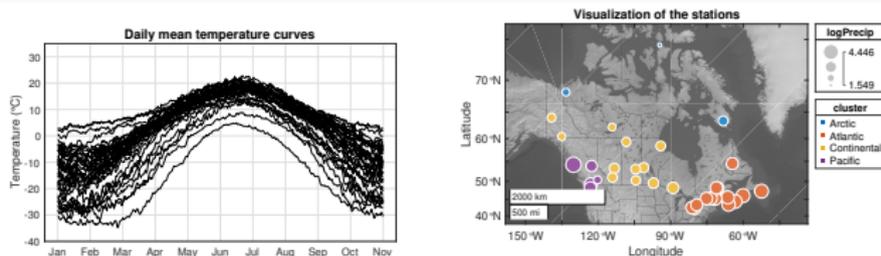


FIGURE – $n = 35$ daily mean temperature measurement curves (X_i 's) in different stations (Left) and the log of precipitation values (Y_i 's) visualized with the climate regions (Z_i 's) (Right).

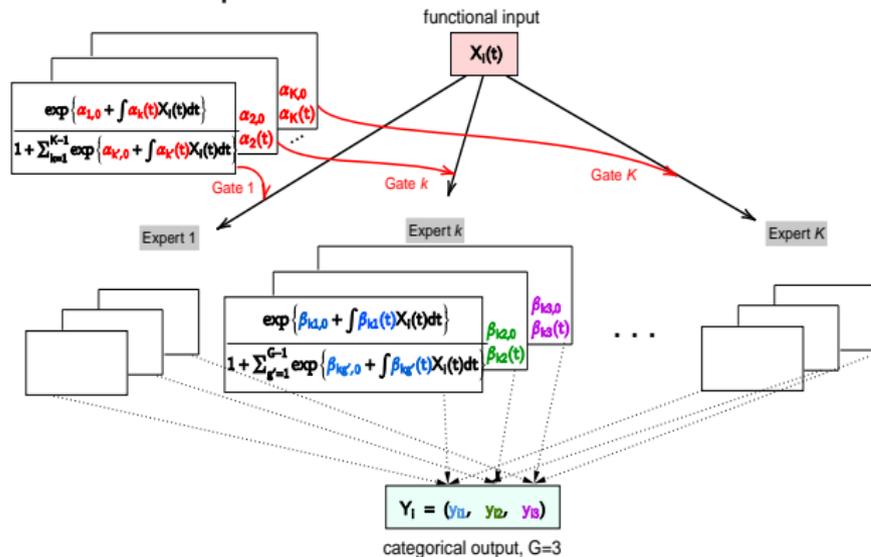
- Relate functional predictors $\{X(t) \in \mathbb{R}; t \in \mathcal{T} \subset \mathbb{R}\}$ to a scalar response $Y \in \mathcal{Y} \subset \mathbb{R}$
- Regression and classification of heterogeneous responses given functional predictors
 - (1) generative functional modeling, sparsity and feature selection (high-dimension)
 - (2) User guideline : keep an interpretable fit

[2.] Functional Mixtures-of-Experts (and Different Learning strategies, in particular)

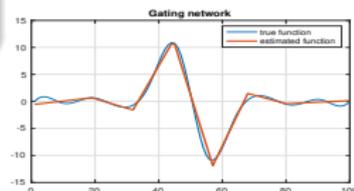
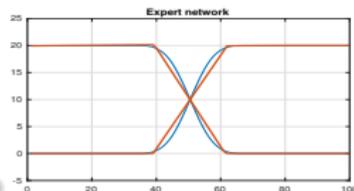
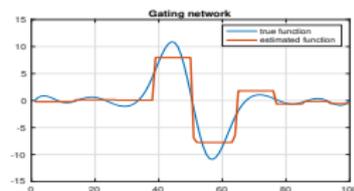
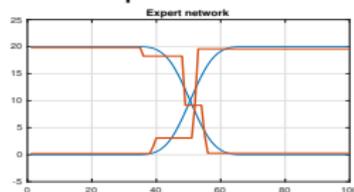
- $Y_i = \beta_{z_i,0} + \int_{\mathcal{T}} X_i(t) \beta_{z_i}(t) dt + \varepsilon_i$ avec $h_z(X_i(\cdot)) = \alpha_{z_i,0} + \int_{\mathcal{T}} X_i(t) \alpha_{z_i}(t) dt$
- Lasso-type Regularized MLE w.r.t the derivatives of the $\alpha(\cdot)$ and $\beta(\cdot)$ functions

[Preprint] Chamroukhi, Pham, Hoang, McLachlan. Functional Mixtures-of-Experts. arXiv :2202.02249, Feb, 2022 (Under Review, Statistics and Computing)

Mixture-of-Experts Architecture



Interpretable fits



$Y_i = \beta_{z_i,0} + \int_{\mathcal{T}} X_i(t) \beta_{z_i}(t) dt + \varepsilon_i$ with $h_z(X_i) = \alpha_{z_i,0} + \int_{\mathcal{T}} X_i(t) \alpha_{z_i}(t) dt$
 l_1 -Regularized MLE w.r.t the derivatives of the $\alpha(\cdot)$ and $\beta(\cdot)$ functions

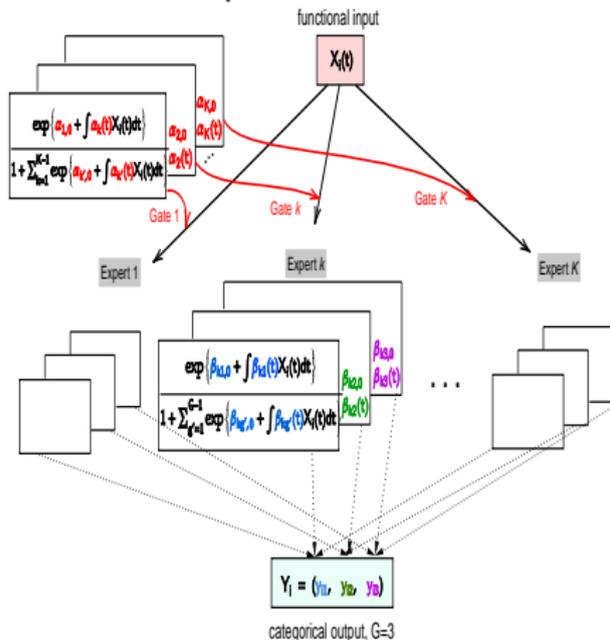
→ produces a meaningful sparse estimates for $\beta_{z_i}(t)$ curves :

$\beta_{z_i}^{(0)}(t) = 0$ implies that $X(t)$ has no effect on Y at t

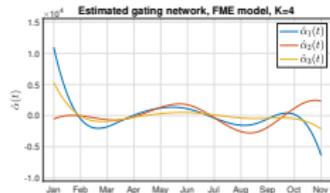
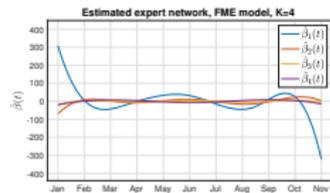
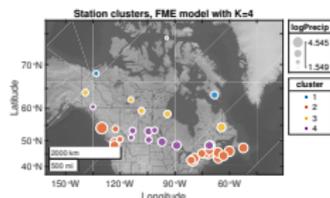
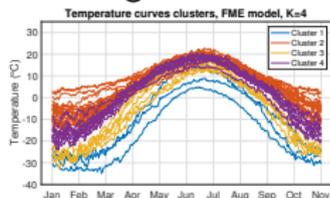
$\beta_{z_i}^{(1)}(t) = 0$ means that $\beta_{z_i}(t)$ is constant at t ,

$\beta_{z_i}^{(0)}(t) = 1$ shows that $\beta_{z_i}(t)$ is a linear function of t , etc.

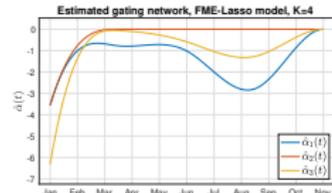
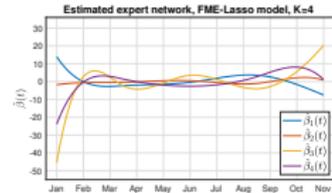
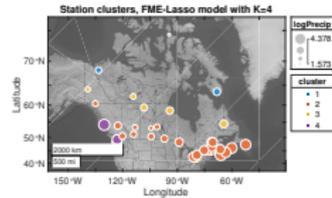
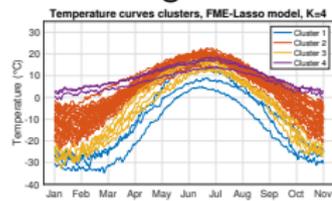
Mixture-of-Experts Architecture



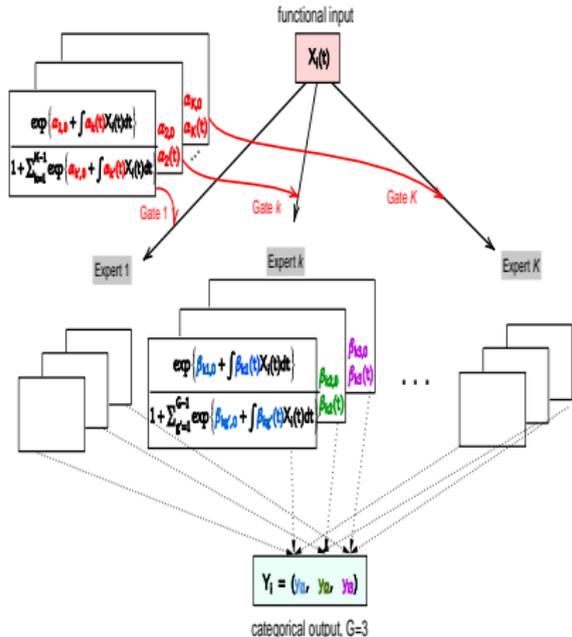
No regularization



LASSO regularization



Mixture-of-Experts Architecture



produces a meaningful sparse estimates for $\beta_{z_i}(t)$ curves :

$\beta_{z_i}^{(0)}(t) = 0$ implies that $X(t)$ has no effect on Y at t

$\beta_{z_i}^{(1)}(t) = 0$ means that $\beta_{z_i}(t)$ is constant at t ,

$\beta_{z_i}^{(0)}(t) = 1$ shows that $\beta_{z_i}(t)$ is a linear function of t , etc.

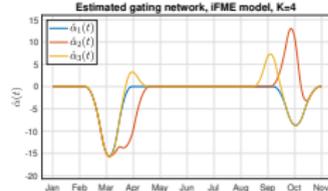
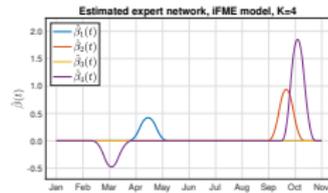
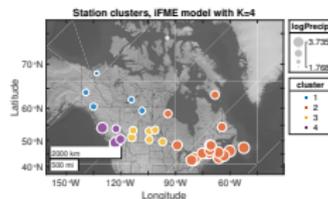
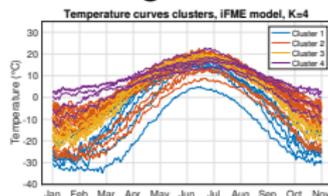
[PhD TN. Pham, 2022]

[arXiv :2202.13934]

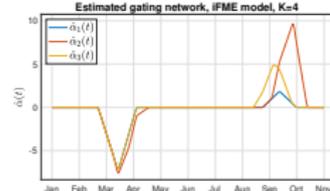
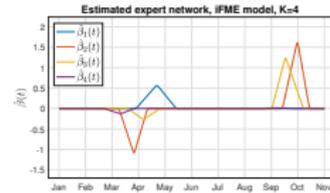
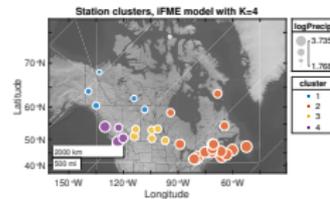
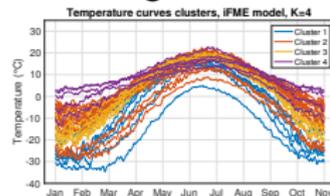
[Submitted, 2023]

[Contract, ANR SMILES]

OUR regularization



OUR regularization



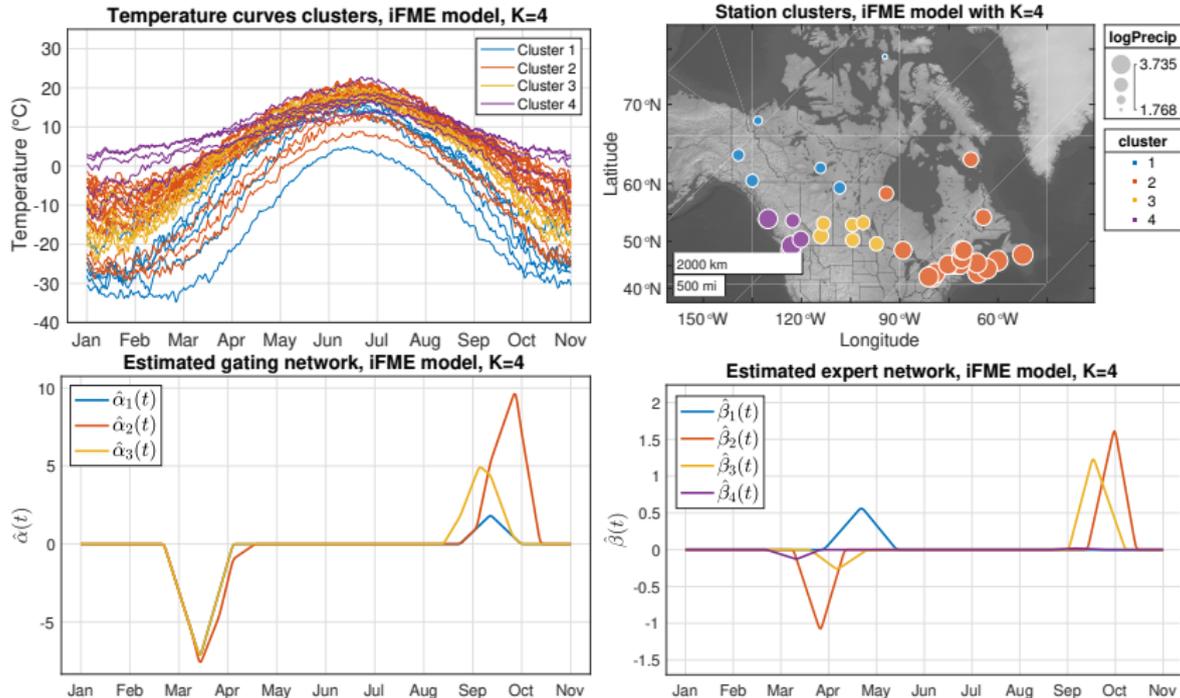


FIGURE – Clustering of temperatures (g.) predicted (log)-precipitations (d.) (Atlantic, Pacific, Continental and Arctic), and parameter functions (bg - bd.)

- Faïcel Chamroukhi, Thien Nhat Pham, Van Ha Hoang and Geoffrey J McLachlan. Functional Mixtures-of-Experts. arXiv preprint arXiv :2202.02249, 2022
- TrungTin Nguyen, Faïcel Chamroukhi, Hien D. Nguyen and Geoffrey J. McLachlan. Approximation of probability density functions via location-scale finite mixtures in Lebesgue spaces. Communications in Statistics - Theory and Methods Taylor & Francis., Vol. 52(14), pages :5048-5059, 2023
- Nguyen, TrungTin, Nguyen, Hien Duy, Chamroukhi, Faïcel and Forbes, Florence. A non-asymptotic approach for model selection via penalization in high-dimensional mixture of experts models. Electronic Journal of Statistics Institute of Mathematical Statistics and Bernoulli Society., Vol. 16(2), pages :4742-4822, 2022
- Chamroukhi, Faïcel, Brivet, Segolene, Savadjiev, Peter, Coates, Mark and Forghani, Reza. DECT-CLUST : Dual-Energy CT Image Clustering and Application to Head and Neck Squamous Cell Carcinoma Segmentation. Diagnostics, Vol. 12(12), pages :3072, 2022
- Nguyen Hien Duy, Nguyen TrungTin, Chamroukhi Faïcel and McLachlan Geoffrey John. Approximations of conditional probability density functions in Lebesgue spaces via mixture of experts models. Journal of Statistical Distributions and Applications SpringerOpen., Vol. 8(1), pages :1-15, 2021
- Nguyen, T Tin, Nguyen, Hien D, Chamroukhi, Faïcel and McLachlan, Geoffrey J. Approximation by finite mixtures of continuous density functions that vanish at infinity. Cogent Mathematics & Statistics Cogent OA., Vol. 7(1), pages :1750861, 2020
- Faïcel Chamroukhi and Bao T. Huynh. Regularized Maximum Likelihood Estimation and Feature Selection in Mixtures-of-Experts Models. Journal de la Soci Franse de Statistique, Vol. 160(1), pages :57-85, March, 2019

Thank you for your attention !