

Model selection by penalization in mixture of experts models with a non-asymptotic approach

TrungTin Nguyen



LABORATOIRE
JEAN KUNTZMANN
MATHÉMATIQUES APPLIQUÉES - INFORMATIQUE



53èmes Journées de Statistique
Lyon, France

Outline and our contributions

- 1 Collection of GLoME models
 - Context and motivating example
 - Boundedness conditions
- 2 Model selection in GLoME and BLoME models
 - Asymptotic approach
 - Non-asymptotic approach with oracle inequalities
- 3 Main positive messages and perspectives

Outline

- 1 Collection of GLoME models
 - Context and motivating example
 - Boundedness conditions
- 2 Model selection in GLoME and BLoME models
 - Asymptotic approach
 - Non-asymptotic approach with oracle inequalities
- 3 Main positive messages and perspectives

- 1 Collection of GLoME models
 - Context and motivating example
 - Boundedness conditions
- 2 Model selection in GLoME and BLoME models
 - Asymptotic approach
 - Non-asymptotic approach with oracle inequalities
- 3 Main positive messages and perspectives

- ✍ **We have:** n random samples $(\mathbf{X}_i, \mathbf{Y}_i)_{i \in [n]} \in (\mathcal{X} \times \mathcal{Y})^n$ with observed values $(\mathbf{x}_i, \mathbf{y}_i)_{i \in [n]}$, $[n] = \{1, \dots, n\}$, arising from an unknown conditional density s_0 .
- ⚙ **Learning:** potentially **nonlinear regression models for high-dimensional heterogeneous data** between output \mathbf{Y} and input \mathbf{X} : **Regression analysis** + **Clustering** + **Model selection** (e.g., number of clusters, complexity in each cluster).
- 👉 **Our proposal:** using **mixture of experts (MoE¹)** regression models due to their flexibility and effectiveness, e.g., several universal approximation theorems. ^{2 3 4}

¹ Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. (1991). Adaptive mixtures of local experts. Neural computation.

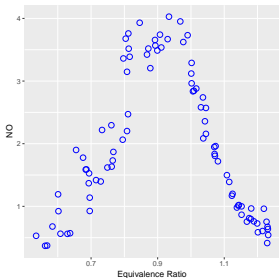
² Nguyen, H. D., Chamroukhi, F., and Forbes, F. (2019). Approximation results regarding the multiple-output Gaussian gated mixture of linear experts model. Neurocomputing.

³ Nguyen, H. D., **Nguyen, T.**, Chamroukhi, F., and McLachlan, G. J. (2021). Approximations of conditional probability density functions in Lebesgue spaces via mixture of experts models. *Journal of Statistical Distributions and Applications*.

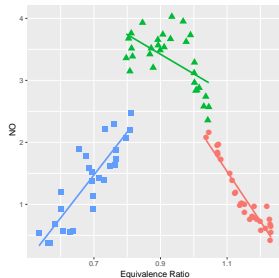
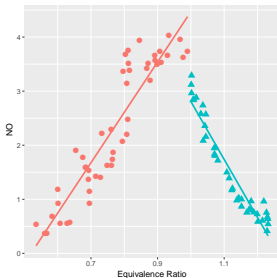
⁴ **Nguyen, T.**, Chamroukhi, F., Nguyen, H. D., and McLachlan, G. J. (2022). Approximation of probability density functions via location-scale finite mixtures in Lebesgue spaces. *Communications in Statistics - Theory and Method*.

Motivating example: Ethanol data set 88 observations

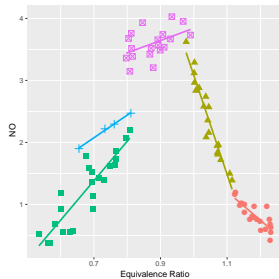
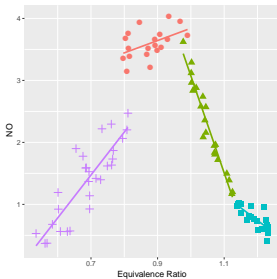
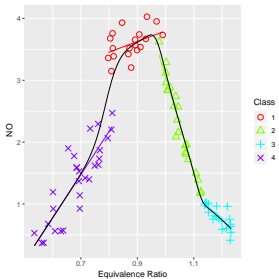
(a) Raw Ethanol data set



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



(b) Our best data-driven MoE model



Definition: GLLiM and GLoME models

$$s_{\psi_{K,d}}(\mathbf{x}|\mathbf{y}) = \sum_{k=1}^K \underbrace{\frac{\pi_k \mathcal{N}_L(\mathbf{y}; \mathbf{c}_k, \mathbf{\Gamma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}_L(\mathbf{y}; \mathbf{c}_j, \mathbf{\Gamma}_j)}}_{\text{Gaussian gating network}} \underbrace{\mathcal{N}_D(\mathbf{x}; \mathbf{v}_{k,d}(\mathbf{y}), \mathbf{\Sigma}_k)}_{\text{Gaussian expert}}.$$

- $\omega = (\pi, \mathbf{c}, \mathbf{\Gamma}) \in (\Pi_{K-1} \times \mathbf{C}_K \times \mathbf{V}'_K) = \Omega_K$, Π_{K-1} : probability simplex, $K \in \mathbb{N}^*$: number of mixture components.
- $d \in \mathbb{N}^*$: mean functions' hyperparameter e.g., degree of polynomial.
- $\psi_{K,d} = (\omega, \mathbf{v}, \mathbf{\Sigma}) \in \Omega_K \times \Upsilon_{K,d} \times \mathbf{V}_K$: model parameter.

High-dimensional data using inverse regression frameworks (GLLiM models⁵): $\mathbf{Y} \equiv \text{input}$, $\mathbf{X} \equiv \text{output}$, $\mathcal{X} \subset \mathbb{R}^D$, $\mathcal{Y} \subset \mathbb{R}^L$, with $D \gg L$ and $D, L \in \mathbb{N}^*$.

⁵Deleforge, A., Forbes, F., and Horaud, R. (2015). High-dimensional regression with gaussian mixtures and partially-latent response variables. Statistics and Computing.

Definition: Gaussian gating networks

$$\mathbf{g}_k(\mathbf{y}; \boldsymbol{\omega}) = \frac{\pi_k \mathcal{N}_L(\mathbf{y}; \mathbf{c}_k, \boldsymbol{\Gamma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}_L(\mathbf{y}; \mathbf{c}_j, \boldsymbol{\Gamma}_j)}, \text{ for every } k \in [K],$$

- $\boldsymbol{\omega} = (\boldsymbol{\pi}, \mathbf{c}, \boldsymbol{\Gamma}) \in (\boldsymbol{\Pi}_{K-1} \times \mathbf{C}_K \times \mathbf{V}'_K) = \boldsymbol{\Omega}_K$,
- $\boldsymbol{\Pi}_{K-1} = \left\{ (\pi_k)_{k \in [K]} \in (\mathbb{R}^+)^K, \sum_{k=1}^K \pi_k = 1 \right\}$,
- \mathbf{C}_K : K -tuples of mean vectors of size $L \times 1$,
- \mathbf{V}'_K : K -tuples of elements in \mathcal{S}_L^{++} ,
- \mathcal{S}_L^{++} : collection of symmetric positive definite matrices on \mathbb{R}^L .

- 1 Collection of GLoME models
 - Context and motivating example
 - Boundedness conditions
- 2 Model selection in GLoME and BLoME models
 - Asymptotic approach
 - Non-asymptotic approach with oracle inequalities
- 3 Main positive messages and perspectives

Mild assumption: Boundedness conditions

- **Gaussian gating parameters:** there exist positive constants $a_\pi, A_c, a_\Gamma, A_\Gamma$ s.t.

$$\tilde{\Omega}_K = \{\omega \in \Omega_K : \forall k \in [K], \|\mathbf{c}_k\|_\infty \leq A_c, \\ a_\Gamma \leq m(\Gamma_k) \leq M(\Gamma_k) \leq A_\Gamma, a_\pi \leq \pi_k\}.$$

- **Gaussian experts linear combination of bounded functions means:**
 $\mathbf{v} = (\mathbf{v}_{k,d})_{k \in [K]} \in \Upsilon_{K,d} = \otimes_{k \in [K]} \Upsilon_{k,d} = \Upsilon_{k,d}^K$, where $\forall k \in [K]$,

$$\Upsilon_{k,d} = \Upsilon_{Bo,d} = \left\{ \mathbf{y} \mapsto \left(\sum_{i=1}^d \alpha_i^{(j)} \theta_{\Upsilon,i}(\mathbf{y}) \right)_{j \in [D]} : \|\alpha\|_\infty \leq T_\Upsilon \right\},$$

Collection of bounded basis functions: $\mathbf{y} \mapsto (\theta_{\Upsilon,i}(\mathbf{y}))_{i \in [d_\Upsilon]}$, $d \in \mathbb{N}^*$,
 $T_\Upsilon \in \mathbb{R}^+$.

Boundedness conditions on Gaussian expert covariance matrices

$$\mathbf{V}_K = \left\{ (\boldsymbol{\Sigma}_k)_{k \in [K]} \equiv \left(B_k \mathbf{P}_k \mathbf{A}_k \mathbf{P}_k^\top \right)_{k \in [K]} : B_- \leq B_k \leq B_+, \right. \\ \left. \mathbf{P}_k \in SO(D), \mathbf{A}_k \in \mathcal{A}(\lambda_-, \lambda_+) \right\} :$$

- $B_k = |\boldsymbol{\Sigma}_k|^{1/D}$: volume, $B_- \in \mathbb{R}^+, B_+ \in \mathbb{R}^+$,
- \mathbf{P}_k : eigenvectors of $\boldsymbol{\Sigma}_k$, $SO(D)$: special orthogonal group of dimension D ,
- \mathbf{A}_k : diagonal matrix of normalized eigenvalues of $\boldsymbol{\Sigma}_k$, $\mathcal{A}(\lambda_-, \lambda_+)$: diagonal matrices \mathbf{A}_k , such that $|\mathbf{A}_k| = 1$ and $\forall i \in [D], \lambda_- \leq (\mathbf{A}_k)_{i,i} \leq \lambda_+$, where $\lambda_-, \lambda_+ \in \mathbb{R}$.

⁶Celeux, G. and Govaert, G. (1995). Gaussian parsimonious clustering models. Pattern Recognition.

Definition: Collection of GLLiM and GLoME models

$$S_{\mathbf{m}} = \left\{ (\mathbf{x}, \mathbf{y}) \mapsto s_{\psi_{K,d}}(\mathbf{x}|\mathbf{y}) = s_{\mathbf{m}}(\mathbf{x}|\mathbf{y}) : \mathbf{m} = (K, d), \right. \\ \left. \psi_{K,d} = (\omega, \mathbf{v}, \Sigma) \in \tilde{\Omega}_K \times \Upsilon_{K,d} \times \mathbf{V}_K = \tilde{\Psi}_{K,d} \right\}.$$

- $\mathbf{m} \in \mathcal{M} = [K_{\max}] \times [d_{\max}], K_{\max}, d_{\max} \in \mathbb{N}^*.$

- 1 Collection of GLoME models
 - Context and motivating example
 - Boundedness conditions
- 2 Model selection in GLoME and BLoME models
 - Asymptotic approach
 - Non-asymptotic approach with oracle inequalities
- 3 Main positive messages and perspectives

Model selection in standard MoE regression models

- ⚙️ **Best data-driven model**: selecting from a collection of MoE models characterized by hyperparameters $\mathbf{m} = (K, d)$.
- **Penalized maximum likelihood estimator (PMLE)**:
 - **MLE is not sufficient**: underestimation of the risk of the estimate \Rightarrow choosing models too complex.
 - **PMLE via adding $\text{pen}(\mathbf{m})$** : compensate bias (too simple model) and variance (too complex model).
- ⚙️ **Our contributions**: establishing **non-asymptotic risk bounds** that take the form of **weak oracle inequalities**, provided that **lower bounds on the penalties** hold true.

Definition: Penalized maximum likelihood estimator (PMLE)

An η' -PMLE $\hat{s}_{\hat{\mathbf{m}}}$ (corresponding the selected model or best data-driven model $S_{\hat{\mathbf{m}}}$ among $(S_{\mathbf{m}})_{\mathbf{m} \in \mathcal{M}}$), defined by

$$\sum_{i=1}^n -\ln(\hat{s}_{\hat{\mathbf{m}}}(\mathbf{x}_i|\mathbf{y}_i)) + \text{pen}(\hat{\mathbf{m}}) \leq \inf_{\mathbf{m} \in \mathcal{M}} \left(\sum_{i=1}^n -\ln(\hat{s}_{\mathbf{m}}(\mathbf{x}_i|\mathbf{y}_i)) + \text{pen}(\mathbf{m}) \right) + \eta',$$

- $\hat{s}_{\hat{\mathbf{m}}}$ is an η -minimizer of the negative log-likelihood (infimum may not be unique or reached) is defined by

$$\sum_{i=1}^n -\ln(\hat{s}_{\hat{\mathbf{m}}}(\mathbf{x}_i|\mathbf{y}_i)) \leq \inf_{s_{\mathbf{m}} \in S_{\mathbf{m}}} \sum_{i=1}^n -\ln(s_{\mathbf{m}}(\mathbf{x}_i|\mathbf{y}_i)) + \eta,$$

- $\text{pen}(\mathbf{m})$: penalty function \leftarrow choosing it is tricky but obviously necessary to compensate variance and bias.

Definition: Loss functions for conditional densities

- **Tensorized Kullback-Leibler divergence $\text{KL}^{\otimes n}$ (conditional densities and random covariate variables):**

$$\text{KL}^{\otimes n}(s, t) = \mathbb{E}_{\mathbf{Y}_{[n]}} \left[\frac{1}{n} \sum_{i=1}^n \text{KL}(s(\cdot | \mathbf{Y}_i), t(\cdot | \mathbf{Y}_i)) \right],$$

if $sdy \ll tdy$, $+\infty$ otherwise. Fixed predictors \Rightarrow no $\mathbb{E}_{\mathbf{Y}_{[n]}}[\cdot]$.

- **Tensorized Jensen-Kullback-Leibler divergence $\text{JKL}_{\rho}^{\otimes n}$ (technical difficulties with conditional densities), given $\rho \in (0, 1)$,**

$$\text{JKL}_{\rho}^{\otimes n}(s, t) = \mathbb{E}_{\mathbf{Y}_{[n]}} \left[\frac{1}{n} \sum_{i=1}^n \frac{1}{\rho} \text{KL}(s(\cdot | \mathbf{Y}_i), (1 - \rho)s(\cdot | \mathbf{Y}_i) + \rho t(\cdot | \mathbf{Y}_i)) \right].$$

- 1 Collection of GLoME models
 - Context and motivating example
 - Boundedness conditions
- 2 Model selection in GLoME and BLoME models
 - Asymptotic approach
 - Non-asymptotic approach with oracle inequalities
- 3 Main positive messages and perspectives

Some asymptotic approaches for model selection in MoE models

- Akaike information criterion (AIC) [Akaike, 1974], Bayesian information criterion (BIC) [Schwarz et al., 1978] and BIC-like approximation of integrated classification likelihood (ICL-BIC) [Biernacki et al., 2000] criteria:

$$\text{pen}_{\text{AIC}}(\mathbf{m}) = \dim(S_{\mathbf{m}}), \quad \text{pen}_{\text{BIC}}(\mathbf{m}) = \frac{\ln(n) \dim(S_{\mathbf{m}})}{2}.$$

$$\text{pen}_{\text{ICL-BIC}}(\mathbf{m}) = \text{pen}_{\text{BIC}}(\mathbf{m}) + \text{ENT}(\mathbf{m}) \leftarrow \text{estimated mean entropy}.$$

- AIC (based on asymptotic theory), BIC, ICL-BIC (based on Bayesian approach):
 - May be wrong in a non-asymptotic context: $\dim(S_{\mathbf{m}})$ and $\text{card}(\mathcal{M})$ depend on and can be much larger than n .
 - No finite sample guarantees.
- Obtain an upper bound on $\mathbb{E} [\text{KL}^{\otimes n}(s_0, \hat{s}_{\mathbf{m}})]$:
 - ✓ Finite sample guarantee.
 - ✗ Strong regularity assumptions of [White, 1982].

- 1 Collection of GLoME models
 - Context and motivating example
 - Boundedness conditions
- 2 Model selection in GLoME and BLoME models
 - Asymptotic approach
 - Non-asymptotic approach with oracle inequalities
- 3 Main positive messages and perspectives

Non-asymptotic upper bound of a single model

➤ **Initial target:**


$$\mathbb{E} [\text{KL}^{\otimes n}(s_0, \hat{s}_m)] \leq \left(\inf_{\psi_m \in \Psi_m} \text{KL}^{\otimes n}(s_0, s_{\psi_m}) + \frac{1}{2n} \dim(S_m) \right) + C_2 \frac{1}{n}.$$


➤ **Our contribution:**

$$\mathbb{E} [\text{JKL}_{\rho}^{\otimes n}(s_0, \hat{s}_m)] \leq C_1 \left(\inf_{\psi_m \in \Psi_m} \text{KL}^{\otimes n}(s_0, s_{\psi_m}) + \frac{\kappa}{n} \mathfrak{D}_m \right) + C_2 \frac{1}{n}.$$

- ① Different divergences: $\text{JKL}_{\rho}^{\otimes n}(s_0, \hat{s}_m) \leq \text{KL}^{\otimes n}(s_0, \hat{s}_m)$.
- ② $C_1 > 1$, κ is a constant that depends on C_1 , $\mathfrak{D}_m \propto \dim(S_m)$.

Theorem: Non-asymptotic oracle inequality⁷

 **Assumptions:** given a deterministic collection $(S_m)_{m \in \mathcal{M}}$ of MoE models, $\rho \in (0, 1)$, $C_1 > 1$,
 $\Xi = \sum_{m \in \mathcal{M}} e^{-z_m} < \infty, z_m \in \mathbb{R}^+, \forall m \in \mathcal{M}$.

 **Conclusion:** there exist constants C and $\kappa(\rho, C_1) > 0$ such that whenever for all $m \in \mathcal{M}$,

$$\text{pen}(\mathbf{m}) \geq \kappa(\rho, C_1) [(C + \ln n) \dim(S_{\mathbf{m}}) + z_{\mathbf{m}}],$$

the η' -PMLE $\hat{s}_{\mathbf{m}}$ satisfies

$$\begin{aligned} \mathbb{E} \left[\text{JKL}_{\rho}^{\otimes n}(s_0, \hat{s}_{\mathbf{m}}) \right] &\leq C_1 \inf_{\mathbf{m} \in \mathcal{M}} \left(\inf_{s_{\mathbf{m}} \in S_{\mathbf{m}}} \text{KL}^{\otimes n}(s_0, s_{\mathbf{m}}) + \frac{\text{pen}(\mathbf{m})}{n} \right) \\ &\quad + \frac{\kappa(\rho, C_1) C_1 \Xi}{n} + \frac{\eta + \eta'}{n}. \end{aligned}$$

⁷ Nguyen, T., Nguyen, H.D., Chamroukhi, F., and Forbes, F. (2022). A non-asymptotic approach for model selection via penalization in high-dimensional mixture of experts. arXiv 2104.02640. Under revision, Electronic Journal of Statistics.

- 1 Collection of GLoME models
 - Context and motivating example
 - Boundedness conditions
- 2 Model selection in GLoME and BLoME models
 - Asymptotic approach
 - Non-asymptotic approach with oracle inequalities
- 3 Main positive messages and perspectives

Main positive messages and perspectives

- 😊 **Our risk assessments are non-asymptotic.**
- 😊 If $\text{pen}(\mathbf{m})$ is properly chosen, then our PMLE behaves in a comparable manner compared to **the best (oracle) model $S_{\mathbf{m}^*}$ in the collection.**
- 😊 **Partially answer** the two following important questions raised in the area of MoE regression models:
 - ① **Which value of K** should be chosen, given the sample size n ,
 - ② Whether it is better to use a **few complex experts** or **combine many simple experts**, given the total number of parameters.
- 😊 **Minimax lower bounds** for MoE regression models, which is only known for mixture models⁸.

⁸Maugis-Rabusseau, C. and Michel, B. (2013). Adaptive density estimation for clustering with Gaussian mixtures. ESAIM: Probability and Statistics.

My Coauthors \in Mixture of French and Australian Experts



Faïcel Chamroukh

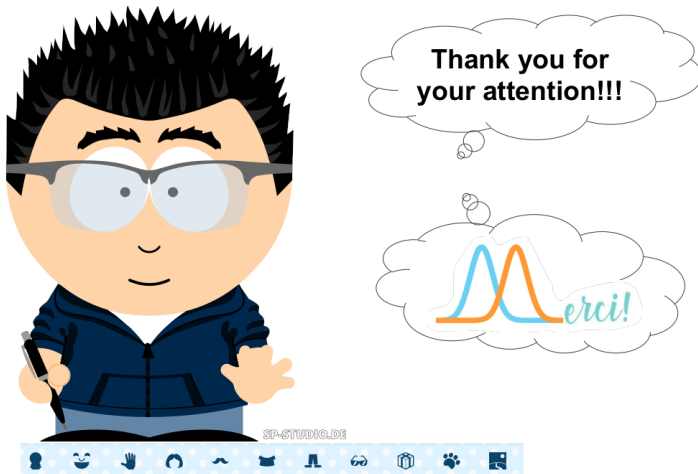


Hien Duy Nguyen



Florence Forbes

“Essentially, all models are wrong, but some are useful”. George E.P. Box (1987).



↑ This is my best data-driven model to approximate myself.