

# Non-asymptotic penalization criteria for model selection in mixture of experts models

**TrungTin Nguyen**

Université de Caen Normandie, France.

<https://trung-tinnguyen.github.io/>

**MiMo 2021: Workshop on Mixture Models**

Joint work with

Faïcel Chamroukhi (Université de Caen Normandie, France),

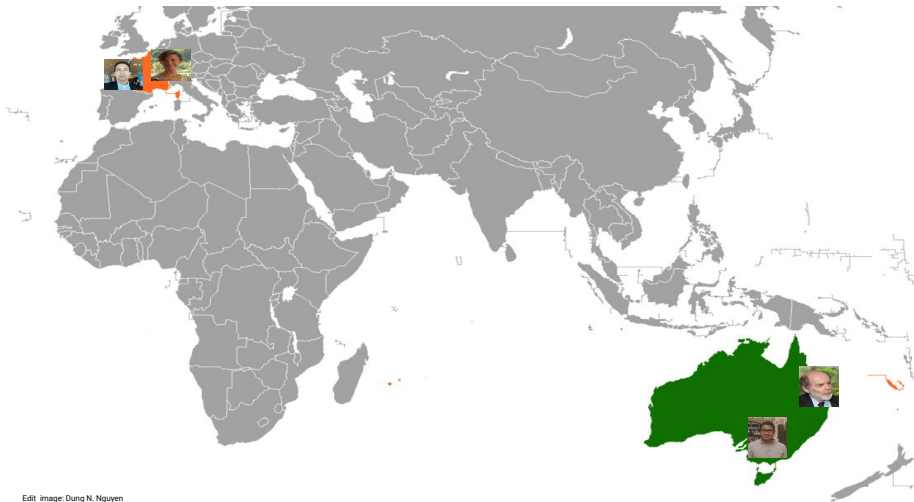
Hien Duy Nguyen (La Trobe University, Australia),

Geoffrey J McLachlan (University of Queensland, Australia),

Florence Forbes (Inria Grenoble-Rhône-Alpes, France).

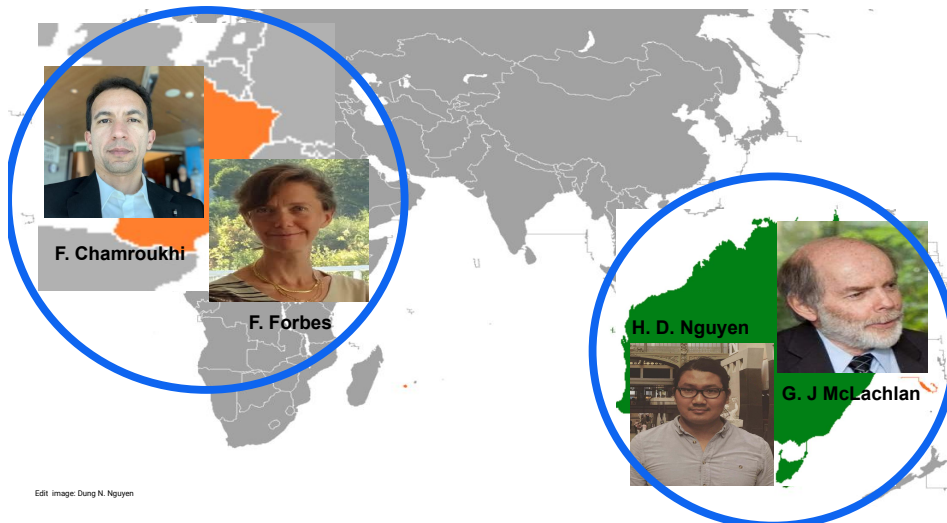
April 09, 2021

# A mixture of French and Australian teams...



Edit image: Dung N. Nguyen

# Joint work with...



- 1 Collection of mixture of experts models
- 2 Non-asymptotic oracle inequality
- 3 Numerical experiment
- 4 Future work on non-asymptotic oracle inequality

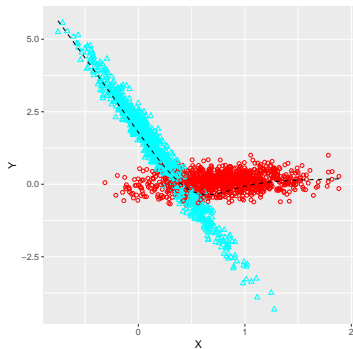
- 1 Collection of mixture of experts models
- 2 Non-asymptotic oracle inequality
- 3 Numerical experiment
- 4 Future work on non-asymptotic oracle inequality

- 1 Collection of mixture of experts models
- 2 Non-asymptotic oracle inequality
- 3 Numerical experiment
- 4 Future work on non-asymptotic oracle inequality

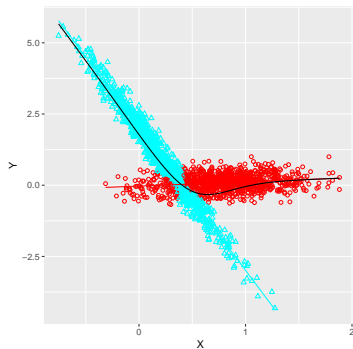
# Nonlinear regression models for heterogeneous data

- We have:  $n$  random samples  $(\mathbf{X}_i, \mathbf{Y}_i)_{i \in [n]} \in (\mathcal{X} \times \mathcal{Y})^n \subset (\mathbb{R}^D \times \mathbb{R}^L)^n$  with the corresponding observed values  $(\mathbf{x}_i, \mathbf{y}_i)_{i \in [n]}$ ,  $[n] := \{1, \dots, n\}$ , arising from an unknown conditional density  $s_0$ .
- Objective: learning **nonlinear regression models for heterogeneous data** between the multivariate response  $\mathbf{Y} = (\mathbf{Y}_j)_{j \in [L]}$ , and the set of covariates  $\mathbf{X} = (\mathbf{X}_j)_{j \in [D]}$ .
- Our proposal: approximating  $s_0$  by a **Gaussian-gated localized mixture of experts (GLoME)** model due to its flexibility and effectiveness.
  - **Model selection problem**: estimating the number of mixture components via penalized maximum likelihood estimators.
  - **Non-asymptotic oracle inequality**: providing a lower bound on the penalty that ensures a weak oracle inequality is satisfied by our estimator.

# Typical realization and regression clustering results



(a) Typical realization: WS case

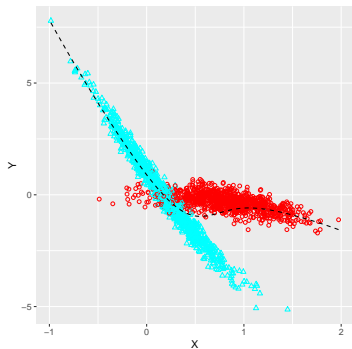


(b) Clustering by GLoME

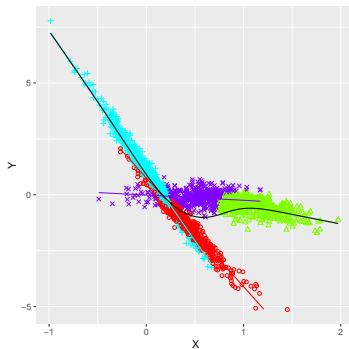
Clustering deduced from the estimated conditional density of GLoME by a MAP principle with 2000 data points of example WS. The dash and solid black curves present the true and estimated mean functions.



# Typical realization and regression clustering results



(a) Typical realization: MS case



(b) Clustering by GLoME

Clustering deduced from the estimated conditional density of GLoME by a MAP principle with 2000 data points of example MS. The dash and solid black curves present the true and estimated mean functions.

GLoME is a mixture of experts (MoE) model  
[Jacobs et al., 1991, Xu et al., 1995, Nguyen and Chamroukhi, 2018]:

- Generalizing the classical finite mixtures and finite mixtures regression models [McLachlan and Peel, 2000].
- Containing a supervised Gaussian locally-linear mapping (GLLiM) model for high-dimensional regression data ( $D \gg L$ ) [Deleforge et al., 2015].
- Approximation capabilities of MoE [Mendes and Jiang, 2012, Ho et al., 2019, Nguyen et al., 2020a].

Motivated by **GLLiM models, an inverse regression framework**, in GLoME models,  $\mathbf{Y}$  becomes the covariates and  $\mathbf{X}$  plays the role of a multivariate response.

## Definition

$$s_{\psi_K}(\mathbf{x}|\mathbf{y}) = \sum_{k=1}^K \underbrace{\mathbf{g}_k(\mathbf{y}; \boldsymbol{\omega})}_{\text{Gaussian gating function}} \underbrace{\Phi_D(\mathbf{x}; \mathbf{v}_k(\mathbf{y}), \boldsymbol{\Sigma}_k)}_{\text{Gaussian expert}},$$

- $K \in \mathbb{N}^*$ : number of mixture components,
- $\psi_K = (\boldsymbol{\omega}, \mathbf{v}, \boldsymbol{\Sigma}) \in \boldsymbol{\Omega}_K \times \boldsymbol{\Upsilon}_K \times \mathbf{V}_K =: \boldsymbol{\Psi}_K$ : model parameter.

## Definition

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

- $\boldsymbol{\omega} = (\boldsymbol{\pi}, \mathbf{c}, \boldsymbol{\Gamma}) \in (\boldsymbol{\Pi}_{K-1} \times \mathbf{C}_K \times \mathbf{V}'_K) =: \boldsymbol{\Omega}_K$ ,
- $\boldsymbol{\Pi}_{K-1} = \left\{ (\pi_k)_{k \in [K]} \in (\mathbb{R}^+)^K, \sum_{k=1}^K \pi_k = 1 \right\}$ ,
- $\mathbf{C}_K$ :  $K$ -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  $\mathbb{R}^L$ .

## Assumptions

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

- $a_\pi, A_c, a_\Gamma, A_\Gamma$ : positive constants,
- $m(\mathbf{A})$  and  $M(\mathbf{A})$ : the modulus of the smallest and largest eigenvalues of a matrix  $\mathbf{A}$ , respectively.

# Boundedness conditions on the Gaussian expert means

## Assumptions ( [Montuelle and Pennec, 2014] )

- *Linear combination of bounded functions:*  $d_{\mathbf{r}} \in \mathbb{N}^*$ ,  $T_{\mathbf{r}} \in \mathbb{R}^+$ ,  
 $\mathbf{v} = (\mathbf{v}_k)_{k \in [K]} \in \mathbf{r}_K = \mathbf{r}_b^K$ ,

$$\mathbf{r}_b = \left\{ \mathbf{y} \mapsto \left( \sum_{i=1}^{d_{\mathbf{r}}} \alpha_i^{(j)} \varphi_{\mathbf{r},i}(\mathbf{y}) \right)_{j \in [D]} =: (\mathbf{v}_j(\mathbf{y}))_{j \in [D]} : \|\alpha\|_{\infty} \leq T_{\mathbf{r}} \right\},$$

$(\varphi_{\mathbf{r},i})_{i \in [d_{\mathbf{r}}]}$  is a collection of bounded functions on  $\mathcal{Y}$ .

- *Polynomial means:*  $d'_{\mathbf{r}} \in \mathbb{N}^*$ ,  $\mathcal{Y} = [0, 1]^L$ ,  $\mathbf{r}_K = \mathbf{r}_p^K$ ,

$$\mathbf{r}_p = \left\{ \mathbf{y} \mapsto \left( \sum_{|\mathbf{r}|=0}^{d'_{\mathbf{r}}} \alpha_{\mathbf{r}}^{(j)} \mathbf{y}^{\mathbf{r}} \right)_{j \in [D]} =: (\mathbf{v}_j(\mathbf{y}))_{j \in [D]} : \|\alpha\|_{\infty} \leq T_{\mathbf{r}} \right\}.$$

# Boundedness conditions on the Gaussian expert covariance matrices

Assumptions ([Celeux and Govaert, 1995])

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

- $B_k = |\boldsymbol{\Sigma}_k|^{1/D}$ : volume,  $B_- \in \mathbb{R}^+, B_+ \in \mathbb{R}^+$ ,
- $\mathbf{P}_k$ : eigenvectors of  $\boldsymbol{\Sigma}_k$  belongs to the special orthogonal  $SO(D)$ ,
- $\mathbf{A}_k$ : diagonal matrix of normalized eigenvalues of  $\boldsymbol{\Sigma}_k$ , such that  $|\mathbf{A}_k| = 1$  and  $0 < \forall i \in [D], \lambda_- \leq (\mathbf{A}_k)_{i,i} \leq \lambda_+$ .

# Collection of GLoME models

Objective: estimating  $s_0$  by conditional densities belonging to the collection of GLoME models  $(S_m)_{m \in \mathcal{M}}$ , defined by

## Definition

$$S_m = \left\{ \mathcal{X} \times \mathcal{Y} \ni (\mathbf{x}, \mathbf{y}) \mapsto s_{\psi_K}(\mathbf{x}|\mathbf{y}) =: s_m(\mathbf{x}|\mathbf{y}) : \right. \\ \left. \psi_K = (\omega, \mathbf{v}, \Sigma) \in \tilde{\Omega}_K \times \mathbf{r}_K \times \mathbf{V}_K =: \tilde{\Psi}_K \right\},$$

- $\mathcal{M} = \{K \in [K_{\max}], K_{\max} \in \mathbb{N}^*\},$
- $\dim(S_m) = \dim(\tilde{\Psi}_K) = \dim(\tilde{\Omega}_K) + \dim(\mathbf{r}_K) + \dim(\mathbf{V}_K).$



# Outline

- 1 Collection of mixture of experts models
- 2 Non-asymptotic oracle inequality**
- 3 Numerical experiment
- 4 Future work on non-asymptotic oracle inequality

# Penalized maximum likelihood estimator (PMLE)

## Definition

An  $\eta'$ -**PMLE**  $\hat{m}$  (corresponding **the best model**  $S_{\hat{m}}$  among  $(S_m)_{m \in \mathcal{M}}$ ) is defined as follows

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

- $\hat{s}_m$  is an  $\eta$ -**minimizer** of the negative log-likelihood (NLL) (infimum may not be unique or reached), defined by

$$\hat{s}_m = \operatorname{argmin}_{s_m \in S_m} \sum_{i=1}^n -\ln(s_m(\mathbf{x}_i | \mathbf{y}_i)) + \eta,$$

- $\operatorname{pen}(m)$ : compensating variance and bias.

## Definition

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

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


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


- **Tensorized Jensen-Kullback-Leibler divergence  $JKL_{\rho}^{\otimes n}$  (technical difficulties with GLoME models)  $\rho \in (0, 1)$ ,**

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

# Main result on a non-asymptotic oracle inequality

Theorem ([[Nguyen et al., 2021](#)])

 **Assumptions:** given a collection  $(S_m)_{m \in \mathcal{M}}$  of GLoME 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  $m \in \mathcal{M}$ ,

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

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

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

# Relationship: soft-max and Gaussian gating functions

$$\begin{aligned}\mathcal{P}_S &= \left\{ \mathbf{y} \mapsto (\mathbf{g}_k(\mathbf{y}; \gamma))_{k \in [K]} = \left( \frac{\exp(\mathbf{a}_k + \mathbf{b}_k^\top \mathbf{y})}{\sum_{l=1}^K \exp(\mathbf{a}_l + \mathbf{b}_l^\top \mathbf{y})} \right)_{k \in [K]}, \gamma \in \Gamma_S \right\}, \\ \Gamma_S &= \left\{ \gamma = ((\mathbf{a}_k)_{k \in [K]}, (\mathbf{b}_k)_{k \in [K]}) \in \mathbb{R}^K \times (\mathbb{R}^L)^K \right\}, \\ \mathcal{P}_G &= \left\{ \mathbf{y} \mapsto (\mathbf{g}_k(\mathbf{y}; \omega))_{k \in [K]} = \left( \frac{\pi_k \Phi_L(\mathbf{y}; \mathbf{c}_k, \Gamma_k)}{\sum_{j=1}^K \pi_j \Phi_L(\mathbf{y}; \mathbf{c}_j, \Gamma_j)} \right)_{k \in [K]}, \omega \in \Omega_K \right\}.\end{aligned}$$

Lemma ([[Nguyen et al., 2020a](#)])

*In general,  $\mathcal{P}_S \subset \mathcal{P}_G$ . If all  $\Gamma_k$ ,  $k \in [K]$ , are identical, then  $\mathcal{P}_G = \mathcal{P}_S$ .*

► Obtaining finite-sample oracle inequality for GLoME model is much more challenging compared to soft-max-gated mixture of experts (SGaME) model [[Montuelle and Pennec, 2014](#)].

# Reparameterization the space of Gaussian gating functions

- Make use of the results for bracketing entropy of logistic weights from SGaME models [[Montuelle and Pennec, 2014](#)].

👉 Reparameterization trick:


$$\mathbf{W}_K = \left\{ \mathbf{y} \mapsto (\ln(\pi_k \Phi_L(\mathbf{y}; \mathbf{c}_k, \mathbf{\Gamma}_k)))_{k \in [K]} =: \mathbf{w}(\mathbf{y}; \boldsymbol{\omega}) : \boldsymbol{\omega} \in \tilde{\Omega}_K \right\},$$
$$\mathcal{P}_K = \left\{ \mathbf{y} \mapsto \left( \frac{e^{\mathbf{w}_k(\mathbf{y})}}{\sum_{l=1}^K e^{\mathbf{w}_l(\mathbf{y})}} \right)_{k \in [K]} =: (g_{\mathbf{w},k}(\mathbf{y}))_{k \in [K]}, \mathbf{w} \in \mathbf{W}_K \right\}.$$

# Asymptotic theory of a single parametric model

**Misspecified case:**  $s_0 \notin S_m$ ,  $\psi_m^* = \operatorname{argmin}_{\psi_m \in \Psi_m} \text{KL}^{\otimes n}(s_0, s_{\psi_m})$ ,


$$S_m = \left\{ (\mathbf{x}, \mathbf{y}) \mapsto s_{\psi_m}(\mathbf{x}|\mathbf{y}) =: s_m(\mathbf{x}|\mathbf{y}) : \psi_m \in \Psi_m \subset \mathbb{R}^{\dim(S_m)} \right\}.$$

Theorem ([White, 1982, Cohen and Pennec, 2011])

 *Assumptions:*  $S_m$  is identifiable, some strong regularity assumptions on  $\psi_m \mapsto s_{\psi_m}$ ,  $\exists \mathbf{A}(\psi_m)$  and  $\mathbf{B}(\psi_m)$ :

$$[\mathbf{A}(\psi_m)]_{k,l} = \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \int \frac{\partial^2 \ln s_{\psi_m}}{\partial \psi_{m,k} \partial \psi_{m,l}}(\mathbf{x}|\mathbf{Y}_i) s_0(\mathbf{x}|\mathbf{Y}_i) d\lambda \right],$$

$$[\mathbf{B}(\psi_m)]_{k,l} = \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \int \frac{\partial \ln s_{\psi_m}}{\partial \psi_{m,k}}(\mathbf{x}|\mathbf{Y}_i) \frac{\partial \ln s_{\psi_m}}{\partial \psi_{m,l}}(\mathbf{x}|\mathbf{Y}_i) s_0(\mathbf{x}|\mathbf{Y}_i) d\lambda \right].$$

 *Conclusion:*  $\mathbb{E} [\text{KL}^{\otimes n}(s_0, \hat{s}_m)]$  is *asymptotically equivalent* to

$$\text{KL}^{\otimes n}(s_0, s_{\psi_m^*}) + \frac{1}{2n} \text{tr} \left( \mathbf{B}(\psi_m^*) \mathbf{A}(\psi_m^*)^{-1} \right).$$

# Asymptotic theory of a single parametric model

**Well-specified case:**  $s_0 \in S_m$ ,  $\psi_m^* = \operatorname{argmin}_{\psi_m \in \Psi_m} \text{KL}^{\otimes n}(s_0, s_{\psi_m})$ .

Theorem ([White, 1982, Cohen and Pennec, 2011])

It holds that

$$s_0 = s_{\psi_m^*}, \mathbf{A}(\psi_m^*) = \mathbf{B}(\psi_m^*).$$

The same assumption in misspecified case:  $\mathbb{E}[\text{KL}^{\otimes n}(s_0, \hat{s}_m)]$  is *asymptotically equivalent* to

$$\underbrace{\text{KL}^{\otimes n}(s_0, s_{\psi_m^*})}_{=0} + \frac{1}{2n} \dim(S_m).$$



# Drawbacks of asymptotic theory

**Well-specified case:**  $s_0 \in S_m$ .

- ★ Problem: **asymptotic normality** of  $\sqrt{n} \left( \hat{\psi}_m - \psi_m^* \right)$  is required!
- ➡ Some previous ideas to handle **non-asymptotic normality**:
  - Extension in non parametric case or non-identifiable model, **Wilk's phenomenon**, [Wilks, 1938].
  - Generalization of the corresponding Chi-Square goodness-of-fit test [Fan et al., 2001].
  - Finite sample deviation of the corresponding empirical quantity in a bounded loss setting [Boucheron and Massart, 2011].

# Non-asymptotic upper bound of a single model

➤ **Initial target:**

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

➤ **Our contribution:** weaker than expected!

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

- ①  $\text{JKL}_{\rho}^{\otimes n}(s_0, \widehat{s}_m) \leq \text{KL}^{\otimes n}(s_0, \widehat{s}_m)$ ,
- ②  $C_1 > 1$ ,  $\kappa$  is a constant that depends on  $C_1$ ,
- ③ Model complexity:  $\mathfrak{D}_m \leftrightarrow \dim(S_m)$ .

❄ **Existence of a corresponding lower bound for GLoME models: still an open question!** (Gaussian mixture models (GMM) [Maugis-Rabusseau and Michel, 2013]).

# Remarks on our weak oracle inequality

$$\mathbb{E} \left[ \text{JKL}_{\rho}^{\otimes n} (s_0, \hat{s}_{\widehat{m}}) \right] \leq C_1 \inf_{m \in \mathcal{M}} \left( \inf_{s_m \in S_m} \text{KL}^{\otimes n} (s_0, s_m) + \frac{\text{pen}(m)}{n} \right) + \frac{\kappa(\rho, C_1) C_1 \Xi}{n} + \frac{\eta + \eta'}{n}.$$

## Potential issues

- 1 Different divergences:  $\text{JKL}_{\rho}^{\otimes n} (s_0, \hat{s}_{\widehat{m}}) \leq \text{KL}^{\otimes n} (s_0, \hat{s}_{\widehat{m}})$ .
- 2  $C_1 > 1$  and misspecified case: as  $n \rightarrow \infty$ , the error bound  $\rightarrow C_1 \inf_{s_m \in S_m} \text{KL}^{\otimes n} (s_0, s_m)$  (potentially large!).
- 3  $\frac{\text{pen}(m)}{n}$  is not directly related to the variance (asymptotic variance in the parametric case  $\dim(S_m)/n$ ).

# Solution for different divergences

$$\mathbb{E} \left[ \text{JKL}_{\rho}^{\otimes n} (s_0, \hat{s}_{\hat{m}}) \right] \leq C_1 \inf_{m \in \mathcal{M}} \left( \inf_{s_m \in S_m} \text{KL}^{\otimes n} (s_0, s_m) + \frac{\text{pen}(m)}{n} \right) + \frac{\kappa(\rho, C_1) C_1 \Xi}{n} + \frac{\eta + \eta'}{n}.$$

- In general:  $\text{JKL}_{\rho}^{\otimes n} (s_0, \hat{s}_{\hat{m}}) \leq \text{KL}^{\otimes n} (s_0, \hat{s}_{\hat{m}})$ .
- If  $\sup_{m \in \mathcal{M}} \sup_{s_m \in S_m} \|s_0/s_m\|_{\infty} < \infty \Leftrightarrow \mathcal{Y}$  is compact,  $s_0$  is compactly supported, the regression functions are uniformly bounded, and a uniform lower bound on the eigenvalues of the covariance matrices, Proposition 1 from [Cohen and Pennec, 2011] implies that

$$\frac{C_{\rho}}{2 + \ln \|s_0/\hat{s}_{\hat{m}}\|_{\infty}} \text{KL}^{\otimes n} (s_0, \hat{s}_{\hat{m}}) \leq \text{JKL}_{\rho}^{\otimes n} (s_0, \hat{s}_{\hat{m}}).$$

# Solution for misspecified model

$$\mathbb{E} \left[ \text{JKL}_{\rho}^{\otimes n} (s_0, \hat{s}_{\hat{m}}) \right] \leq C_1 \inf_{m \in \mathcal{M}} \left( \inf_{s_m \in S_m} \text{KL}^{\otimes n} (s_0, s_m) + \frac{\text{pen}(m)}{n} \right) + \frac{\kappa(\rho, C_1) C_1 \Xi}{n} + \frac{\eta + \eta'}{n}.$$

❄  $C_1 = 1$  with  $\text{KL}^{\otimes n}$  loss: still an open question!

👉 Bias  $\inf_{s_m \in S_m} \text{KL}^{\otimes n} (s_0, s_m)$ : small for  $\mathcal{M}$  well-chosen via approximation capabilities of MoE and GMM models  
[[Nguyen et al., 2019](#), [Nguyen et al., 2020c](#), [Nguyen et al., 2020b](#), [Nguyen et al., 2020a](#)].

# Outline

- 1 Collection of mixture of experts models
- 2 Non-asymptotic oracle inequality
- 3 Numerical experiment**
- 4 Future work on non-asymptotic oracle inequality

# Procedure for collection of supervised GLLiM models

**Goal:** look for the best model among  $(S_m^*)_{m \in \mathcal{M}}$ ,  $\mathcal{M} = [K_{\max}]$  based on  $(\mathbf{x}_i, \mathbf{y}_i)_{i \in [n]}$  arising from an forward conditional density  $s_0^*$ :

- 1 Each  $m \in \mathcal{M}$ : estimate the forward MLE  $(\hat{s}_m^*(\mathbf{y}_i | \mathbf{x}_i))_{i \in [n]}$  by inverse MLE  $\hat{s}_m$  via an **inverse regression trick** by GLLiM-EM algorithm (**xLLiM** package).
- 2 Calculate  $\eta'$ -PMLE  $\hat{m}$  with  $\text{pen}(m) = \kappa \dim(S_m^*)$ .

★ **Large enough but not explicit value** for  $\kappa$ !

- Asymptotic criteria: AIC ( $\kappa = 1$ ) and BIC ( $\kappa = \frac{\ln n}{2}$ ) [[Akaike, 1974](#), [Schwarz et al., 1978](#)].
- **Non-asymptotic criterion**: our finite-sample oracle inequality, strong justification for **slope heuristic approach** (**capushe** package) in a finite sample setting [[Birgé and Massart, 2007](#), [Baudry et al., 2012](#)].

- $L = D = 1$ : behavior of  $\text{JKL}_{\rho}^{\otimes n}(s_0^*, \widehat{s}_m^*)$  and convergence rates of error terms.
- $D \gg L$ : dimensionality reduction capability of GLLiM in high-dimensional regression data [Deleforge et al., 2015].

▮▮▮ **Well-Specified (WS):**  $s_0^* \in S_m^*$ ,

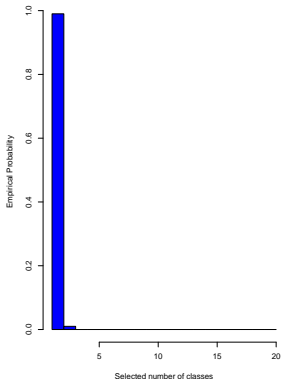
$$s_0^*(y|x) = \frac{\Phi(x; 0.2, 0.1)\Phi(y; -5x + 2, 0.09) + \Phi(x; 0.8, 0.15)\Phi(y; 0.1x, 0.09)}{\Phi(x; 0.2, 0.1) + \Phi(x; 0.8, 0.15)}.$$

▮▮▮ **Misspecified (MS):**  $s_0^* \notin S_m^*$ ,

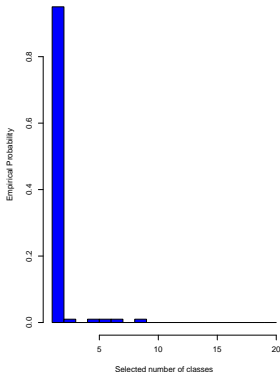
$$s_0^*(y|x) = \frac{\Phi(x; 0.2, 0.1)\Phi(y; x^2 - 6x + 1, 0.09) + \Phi(x; 0.8, 0.15)\Phi(y; -0.4x^2, 0.09)}{\Phi(x; 0.2, 0.1) + \Phi(x; 0.8, 0.15)}.$$



# Histogram of selected $K$ using slope heuristic over 100 trials



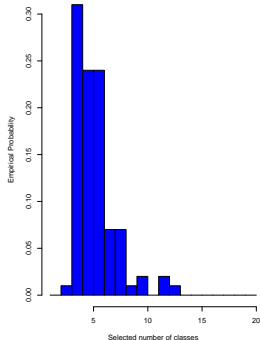
(a) 2000 data points



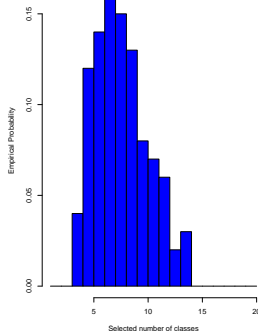
(b) 10000 data points

Comparison histograms of selected  $K$  in **WS case** using jump criterion over 100 trials between 2000 and 10000 data points.

# Histogram of selected $K$ using slope heuristic over 100 trials



(a) 2000 data points



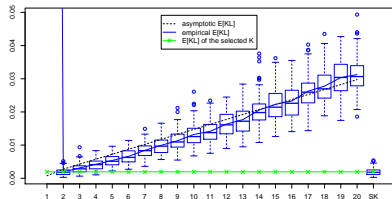
(b) 10000 data points

Comparison histograms of selected  $K$  in **MS case** using jump criterion over 100 trials between 2000 and 10000 data points.

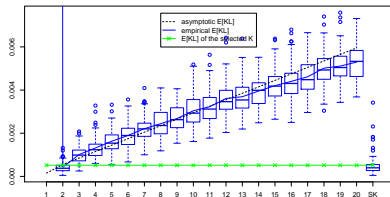
The bias-variance trade-off differs between the two examples:

- **WS case:** since the true density belongs to the model, the best choice is  $K = 2$  even for large  $n$ .
- **MS case:** best choice  $K$  should balance a model **approximation error term** and a **variance** one, *i.e.*, the larger  $n$  the more complex the model and thus  $K$ .

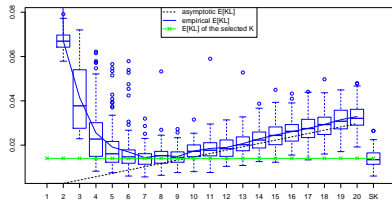
# Box-plot of Kullback-Leibler divergence over 100 trials



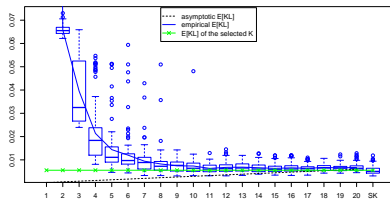
(a) WS with  $n = 2000$



(b) WS with  $n = 10000$



(c) MS with  $n = 2000$



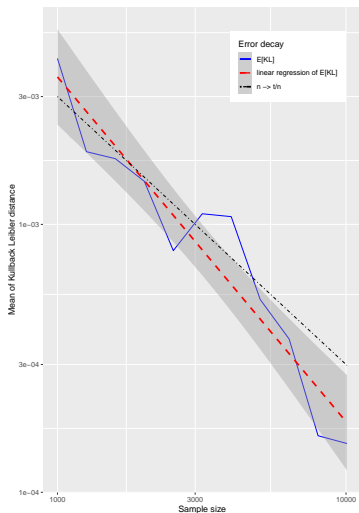
(d) MS with  $n = 10000$

# Empirical behavior of weak oracle inequality

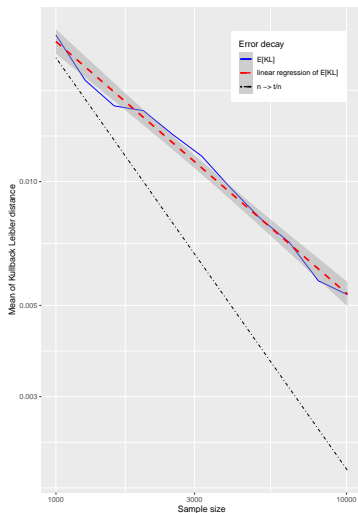
$$\mathbb{E} [\text{JKL}_{\rho}^{\otimes n}(s_0, \hat{s}_{\hat{m}})] \leq C_1 \inf_{m \in \mathcal{M}} \left( \inf_{s_m \in S_m} \text{KL}^{\otimes n}(s_0, s_m) + \frac{\text{pen}(m)}{n} \right) + \frac{\kappa(\rho, C_1) C_1 \Xi}{n} + \frac{\eta + \eta'}{n}.$$

- No known formula for  $\text{JKL}_{\rho}^{\otimes n}(s_0, \hat{s}_{\hat{m}}) \Rightarrow$  Monte Carlo method.
- Empirical mean  $\text{JKL}_{\rho}^{\otimes n}(s_0, \hat{s}_{\hat{m}}) \leq$  Empirical mean  $\text{KL}^{\otimes n}(s_0, \hat{s}_m)$ ,  $m \in \mathcal{M} = [20]$  over 55 trials.
- Empirical mean  $\text{KL}^{\otimes n}(s_0, \hat{s}_m) \sim \frac{\dim(S_m)}{2n}$  (shown by a dotted line): **expected behavior in asymptotic theory in WS case!**

# Rate of error decay in a log-log scale, using 30 trials



(a) WS: free regression's slope  $\approx -1.287$  and  $t = 3$ .



(b) MS: free regression's slope  $\approx -0.6120$ ,  $t = 20$ .

# Outline

- 1 Collection of mixture of experts models
- 2 Non-asymptotic oracle inequality
- 3 Numerical experiment
- 4 Future work on non-asymptotic oracle inequality

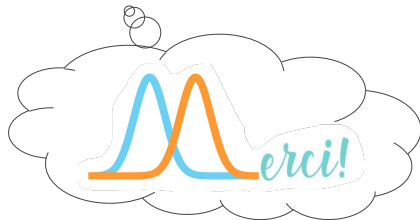
- ❶ **Minimax analysis** (lower bound) for GLoME and SGaME models (GMM [[Maugis-Rabusseau and Michel, 2013](#)]).
- ❷ **Improvement on upper bound:**  $C_1 > 1 \rightarrow C_1 = 1$ .



This is actually a realistic model selection problem...



**Thank you for  
your attention!!!**



Edit image: Dung N. Nguyen

# Universal approximation theorem of MO-MoLE models

Multiple-output Gaussian gated mixture of linear experts (MO-MoLE) models and the class of MO-continuous functions, given  $\mathcal{Y}$  is a compact set,

$$\mathcal{M}_K(\mathcal{Y}) = \left\{ \mathcal{Y} \ni \mathbf{y} \mapsto \sum_{k=1}^K \mathbf{g}_k(\mathbf{y}; \boldsymbol{\omega}) \left[ \mathbf{a}_k + \mathbf{B}_k^\top \mathbf{y} \right] \right\},$$
$$\mathcal{C}_L(\mathcal{Y}) = \left\{ \mathcal{Y} \ni \mathbf{y} \mapsto \mathbf{m}(\mathbf{y}) = (\mathbf{m}_j(\mathbf{y}))_{j \in [L]} : \mathbf{m}_j \in \mathcal{C}(\mathcal{Y}), j \in [L] \right\}.$$

Theorem ([[Nguyen et al., 2019](#)])

For all  $\mathbf{m}^0 \in \mathcal{C}_L(\mathcal{Y})$ , there exists  $\{\mathbf{m}_K\}_{K \in \mathbb{N}^*} \subset \bigcup_{K \in \mathbb{N}^*} \mathcal{M}_K(\mathcal{Y})$ ,

$$\sum_{j=1}^L \sup_{\mathbf{y} \in \mathcal{Y}} |\mathbf{m}_j^0(\mathbf{y}) - \mathbf{m}_{K,j}(\mathbf{y})| \xrightarrow{K \rightarrow \infty} 0.$$

★  $\mathcal{M}_K(\mathcal{Y}) \not\subseteq S_m!$

# Universal approximation theorem of GMMs

Given a PDF  $\varphi$  (e.g., standard multivariate normal distribution (MND)),

$$\mathcal{S}^\varphi = \bigcup_{K \in \mathbb{N}^*} \mathcal{S}_K^\varphi, \text{ where } \mathcal{S}_K^\varphi = \left\{ \mathcal{Y} \ni \mathbf{y} \mapsto s_K^\varphi(\mathbf{y}) = \sum_{k=1}^K \frac{\pi_k}{\sigma_k^L} \varphi\left(\frac{\mathbf{y} - \mathbf{v}_k}{\sigma_k}\right), \right. \\ \left. \mathbf{v}_k \in \mathbb{R}^L, \sigma_k \in \mathbb{R}^+, \boldsymbol{\pi} \in \boldsymbol{\Pi}_{K-1} \right\}.$$

Theorem ([[Nguyen et al., 2020c](#), [Nguyen et al., 2020b](#)])

- Given any PDFs  $s_0, \varphi \in \mathcal{C}$  and a compact set  $\mathcal{Y} \subset \mathbb{R}^L$ , there exists  $\{s_K^\varphi\}_{K \in \mathbb{N}^*} \subset \mathcal{S}^\varphi$ ,  $\lim_{K \rightarrow \infty} \sup_{\mathbf{y} \in \mathcal{Y}} |s_0(\mathbf{y}) - s_K^\varphi(\mathbf{y})| = 0$ .
- For  $p \in [1, \infty)$ , if  $s_0 \in \mathcal{L}_p$  and  $\varphi \in \mathcal{L}_\infty$ , there exists  $\{s_K^\varphi\}_{K \in \mathbb{N}^*} \subset \mathcal{S}^\varphi$ ,  $\lim_{K \rightarrow \infty} \|s_0 - s_K^\varphi\|_{\mathcal{L}_p} = 0$ .

★  $\mathcal{S}_K^\varphi \not\subset \mathcal{S}_m!$

# Essentially bounded and Lebesgue conditional PDF

## Definition

- Essentially bounded function on  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ :

$$\mathcal{L}_\infty(\mathcal{Z}) = \left\{ f : \underbrace{\inf \{a \geq 0 : \lambda(\{\mathbf{z} \in \mathcal{Z} : |f(\mathbf{z})| > a\}) = 0\}}_{:= \|f\|_{\infty, \mathcal{Z}}} < \infty \right\}.$$

- Lebesgue conditional PDF:  $\mathcal{F}_p = \mathcal{F} \cap \mathcal{L}_p$ ,  $p \in [1, \infty)$ ,

$$\mathcal{F} = \left\{ f : \mathcal{Z} \rightarrow [0, \infty), \int_{\mathcal{Y}} f(\mathbf{x}|\mathbf{y}) d\lambda(\mathbf{x}) = 1 \right\},$$
$$\mathcal{L}_p(\mathcal{Z}) = \left\{ f := \underbrace{\left( \int_{\mathcal{Z}} |f(\mathbf{z})|^p d\lambda(\mathbf{z}) \right)^{1/p}}_{:= \|f\|_{p, \mathcal{Z}}} < \infty \right\}.$$

# Approximation class: isotropic SGaME and GLLiM

- Location-scale family: given a PDF  $\varphi$ ,  $\mathbf{v} \in \mathcal{X}$ ,  $\sigma \in R^+$ ,

$$\mathcal{E}_\varphi = \left\{ \mathbf{x} \mapsto \frac{1}{\sigma^D} \varphi \left( \frac{\mathbf{x} - \mathbf{v}}{\sigma} \right) = \Phi_D(\mathbf{x}; \mathbf{v}, \sigma) \right\}.$$

- Isotropic **SGaME** models:

$$\mathcal{S}_S^\varphi = \left\{ (\mathbf{x}, \mathbf{y}) \mapsto s_K^\varphi(\mathbf{x}|\mathbf{y}) = \sum_{k=1}^K \mathbf{g}_k(\mathbf{y}; \boldsymbol{\gamma}) \Phi_D(\mathbf{x}; \mathbf{v}_k, \sigma_k), \right. \\ \left. \Phi_D \in \mathcal{E}_\varphi \cap \mathcal{L}_\infty, \mathbf{g}_k(\cdot; \boldsymbol{\gamma}) \in \mathcal{P}_S^K, K \in \mathbb{N}^* \right\}.$$

- Isotropic **GLLiM** model ( $\subset S_m$ ) when  $\varphi$  is standard MND:

$$\mathcal{S}_G^\varphi = \left\{ (\mathbf{x}, \mathbf{y}) \mapsto s_K^\varphi(\mathbf{x}|\mathbf{y}) = \sum_{k=1}^K \mathbf{g}_k(\mathbf{y}; \boldsymbol{\omega}) \varphi(\mathbf{x}; \mathbf{v}_k, \sigma_k), \right. \\ \left. \Phi_D \in \mathcal{E}_\varphi \cap \mathcal{L}_\infty, \mathbf{g}_k(\cdot; \boldsymbol{\omega}) \in \mathcal{P}_G^K, K \in \mathbb{N}^* \right\}.$$

## Theorem ([[Nguyen et al., 2020a](#)])

- (a) Given  $\varphi \in \mathcal{F} \cap \mathcal{C}$ , for any target  $s_0 \in \mathcal{F}_p \cap \mathcal{C}_b^u$ , there exist sequences  $\{s_K^\varphi\}_{K \in \mathbb{N}^*} \subset \mathcal{S}_S^\varphi$  and  $\{s_K'^\varphi\}_{K \in \mathbb{N}^*} \subset \mathcal{S}_G^\varphi$  such that

$$\lim_{K \rightarrow \infty} \|s_0 - s_K^\varphi\|_{\mathcal{L}_p} = 0,$$

$$\lim_{K \rightarrow \infty} \|s_0 - s_K'^\varphi\|_{\mathcal{L}_p} = 0.$$

- (b) Given  $\varphi \in \mathcal{F} \cap \mathcal{C}_b^u$ , for any target  $s_0 \in \mathcal{F} \cap \mathcal{C}_b^u$ ,  $L = 1$ , and  $0 < \lambda(\mathcal{X}) < \infty$ , there exist sequences  $\{s_K^\varphi\}_{K \in \mathbb{N}^*} \subset \mathcal{S}_S^\varphi$  and  $\{s_K'^\varphi\}_{K \in \mathbb{N}^*} \subset \mathcal{S}_G^\varphi$  such that

$$\lim_{K \rightarrow \infty} s_K^\varphi = s_0 \text{ almost uniformly,}$$

$$\lim_{K \rightarrow \infty} s_K'^\varphi = s_0 \text{ almost uniformly.}$$



Akaike, H. (1974).

A new look at the statistical model identification.

*IEEE transactions on automatic control*, 19(6):716–723.

(Cited on page [31](#).)



Baudry, J.-P., Maugis, C., and Michel, B. (2012).

Slope heuristics: overview and implementation.

*Statistics and Computing*, 22(2):455–470.

(Cited on page [31](#).)



Birgé, L. and Massart, P. (2007).

Minimal penalties for Gaussian model selection.

*Probability theory and related fields*, 138(1-2):33–73.

(Cited on page [31](#).)



Boucheron, S. and Massart, P. (2011).  
A high-dimensional Wilks phenomenon.  
*Probability Theory and Related Fields*, 150(3):405–433.  
(Cited on page 25.)



Celeux, G. and Govaert, G. (1995).  
Gaussian parsimonious clustering models.  
*Pattern recognition*, 28(5):781–793.  
(Cited on page 15.)



Cohen, S. and Pennec, E. L. (2011).  
Conditional density estimation by penalized likelihood model selection  
and applications.  
*Technical report, INRIA*.  
(Cited on pages 23, 24, and 28.)





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

*Statistics and Computing*, 25(5):893–911.

(Cited on pages 10 and 32.)



Fan, J., Zhang, C., and Zhang, J. (2001).  
Generalized Likelihood Ratio Statistics and Wilks Phenomenon.

*The Annals of Statistics*, 29(1):153–193.




(Cited on page 25.)



Ho, N., Yang, C.-Y., and Jordan, M. I. (2019).  
Convergence Rates for Gaussian Mixtures of Experts.

*arXiv preprint arXiv:1907.04377*.

(Cited on page 10.)

-  Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. (1991).  
Adaptive Mixtures of Local Experts.  
*Neural Computation*, 3:79–87.  
(Cited on page 10.)
-  Maugis-Rabusseau, C. and Michel, B. (2013).  
Adaptive density estimation for clustering with Gaussian mixtures.  
*ESAIM: Probability and Statistics*, 17:698–724.  
(Cited on pages 26 and 40.)
-  McLachlan, G. J. and Peel, D. (2000).  
*Finite Mixture Models*.  
John Wiley & Sons.  
(Cited on page 10.)



Mendes, E. F. and Jiang, W. (2012).  
On Convergence Rates of Mixtures of Polynomial Experts.  
*Neural Computation*, 24(11):3025–3051.  
(Cited on page 10.)



Montuelle, L. and Pennec, E. L. (2014).  
Mixture of Gaussian regressions model with logistic weights, a  
penalized maximum likelihood approach.  
*Electronic Journal of Statistics*, 8(1):1661–1695.  
(Cited on pages 14, 21, and 22.)



Nguyen, H. D. and Chamroukhi, F. (2018).  
Practical and theoretical aspects of mixture-of-experts modeling: An  
overview.  
*WIREs Data Mining and Knowledge Discovery*, 8(4):e1246.  
(Cited on page 10.)



Nguyen, H. D., Chamroukhi, F., and Forbes, F. (2019).

Approximation results regarding the multiple-output Gaussian gated mixture of linear experts model.

*Neurocomputing*, 366:208–214.

(Cited on pages 29 and 42.)



Nguyen, H. D., Nguyen, T., Chamroukhi, F., and McLachlan, G. (2020a).

Approximations of conditional probability density functions in Lebesgue spaces via mixture of experts models.

*arXiv preprint arXiv:2012.02385*.

(Cited on pages 10, 21, 29, and 46.)



Nguyen, T., Chamroukhi, F., Nguyen, H. D., and McLachlan, G. J. (2020b).

Approximation of probability density functions via location-scale finite mixtures in Lebesgue spaces.

*arXiv preprint arXiv:2008.09787.*

(Cited on pages 29 and 43.)



Nguyen, T. T., Nguyen, H. D., Chamroukhi, F., and Forbes, F. (2021).

A non-asymptotic penalization criterion for model selection in mixture of experts models.

*arXiv preprint arXiv:2104.02640.*

(Cited on page 20.)

# References VIII



Nguyen, T. T., Nguyen, H. D., Chamroukhi, F., and McLachlan, G. J. (2020c).

Approximation by finite mixtures of continuous density functions that vanish at infinity.

*Cogent Mathematics & Statistics*, 7(1):1750861.

(Cited on pages 29 and 43.)



Schwarz, G. et al. (1978).

Estimating the dimension of a model.

*The annals of statistics*, 6(2):461–464.

(Cited on page 31.)



White, H. (1982).

Maximum Likelihood Estimation of Misspecified Models.

*Econometrica*, 50(1):1–25.

(Cited on pages 23 and 24.)



Wilks, S. S. (1938).

The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses.

*The Annals of Mathematical Statistics*, 9(1):60–62.

(Cited on page 25.)



Xu, L., Jordan, M., and Hinton, G. E. (1995).

An Alternative Model for Mixtures of Experts.

In Tesauro, G., Touretzky, D., and Leen, T., editors, *Advances in Neural Information Processing Systems*, volume 7. MIT Press.

(Cited on page 10.)