

Mixture of Gaussian Regressions with logistic weights and conditional density estimation

E. Le Pennec

(CMAP - École Polytechnique)

and

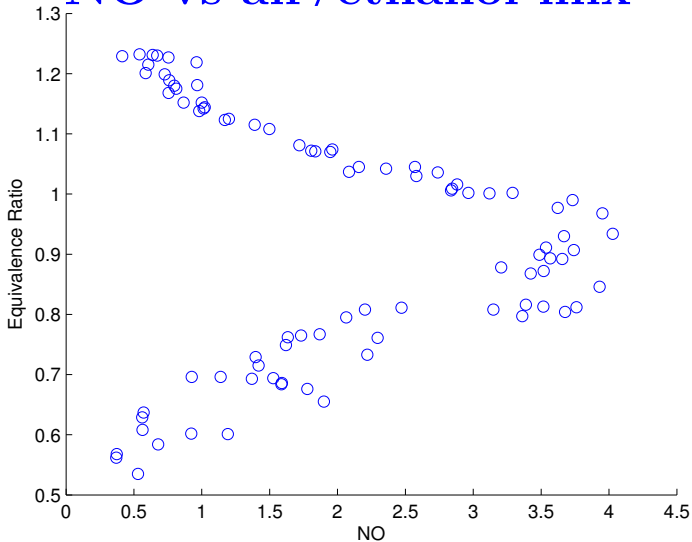
L. Montuelle (SELECT - Inria Saclay / Université Paris Sud)

Alea

11 march 2014

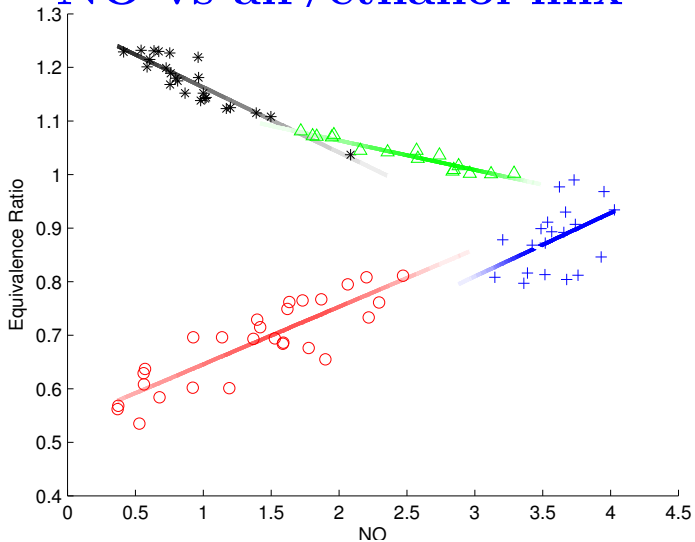


NO vs air/ethanol mix



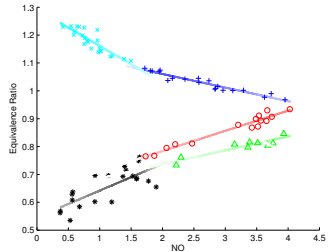
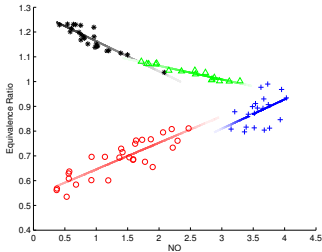
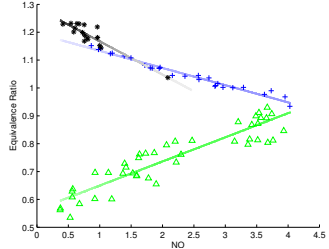
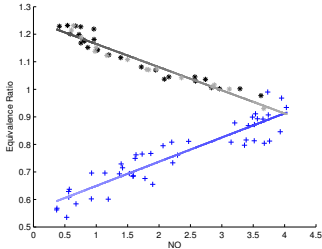
● Regression

NO vs air/ethanol mix



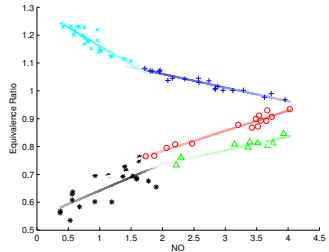
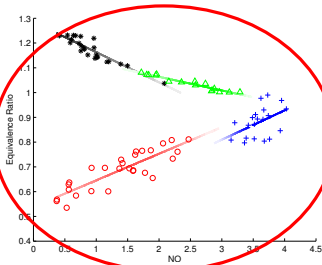
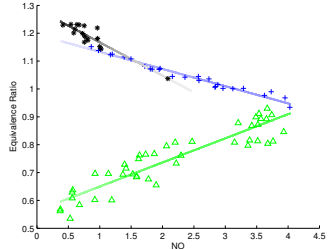
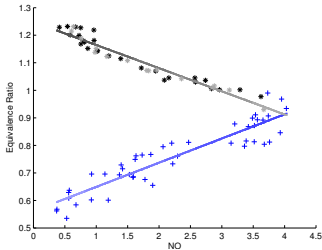
● Regression Mixture

NO vs air/ethanol mix



- Regression Mixture
- Mixtures of Gaussian regressions with logistic weights

NO vs air/ethanol mix



- Regression Mixture
- Mixtures of Gaussian regressions with logistic weights
- Model selection

Statistical modeling

- **Data:** $(X_i, Y_i)_{i \leq n} \in [0; 1]^d \times \mathbb{R}^p$

- $X_i \perp X_j$
- $Y_i | (X_k)_k \perp Y_j | (X_k)_k$
- $Y | X$ has a density s_0 w.r.t. Lebesgue measure

- **Regression** = specific modeling of the conditional density $s_0(\cdot | x)$

- **Gaussian regression mixture with logistic weights:**

$$s_{\mathbf{K}, \mathbf{v}, \mathbf{\Sigma}, \mathbf{w}}(y | x) = \sum_{k=1}^{\mathbf{K}} \pi_{\mathbf{w}, k}(x) \Phi_{\mathbf{v}_k(x), \mathbf{\Sigma}_k(x)}(y),$$

- with
- $\pi_{\mathbf{w}, k}(x) = \frac{e^{\mathbf{w}_k(x)}}{\sum_{k'=1}^{\mathbf{K}} e^{\mathbf{w}_{k'}(x)}}$, logistic weights
 - $\Phi_{\mathbf{v}_k(x), \mathbf{\Sigma}_k(x)}$ density of $\mathcal{N}(\mathbf{v}_k(x), \mathbf{\Sigma}_k(x))$

- **Parameters:**

- \mathbf{K} : number of components
- \mathbf{v} and $\mathbf{\Sigma}$: \mathbf{K} regression functions and covariance matrices functions
- \mathbf{w} : \mathbf{K} weights functions defining the mixture proportions

Models

- **Gaussian regression mixture with logistic weights:**

$$s_{K,v,\Sigma,w}(y|x) = \sum_{k=1}^K \pi_{w,k}(x) \Phi_{v_k(x),\Sigma_k(x)}(y),$$

with

- $\pi_{w,k}(x) = \frac{e^{w_k(x)}}{\sum_{k'=1}^K e^{w_{k'}(x)}}$, logistic weights
- $\Phi_{v_k(x),\Sigma_k(x)}$ density of $\mathcal{N}(v_k(x), \Sigma_k(x))$

- **Parameters:** $\theta = (K, v, \Sigma, w)$

- K : number of components
- v and Σ : K regression functions and covariance matrices functions
- w : K weights functions defining the mixture proportions

- **Model** $S_m = \{s_\theta, \theta \in \Theta_m\}$ with $\Theta_m = \{K\} \otimes \mathcal{T}_K \otimes \mathcal{V}_K \otimes \mathcal{W}_K$:

- K : number of components.
- \mathcal{T}_K and \mathcal{V}_K : sets for the K -tuple of regressions functions and covariance matrices functions.
- \mathcal{W}_K : sets for for the K -tuple of weights functions.

- **Typical choice:**

- \mathcal{T}_K and \mathcal{W}_K : tensorial product of polynomial sets of low degree.
- \mathcal{W}_K : constant covariance structures independent of X .

Maximum likelihood and penalization

- **Model** $S_m = \{s_\theta, \theta \in \Theta_m\}$ with $\Theta_m = \{K\} \otimes \mathcal{T}_K \otimes \mathcal{V}_K \otimes \mathcal{W}_K$:
 - K : number of components.
 - \mathcal{T}_K and \mathcal{V}_K : sets for the K -tuple of regressions functions and covariance matrices functions.
 - \mathcal{W}_K : sets for for the K -tuple of weights functions.
- **Maximum likelihood estimation** within each model:

$$\hat{s}_m = \operatorname{argmax}_{\theta \in \Theta_m} - \sum_{i=1}^n \ln s_\theta(Y_i|X_i)$$

- **Model selection** by a penalization proportional to the dimension:

$$\hat{m} = \operatorname{argmin}_{m \in \mathcal{M}} \sum_{k=1}^K - \ln \hat{s}_m(Y_i|X_i) + \kappa \dim \Theta_m$$

- Usual complexity/fidelity tradeoff.

Contributions

Characterization of the theoretical performances

- **Penalty choice:** $\text{pen}(m) = \kappa(C + \ln n) \dim(S_m)$.
- **Oracle inequality:**

$$\mathbb{E} \left[JKL_{\rho}^{\otimes n}(s_0, \widehat{s}_{\widehat{m}}) \right] \leq C_1 \inf_{m \in \mathcal{M}} \left(\inf_{s_m \in S_m} KL^{\otimes n}(s_0, s_m) + \frac{\text{pen}(m)}{n} \right) + \frac{C_2}{n}$$

Numerical implementation of the penalized maximum likelihood

- **EM type minimization scheme** with a focus on initialization issues.
- Practical scheme for the **penalty calibration** with the slope heuristic approach.

Conditional density and selection

- **General framework:** observation of (X_i, Y_i) with X_i independent and Y_i cond. independent of law of density $s_0(y|X_i)$.
- **Goal:** estimation of $s_0(y|x)$.
- **Penalized model selection principle:**
 - choice of a collection of cond. dens. models $S_m = \{s_m(y|x)\}$ with $m \in \mathcal{S}$,
 - Maximum likelihood estimation of a cond. density \hat{s}_m for each model S_m :

$$\hat{s}_m = \operatorname{argmin}_{s_m \in S_m} - \sum_{i=1}^n \ln s_m(Y_i|X_i)$$

- Selection of a model \hat{m} by
$$\hat{m} = \operatorname{argmin}_{m \in \mathcal{S}} - \sum_{i=1}^n \ln \hat{s}_m(Y_i|X_i) + \operatorname{pen}(m).$$

with $\operatorname{pen}(m)$ well chosen.

- **Typical oracle inequality result:**

$$\mathbb{E} \left[d^2(s_0, \hat{s}_{\hat{m}}) \right] \leq C \inf_{m \in \mathcal{S}} \left(\inf_{s_m \in S_m} KL(s_0, s_m) + \frac{\operatorname{pen}(m)}{n} \right) + \frac{C'}{n}.$$

- **Short bib.:** Rosenblatt, Fan et al., de Gooijer and Zerom, Efromovitch, Brunel, Comte, Lacour... / Plugin, direct estimation, L^2 , minimax, censure...

Ideal oracle inequality

- **Oracle inequality:**

$$\mathbb{E} [KL^{\otimes n}(s_0, \widehat{s}_m)] \leq C_1 \inf_{s_m \in \mathcal{S}} \left(\underbrace{\inf_{s_m \in \mathcal{S}_m} KL^{\otimes n}(s_0, s_m)}_{\text{Bias term}} + \underbrace{\frac{\Delta(m)}{n}}_{\text{Variance term}} \right)$$

as soon as $\text{pen}(m)$ is *large enough*

- **Divergence adapted to the conditional density setting:**

- Divergence on the product density conditioned on the design (Kolaczyk, Bigot).
- *Tensorization* principle and expectation on the design: design:

$$KL \rightarrow KL^{\otimes n}(s, s') = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n KL(s(\cdot|X_i), s'(\cdot|X_i)) \right]$$

- Much **more information** using the second approach because losses used are *larger*.
- Ability to handle **independent but non i.i.d. case** and **integrated loss**.
- Classical density estimation theorem if $s(\cdot|X_i) = s(\cdot)$.

Notations

- Let for any function $g(x, y)$,

- $P_n^{\otimes n}(g)$: its **empirical process** $P_n^{\otimes n}(g) = \frac{1}{n} \sum_{i=1}^n g(X_i, Y_i)$.

- $P^{\otimes n}(g)$: its **expectation** $P^{\otimes n}(g) = \mathbb{E} [P_n^{\otimes n}(g)] = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n g(X_i, Y_i) \right]$.

- $\nu_n^{\otimes n}(g) = P_n^{\otimes n}(g) - P^{\otimes n}(g)$: its **recentered** process.

- Maximum likelihood estimate:**

$$\begin{aligned}\hat{s}_m &= \operatorname{argmin}_{s_m \in S_m} \sum_{i=1}^n -\ln s_m(Y_i | X_i) = \operatorname{argmin}_{s_m \in S_m} P_n^{\otimes n}(-\ln s_m) \\ &= \operatorname{argmin}_{s_m \in S_m} P_n^{\otimes n}\left(-\ln \frac{s_m}{s_0}\right)\end{aligned}$$

- Best projection:**

$$\begin{aligned}\tilde{s}_m &= \operatorname{argmin}_{s_m \in S_m} KL^{\otimes n}(s_0, s_m) = \operatorname{argmin}_{s_m \in S_m} P^{\otimes n}\left(-\ln \frac{s_m}{s_0}\right) \\ &= \operatorname{argmin}_{s_m \in S_m} P^{\otimes n}(-\ln s_m)\end{aligned}$$

Ideal penalty

- By definition:

$$KL^{\otimes n}(s_0, \hat{s}_m) = P_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right) - \underbrace{\nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right)}_{\text{pen}_{\text{id}}(m)/n}$$

- With the **ideal penalty** $\text{pen}_{\text{id}}(m)$:

$$\begin{aligned} KL^{\otimes n}(s_0, \hat{s}_{\hat{m}}) &= P_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) + \frac{\text{pen}_{\text{id}}(\hat{m})}{n} \\ &\leq \inf_m P_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right) + \frac{\text{pen}_{\text{id}}(m)}{n} \leq \inf_m KL^{\otimes n}(s_0, \hat{s}_m) \\ &\leq \inf_m (KL^{\otimes n}(s_0, \tilde{s}_m) + (KL^{\otimes n}(s_0, \hat{s}_m) - KL^{\otimes n}(s_0, \tilde{s}_m))) \end{aligned}$$

- Ideal penalty oracle inequality:**

$$\mathbb{E} [KL^{\otimes n}(s_0, \hat{s}_{\hat{m}})] \leq \inf_{S_m \in \mathcal{S}} \left(\underbrace{KL^{\otimes n}(s_0, \tilde{s}_m)}_{\text{Bias term}} + \underbrace{\mathbb{E} [KL^{\otimes n}(s_0, \hat{s}_m) - KL^{\otimes n}(s_0, \tilde{s}_m)]}_{\text{Variance term}} \right)$$

Non ideal penalization

- By construction

$$\begin{aligned} KL^{\otimes n}(s_0, \hat{s}_{\hat{m}}) &= P_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) \\ &= P_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) + \frac{\text{pen}(\hat{m})}{n} \\ &\quad - \nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \frac{\text{pen}(\hat{m})}{n} \\ &\leq \min_{S_m \in \mathcal{S}} \left(P_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right) + \frac{\text{pen}(m)}{n} \right) \\ &\quad - \nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \frac{\text{pen}(\hat{m})}{n} \end{aligned}$$

Non ideal penalization

- Using $\tilde{s}_m = \operatorname{argmin}_{s_m \in \mathcal{S}_m} KL^{\otimes n}(s_0, s_m)$:

$$KL^{\otimes n}(s, \hat{s}_{\hat{m}}) \leq \min_{s_m \in \mathcal{S}} \left(P_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right) + \frac{\operatorname{pen}(m)}{n} \right) \\ - \nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \frac{\operatorname{pen}(\hat{m})}{n}$$

$$KL^{\otimes n}(s_0, \hat{s}_{\hat{m}}) \leq \min_{s_m \in \mathcal{S}} \left(P_n^{\otimes n} \left(-\ln \frac{\tilde{s}_m}{s_0} \right) + \frac{\operatorname{pen}(m)}{n} \right) \\ - \nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \frac{\operatorname{pen}(\hat{m})}{n}$$

Non ideal penalization

- **Summary:**

$$KL^{\otimes n}(s_0, \hat{s}_{\hat{m}}) \leq \min_{s_m \in \mathcal{S}} \left(P_n^{\otimes n} \left(-\ln \frac{\tilde{s}_m}{s_0} \right) + \frac{\text{pen}(m)}{n} \right) \\ - \nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \frac{\text{pen}(\hat{m})}{n}$$

- **Oracle inequality** up to something:

$$\mathbb{E} [KL^{\otimes n}(s_0, \hat{s}_{\hat{m}})] \leq \min_{s_m \in \mathcal{S}} \left(KL^{\otimes n}(s_0, \tilde{s}_m) + \frac{\text{pen}(m)}{n} \right) \\ + \mathbb{E} \left[-\nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \frac{\text{pen}(\hat{m})}{n} \right]$$

- If $\mathbb{E} \left[-\nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \frac{\text{pen}(\hat{m})}{n} \right] \leq 0$ then **exact** oracle inequality!
- If $\mathbb{E} \left[-\nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \frac{\text{pen}(\hat{m})}{n} - \epsilon KL^{\otimes n}(s_0, \hat{s}_{\hat{m}}) \right] \leq 0$ then **inexact** oracle inequality.

Kullback-Leibler and extension

- **Issue** in the previous approach: control of

$$\nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right)$$

hard due to the **unboundedness** of $-\ln \frac{\hat{s}_m}{s_0}$

- **Trick:** replace this quantity by the bounded one

$$-\frac{1}{\rho} \ln \frac{\rho \hat{s}_m + (1 - \rho)s_0}{s_0} \leq -\frac{1}{\rho} \ln(1 - \rho)$$

- By convexity,

$$-\frac{1}{\rho} \ln \frac{\rho \hat{s}_m + (1 - \rho)s_0}{s_0} \leq -\ln \frac{\hat{s}_m}{s_0}$$

- **Jensen-Kullback-Leibler divergence:**

$$\begin{aligned} JKL_{\rho}^{\otimes n}(s_0, \hat{s}_m) &= P^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \hat{s}_m + (1 - \rho)s_0}{s_0} \right) = \frac{1}{\rho} KL^{\otimes n}(s_0, \rho \hat{s}_m + (1 - \rho)s_0) \\ &\leq P^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right) = KL^{\otimes n}(s_0, \hat{s}_m) \end{aligned}$$

$JKL^{\otimes n}$ and non ideal penalization

- By construction

$$\begin{aligned} JKL_{\rho}^{\otimes n}(s_0, \hat{s}_{\hat{m}}) &= P_n^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \hat{s}_{\hat{m}} + (1-\rho)s_0}{s_0} \right) \\ &\quad - \nu_n^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \hat{s}_{\hat{m}} + (1-\rho)s_0}{s_0} \right) \\ &\leq P_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) - \nu_n^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \hat{s}_{\hat{m}} + (1-\rho)s_0}{s_0} \right) \\ &\leq P_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) + \frac{\text{pen}(\hat{m})}{n} \\ &\quad - \nu_n^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \hat{s}_{\hat{m}} + (1-\rho)s_0}{s_0} \right) - \frac{\text{pen}(\hat{m})}{n} \\ &\leq \min_{S_m \in \mathcal{S}} \left(P_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right) + \frac{\text{pen}(m)}{n} \right) \\ &\quad - \nu_n^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \hat{s}_{\hat{m}} + (1-\rho)s_0}{s_0} \right) - \frac{\text{pen}(\hat{m})}{n} \end{aligned}$$

$JKL^{\otimes n}$ and non ideal penalization

- Using $\tilde{s}_m = \operatorname{argmin}_{s_m \in \mathcal{S}_m} KL^{\otimes n}(s_0, s_m)$:

$$JKL_{\rho}^{\otimes n}(s_0, \hat{s}_{\hat{m}}) \leq \min_{s_m \in \mathcal{S}} \left(P_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right) + \frac{\operatorname{pen}(m)}{n} \right) \\ - \nu_n^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \hat{s}_{\hat{m}} + (1-\rho)s_0}{s_0} \right) - \frac{\operatorname{pen}(\hat{m})}{n}$$

$$JKL_{\rho}^{\otimes n}(s_0, \hat{s}_{\hat{m}}) \leq \min_{s_m \in \mathcal{S}} \left(P_n^{\otimes n} \left(-\ln \frac{\tilde{s}_m}{s_0} \right) + \frac{\operatorname{pen}(m)}{n} \right) \\ - \nu_n^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \hat{s}_{\hat{m}} + (1-\rho)s_0}{s_0} \right) - \frac{\operatorname{pen}(\hat{m})}{n}$$

$JKL^{\otimes n}$ and non ideal penalization

- **Summary:**

$$JKL_{\rho}^{\otimes n}(s_0, \widehat{s}_{\widehat{m}}) \leq \min_{S_m \in \mathcal{S}} \left(P_n^{\otimes n} \left(-\ln \frac{\widetilde{s}_m}{s_0} \right) + \frac{\text{pen}(m)}{n} \right) \\ - \nu_n^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \widehat{s}_{\widehat{m}} + (1 - \rho)s_0}{s_0} \right) - \frac{\text{pen}(\widehat{m})}{n}$$

- **Oracle inequality up to something:**

$$\mathbb{E} \left[JKL_{\rho}^{\otimes n}(s_0, \widehat{s}_{\widehat{m}}) \right] \leq \min_{S_m \in \mathcal{S}} \left(KL^{\otimes n}(s_0, \widetilde{s}_m) + \frac{\text{pen}(m)}{n} \right) \\ + \mathbb{E} \left[-\nu_n^{\otimes n} \left(-\frac{1}{\rho} \ln \frac{\rho \widehat{s}_{\widehat{m}} + (1 - \rho)s_0}{s_0} \right) - \frac{\text{pen}(\widehat{m})}{n} \right]$$

- Under some assumptions on the model collection, **it exists a penalty such that**

$$\mathbb{E} \left[-\nu_n^{\otimes n} \left(-\ln \frac{\widehat{s}_{\widehat{m}}}{s_0} \right) - \frac{\text{pen}(\widehat{m})}{n} - \epsilon JKL_{\rho}^{\otimes n}(s_0, \widehat{s}_{\widehat{m}}) \right] \leq 0$$

- For such a penalty, one has an **inexact** oracle inequality.

Theorem

Assumption (H): For every model S_m in the collection \mathcal{S} , there is a non-decreasing function $\phi_m(\delta)$ such that $\delta \mapsto \frac{1}{\delta}\phi_m(\delta)$ is non-increasing on $(0, +\infty)$ and for every $\sigma \in \mathbb{R}^+$ and every $s_m \in S_m$

$$\int_0^\sigma \sqrt{H_{[\cdot], d^{\otimes n}}(\epsilon, S_m(s_m, \sigma))} d\epsilon \leq \phi_m(\sigma).$$

Assumption (K): There is a family $(x_m)_{m \in \mathcal{M}}$ of non-negative number such that

$$\sum_{m \in \mathcal{M}} e^{-x_m} \leq \Sigma < +\infty$$

Theorem

Assume we observe (X_i, Y_i) with unknown conditional s_0 . Let $\mathcal{S} = (S_m)_{m \in \mathcal{M}}$ a at most countable collection of conditional density sets. Assume Assumptions (H), (K) and (S) hold.

Let \hat{s}_m be a δ -log-likelihood minimizer in S_m :

$$\sum_{i=1}^n -\ln(\hat{s}_m(Y_i|X_i)) \leq \inf_{s_m \in S_m} \left(\sum_{i=1}^n -\ln(s_m(Y_i|X_i)) \right) + \delta$$

Then for any $\rho \in (0, 1)$ and any $C_1 > 1$, there is a constant κ_0 depending only on ρ and C_1 such that, as soon as for every index $m \in \mathcal{M}$ $\text{pen}(m) \geq \kappa(\mathfrak{D}_m + x_m)$ with $\kappa > \kappa_0$

where $\mathfrak{D}_m = n\sigma_m^2$ with σ_m the unique root of $\frac{1}{\sigma}\phi_m(\sigma) = \sqrt{n}\sigma$,

the penalized likelihood estimate $\hat{s}_{\hat{m}}$ with \hat{m} defined by

$$\hat{m} = \underset{m \in \mathcal{M}}{\text{argmin}} \sum_{i=1}^n -\ln(\hat{s}_m(Y_i|X_i)) + \text{pen}(m)$$

satisfies $\mathbb{E} \left[JKL_{\rho}^{\otimes n}(s_0, \hat{s}_{\hat{m}}) \right] \leq C_1 \left(\inf_{S_m \in \mathcal{S}} \left(\inf_{s_m \in S_m} K L^{\otimes n}(s_0, s_m) + \frac{\text{pen}(m)}{n} \right) + \frac{\kappa_0 \Sigma + \delta}{n} \right).$

Simplified Theorem...

- **Oracle inequality:**

$$\mathbb{E} \left[JKL_{\rho}^{\otimes n}(s_0, \widehat{s}_m) \right] \leq C_1 \left(\inf_{S_m \in \mathcal{S}} \left(\inf_{s_m \in S_m} KL^{\otimes n}(s_0, s_m) + \frac{\text{pen } m}{n} \right) + \frac{\kappa_0 \Sigma + \delta}{n} \right)$$

as soon as

$$\text{pen}(m) \geq \kappa (\mathfrak{D}_m + x_m) \quad \text{with } \kappa > \kappa_0,$$

where \mathfrak{D}_m measure the **complexity of the model** S_m (entropy term) and x_m the **coding cost** within the collection.

- \mathfrak{D}_m linked to the **bracketing entropy** of S_m with respect to the tensorized Hellinger distance $d^{2 \otimes n}$.
- Often $\mathfrak{D}_m \propto (\log n) \dim(S_m) \dots$

Penalty and complexities

- Control required on

$$-\nu_n^{\otimes n} \left(-\ln \frac{\widehat{s}_m}{s_0} \right) - \frac{\text{pen}(\widehat{m})}{n} - \epsilon JKL_{\rho}^{\otimes n}(s_0, \widehat{s}_m)$$

through a **supremum**!

- **Control in expectation** requires a $\text{pen}(m)$ taking into account
 - the **intrinsic complexity of the model**,
 - the **complexity of the collection**.
- Here:
 - **Model complexity**: entropy complexity \mathfrak{D}_m defined from the *bracketing entropy* $H_{[\cdot], d^{\otimes n}}(\epsilon, S_m)$ of S_m with respect to the tensorized Hellinger distance $d^{2^{\otimes n}}$.
 - **Collection (coding)**: Kraft type inequality $\sum_{m \in \mathcal{S}} e^{-x_m} \leq \Sigma < +\infty$
- **Classical constraint on the penalty**

$$\text{pen}(m) \geq \kappa (\mathfrak{D}_m + x_m) \quad \text{with } \kappa > \kappa_0.$$

- Often $\mathfrak{D}_m \propto (\ln(n)) \dim(S_m)$ and thus **classical penalization by dimension** setting...

Brackets and complexity

- **Bracketing entropy:** $H_{[\cdot], d^{\otimes n}}(\epsilon, S) =$ logarithm of the minimum number of brackets $[t_i^-, t_i^+]$ such that

- $\forall i, d^{\otimes n}(t_i^-, t_i^+) \leq \epsilon$

- $\forall s \in S, \exists i, t_i^- \leq s \leq t_i^+$

where $d^{\otimes n} = \sqrt{d^{2 \otimes n}} = \sqrt{\mathbb{E} \left[\frac{1}{n} \sum d^2(s(\cdot|X_i), s'(\cdot|X_i)) \right]}$ is the tensorized Hellinger distance.

- **Assumption (H):** for all model S_m , there is a non decreasing $\phi_m(\delta)$ such that $\delta \mapsto \frac{1}{\delta} \phi_m(\delta)$ is non increasing $(0, +\infty)$ and such that for all $\sigma \in \mathbb{R}^+$ and all $s_m \in S_m$

$$\int_0^\sigma \sqrt{H_{[\cdot], d^{\otimes n}}(\epsilon, S_m)} d\epsilon \leq \phi_m(\sigma),$$

- **Complexity** \mathcal{D}_m def. as $n\sigma_m^2$ with σ_m unique root of $\phi_m(\sigma) = \sqrt{n}\sigma^2$.
- **Key: Dudley type integral** and optimization of a deviation bound.
- Typically, $H_{[\cdot], d^{\otimes n}}(\epsilon, S_m) \sim \dim S_m (C + \log 1/\epsilon)$ which implies $\mathcal{D}_m \propto (\ln n) \dim(S_m) \dots$

Gaussian regression mixtures

- **Model** $S_m = \{s_\theta, \theta \in \Theta_m\}$ with $\Theta_m = \{K\} \otimes \Upsilon_K \otimes V_K \otimes W_K$:
 - Υ_K and V_K : sets for the K -tuple of regressions functions and covariance matrices functions.
 - W_K : sets for for the K -tuple of weights functions.
- **Structural assumptions:**
 - V_K is a set of covariance matrices independent of the covariate,
 - Υ_K and W_K are such that

$$H_{\max_{k=1}^K \sup_x \|\cdot_k(x)\|}(\delta, W_K) \leq \dim(W_K) \left(C_W + \ln \frac{1}{\delta} \right)$$

$$H_{\max_{k=1}^K \sup_x \|\cdot_k(x)\|_2}(\delta, \Upsilon_K) \leq \dim(\Upsilon_K) \left(C_\Upsilon + \ln \frac{1}{\delta} \right)$$

- Satisfied for instance if Υ_K and W_K are K -tuples of **polynomials with bounded coefficients** and x is bounded.
- **Th:** Under this assumption, if $\text{pen}(m) = \kappa(C + \ln n) \dim(S_m)$ then

$$\mathbb{E} \left[JKL_{\rho}^{\otimes n}(s_0, \hat{s}_m) \right] \leq C_1 \inf_{m \in \mathcal{M}} \left(\inf_{s_m \in S_m} KL^{\otimes n}(s_0, s_m) + \frac{\text{pen}(m)}{n} \right) + \frac{C_2}{n}$$

- **Key:** upper bound of the bracketing entropy $H_{[\cdot], d^{\otimes n}}(\epsilon, S_m)$.

Bracketing entropy decomposition

- **Model:**

$$S_m = \left\{ \sum_{k=1}^K \pi_{w,k}(x) \Phi_{v_k(x), \Sigma_k(x)}(y), (K, v, \Sigma, w) \in \Theta_m \right\}$$

with $\Theta_m = \{K\} \otimes \mathcal{T}_K \otimes \mathcal{V}_K \otimes \mathcal{W}_K$

- **Weight and regression models:**

$$\mathcal{W}_K = \left\{ (\pi_{w,k}(x))_{k=1}^K, w \in \mathcal{W}_K \right\}$$

$$\mathcal{R}_K = \left\{ \left(\Phi_{v_k(x), \Sigma_k(x)}(y) \right)_{k=1}^K, (v, \Sigma) \in \mathcal{T}_K \times \mathcal{V}_K \right\}$$

- **Splitting properties:**

$$H_{[\cdot], d^{\otimes n}}(\delta, S_m) \leq H_{[\cdot], \sup_x \max_k d} \left(\frac{\delta}{5}, \mathcal{R}_K \right) + H_{[\cdot], \sup_x d} \left(\frac{\delta}{5}, \mathcal{W}_K \right)$$

Bracketing entropy decomposition

- **Gaussian K -tuple bracketing entropy:**

$$\begin{aligned} H_{[\cdot], \sup_x \max_k d} \left(\frac{\delta}{5}, \mathcal{R}_K \right) &\leq H_{\max_{k=1}^K \sup_x \|\cdot_k(x)\|_2} (\epsilon_1 \delta, \Upsilon_K) + H_d(\epsilon_2 \delta, V_K) \\ &\leq (\dim(\Upsilon_K) + \dim(V_K)) \left(C_{\mathcal{R}} + \ln \frac{1}{\delta} \right) \end{aligned}$$

- **Logistic weight K -tuple bracketing entropy:**

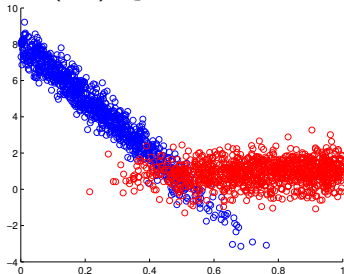
$$H_{[\cdot], \sup_x d} \left(\frac{\delta}{5}, \mathcal{W}_K \right) \leq H_{\max_k \sup_x \|\cdot_k\|} \left(\frac{\epsilon_3 \delta}{\sqrt{K}}, W_K \right) \leq \dim(W_K) \left(C_{\mathcal{W}} + \ln \frac{1}{\delta} \right)$$

- **Bracketing entropy bound if $K \leq K_{\max}$:**

$$\begin{aligned} H_{[\cdot], d^{\otimes n}}(\delta, S_m) &\leq H_{[\cdot], \sup_x \max_k d} \left(\frac{\delta}{5}, \mathcal{R}_K \right) + H_{[\cdot], \sup_x d} \left(\frac{\delta}{5}, \mathcal{W}_K \right) \\ &\leq (\dim(\Upsilon_K) + \dim(V_K) + \dim(W_K)) \left(C + \ln \frac{1}{\delta} \right) \\ &\leq \dim(S_m) \left(C + \ln \frac{1}{\delta} \right) \end{aligned}$$

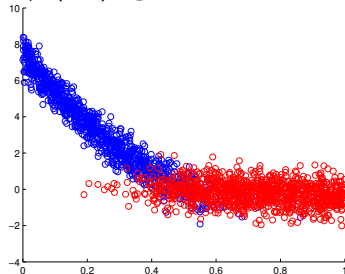
Numerical experiments

❶ $s_0 \in (S_m)_{m \in \mathcal{M}}$



Well specified

❷ $s_0 \notin (S_m)_{m \in \mathcal{M}}$



Misspecified

2 000 points

● Models S_m used:

● Affine models for the weights and the regressions:

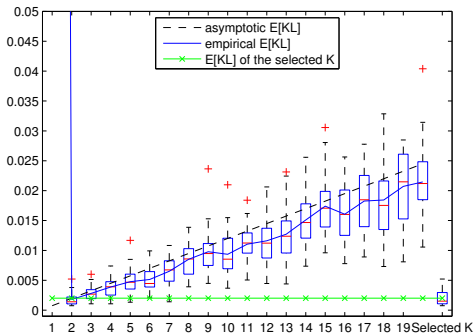
$$\mathcal{T}_K = \mathcal{W}_K = \{(a_k x + b_k)_{k=1}^K, (a, b) \in \mathbb{R}^{K \times 2}\}$$

● Free variance: $V_K = \mathbb{R}_+^K$

● Only choice is the number of components K

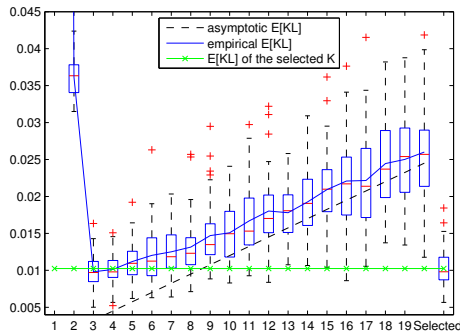
KL risk 2 000 points

❶ $s_0 \in (S_m)_{m \in \mathcal{M}}$



Well specified

❷ $s_0 \notin (S_m)_{m \in \mathcal{M}}$

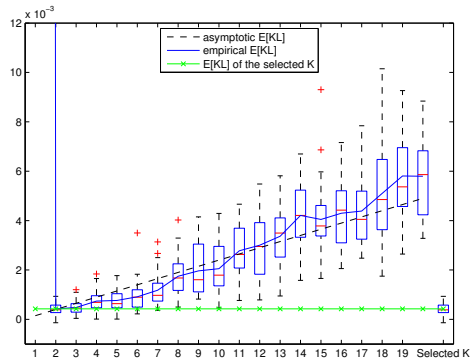


Misspecified

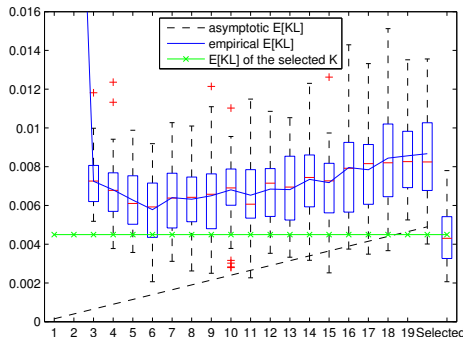
KL risk 10 000 points

❶ $s_0 \in (S_m)_{m \in \mathcal{M}}$

❷ $s_0 \notin (S_m)_{m \in \mathcal{M}}$



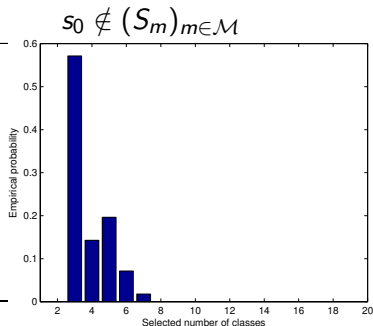
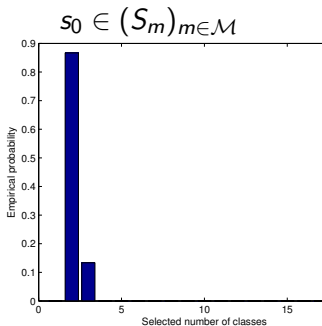
Well specified



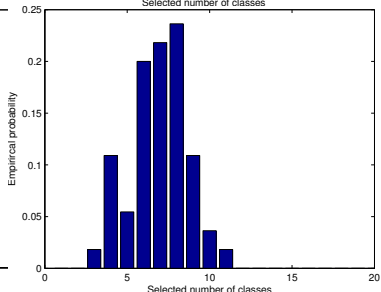
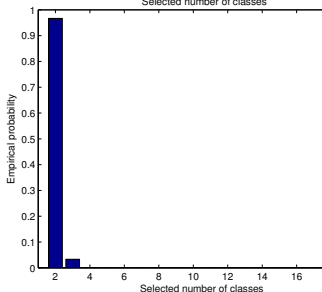
Misspecified

Histograms of the selected K

2 000



10 000



Numerical optimization

- **Model** $S_m = \{s_\theta, \theta \in \Theta_m\}$ with $\Theta_m = \{K\} \otimes \mathcal{T}_K \otimes V_K \otimes W_K$:
 - K : number of components.
 - \mathcal{T}_K and V_K : sets for the K -tuple of regressions functions and covariance matrices functions.
 - W_K : sets for the K -tuple of weights functions.
- **Maximum likelihood estimation:**

$$\hat{s}_m = \operatorname{argmin}_{\theta \in \Theta_m} - \sum_{i=1}^N \ln s_\theta(Y_i|X_i)$$

- **Penalized model selection:**

$$\hat{m} = \operatorname{argmin}_m - \sum_{i=1}^N \ln \hat{s}_m(Y_i|X_i) + \kappa \dim \Theta_m$$

- Model selection computed by **exhaustive exploration**.
- Focus on maximum likelihood estimation!

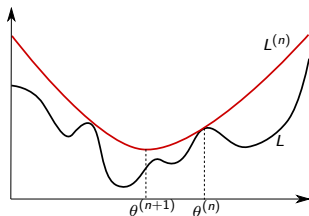
Maximum likelihood estimation

- **Model** $S_m = \{s_\theta, \theta \in \Theta_m\}$ with $\Theta_m = \{K\} \otimes \Upsilon_K \otimes V_K \otimes W_K$:
 - K : number of components.
 - Υ_K and V_K : sets for the K -tuple of regressions functions and covariance matrices functions.
 - W_K : sets for the K -tuple of weights functions.
- **Maximum likelihood estimation:**

$$\hat{s}_m = \underset{\theta \in \Theta_m}{\operatorname{argmin}} - \underbrace{\sum_{i=1}^N \ln \left(\sum_{k=1}^K \pi_{w,k}(X_i) \Phi_{v_k(x_i), \Sigma_k(x_i)}(Y_i) \right)}_{L(\theta)}$$

- **Non convex** minimization problem!
- **Majorization/Minimization** approach

MM approach



- **Iterative** approach to minimize $L(\theta)$ by minimizing a **sequence of (convex) proxies** of L .
- **Majorization/Minimization:**
 - **Current estimate** of the minimizer: $\theta^{(n)}$
 - Construction of a **Majorization** $L^{(n)}$ of L such that $L^{(n)}(\theta^{(n)}) = L(\theta^{(n)})$ with $L^{(n)}$ easy to minimize (convex for example).
 - Computation of a **Minimizer**

$$\theta^{(n+1)} = \operatorname{argmin} L^{(n)}(\theta)$$

- By construction, $L(\theta^{(n+1)}) \leq L(\theta^{(n)})$!
- Very **generic methodology**...
- Minimization can be replaced by a diminution...

Maximum Likelihood and EM

- Back to **our** maximum likelihood:

$$L(\theta) = L(\textcolor{red}{K}, \textcolor{blue}{v}, \textcolor{blue}{\Sigma}, \textcolor{red}{w}) = - \sum_{i=1}^N \ln \left(\sum_{k=1}^{\textcolor{red}{K}} \pi_{\textcolor{red}{w},k}(X_i) \Phi_{\textcolor{blue}{v}_k(X_i), \textcolor{blue}{\Sigma}_k(X_i)}(Y_i) \right)$$

- **EM**: specific case of MM for this type of mixture.
- **(Conditional) Expectation**: at step n , we let

$$P_k^{i,(n)} = P \left(k_i = k \mid X_i, Y_i, \textcolor{red}{K}, \textcolor{blue}{v}^{(n)}, \textcolor{blue}{\Sigma}^{(n)}, \textcolor{red}{w}^{(n)} \right)$$

and
$$L^{(n)}(\textcolor{red}{K}, \textcolor{blue}{v}, \textcolor{blue}{\Sigma}, \textcolor{red}{w}) = - \sum_{i=1}^N \sum_{k=1}^{\textcolor{red}{K}} P_k^{i,(n)} \ln \left(\pi_{\textcolor{red}{w},k}(X_i) \Phi_{\textcolor{blue}{v}_k(X_i), \textcolor{blue}{\Sigma}_k(X_i)}(Y_i) \right).$$

- **Maj. prop.:** $L \leq L^{(n)} + \text{Cst}^{(n)}$ with equ. at $\theta = (\textcolor{red}{K}, \textcolor{blue}{v}^{(n)}, \textcolor{blue}{\Sigma}^{(n)}, \textcolor{red}{w}^{(n)})$
- **Separability** in $(\textcolor{blue}{v}^{(n)}, \textcolor{blue}{\Sigma}^{(n)})$ and $\textcolor{red}{w}^{(n)}$:

$$L^{(n)}(\textcolor{red}{K}, \textcolor{blue}{v}, \textcolor{blue}{\Sigma}, \textcolor{red}{w}) = \left(- \sum_{i=1}^N \sum_{k=1}^{\textcolor{red}{K}} P_k^{i,(n)} \ln \Phi_{\textcolor{blue}{v}_k(X_i), \textcolor{blue}{\Sigma}_k(X_i)}(Y_i) \right) + \left(- \sum_{i=1}^N \sum_{k=1}^{\textcolor{red}{K}} P_k^{i,(n)} \ln \pi_{\textcolor{red}{w},k}(X_i) \right)$$

Minimization of $L^{(n)}$

- **Separability** in $(v^{(n)}, \Sigma^{(n)})$ and $w^{(n)}$:

$$L^{(n)}(\mathbf{K}, \mathbf{v}, \mathbf{\Sigma}, \mathbf{w}) = \left(- \sum_{i=1}^N \sum_{k=1}^{\mathbf{K}} P_k^{i,(n)} \ln \Phi_{\mathbf{v}_k(X_i), \mathbf{\Sigma}_k(X_i)}(Y_i) \right) + \left(- \sum_{i=1}^N \sum_{k=1}^{\mathbf{K}} P_k^{i,(n)} \ln \pi_{\mathbf{w},k}(X_i) \right)$$

- For the **regression parameters** $(v^{(n)}, \Sigma^{(n)})$:

- \mathbf{K} weighted linear regressions: $-\sum_{i=1}^N P_k^{i,(n)} \ln \Phi_{\mathbf{v}_k(X_i), \mathbf{\Sigma}_k(X_i)}(Y_i)$

- Explicit formulas!

- For the **weight parameters** $w^{(n)}$:

- Single \mathbf{K} modality logistic regression: $-\sum_{i=1}^N \sum_{k=1}^{\mathbf{K}} P_k^{i,(n)} \ln \pi_{\mathbf{w},k}(X_i)$

- Iterative minimization scheme (Newton = Iterative Reweighted Least Square)

Initialization

- **Very important issue!**
- For the **weights**: initialization to uniform weights seems sufficient.
- For the **means**:
 - Comparison between **several** strategies
 - Naive purely random initialization
 - Small-EM: Random initialization followed by a few minimization steps and selection
 - Advanced Small-EM: Initialization based on a first $2D$ clustering followed by a few minimization steps and selection
 - Advanced Small-EM 2 : Initialization based on a random drawing of lines between points clustering followed by a few minimization steps and selection
 - **Criterion**: lowest likelihood for a given amount of time!
 - Similar results in term on expectation but **different behaviors in term of dispersion**:
 - Too simple strategies fail sometimes to provide a satisfactory answer while too complex ones may not explore sufficient local maxima.
 - **Winner**: Advanced Small-EM 2 with 3 minimizations steps and 50 candidates.

Newton-EM Algorithm

- **Initialization with Advanced Small-EM 2:**

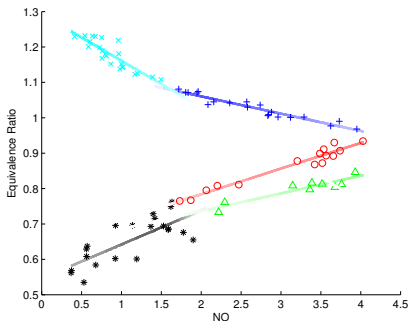
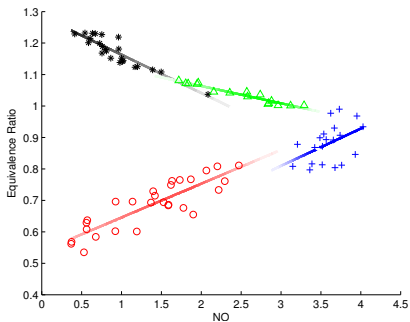
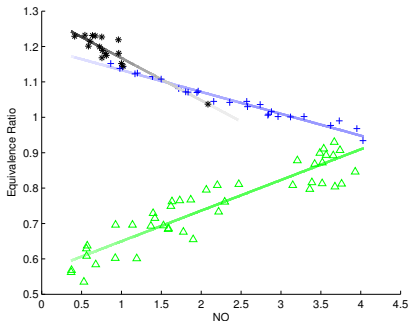
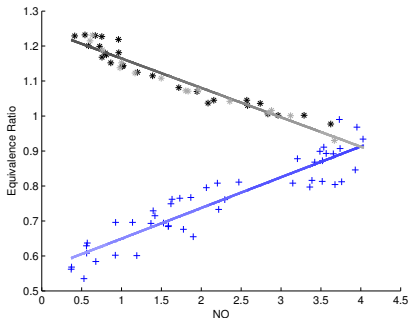
- Initialization based on a random drawing of lines between points clustering
- 3 minimizations steps
- selection among 50 candidates

- **Iterate** until convergence:

- **Newton steps** over weights $w^{(n)}$ if the likelihood increases (up to 5 times)
- K **linear regressions** to update mean and variance ($v^{(n)}, \Sigma^{(n)}$) in each class

- **Note:** Initialization issues in high dimension with this scheme!

Numerical results



Penalization strategy

- **Penalized model selection:**

$$\hat{m} = \operatorname{argmin}_{m \in \mathcal{M}} \sum_{k=1}^K -\ln \hat{s}_m(Y_i|X_i) + \operatorname{pen}(m)$$

- **Theoretical analysis:**

$$\operatorname{pen}(m) = \kappa(C + \ln n) \dim \Theta_m$$

- κ and C are only **loosely upper bounded!**
- **In practice**, use $\operatorname{pen}(m) = \kappa \dim \Theta_m$ with κ chosen appropriately.
- **Classical choice:**
 - BIC: $\kappa = \log n/2$
 - AIC: $\kappa = 1$
- Here: **Jump/slope heuristic** = data driven choice of κ

Ideal penalty

- By definition:

$$KL^{\otimes n}(s_0, \hat{s}_m) = P_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right) - \underbrace{\nu_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right)}_{\text{pen}_{\text{id}}(m)/n}$$

- With the **ideal penalty** $\text{pen}_{\text{id}}(m)$:

$$\begin{aligned} KL^{\otimes n}(s_0, \hat{s}_{\hat{m}}) &= P_n^{\otimes n} \left(-\ln \frac{\hat{s}_{\hat{m}}}{s_0} \right) + \frac{\text{pen}_{\text{id}}(\hat{m})}{n} \\ &\leq \inf_m P_n^{\otimes n} \left(-\ln \frac{\hat{s}_m}{s_0} \right) + \frac{\text{pen}_{\text{id}}(m)}{n} \leq \inf_m KL^{\otimes n}(s_0, \hat{s}_m) \\ &\leq \inf_m (KL^{\otimes n}(s_0, \tilde{s}_m) + (KL^{\otimes n}(s_0, \hat{s}_m) - KL^{\otimes n}(s_0, \tilde{s}_m))) \end{aligned}$$

- Ideal penalty oracle inequality:**

$$\mathbb{E} [KL^{\otimes n}(s_0, \hat{s}_{\hat{m}})] \leq \inf_{S_m \in \mathcal{S}} \left(\underbrace{KL^{\otimes n}(s_0, \tilde{s}_m)}_{\text{Bias term}} + \underbrace{\mathbb{E} [KL^{\otimes n}(s_0, \hat{s}_m) - KL^{\otimes n}(s_0, \tilde{s}_m)]}_{\text{Variance term}} \right)$$

Jump/Slope heuristic

- **Ideal penalty decomposition:**

$$\begin{aligned}\frac{\text{pen}_{\text{id}}(m)}{n} &= -\nu_n^{\otimes n} \left(-\ln \frac{\widehat{s}_m}{s_0} \right) \\ &= \nu_n^{\otimes n} \left(-\ln \frac{\widetilde{s}_m}{\widehat{s}_m} \right) - \underbrace{\nu_n^{\otimes n} (-\ln \widetilde{s}_m)}_{\text{independent of } m} + \underbrace{\nu_n^{\otimes n} (-\log s_0)}_{\text{independent of } m}.\end{aligned}$$

- **Jump/Slope heuristic:**

- Concentration: $\nu_n^{\otimes n} (-\ln \widetilde{s}_m) \ll \nu_n^{\otimes n} \left(-\ln \frac{\widetilde{s}_m}{\widehat{s}_m} \right)$
- Symmetry: $P_n^{\otimes n} (-\ln(\widetilde{s}_m/\widehat{s}_m)) \sim P_n^{\otimes n} (-\ln(\widehat{s}_m/\widetilde{s}_m))$

- Resulting **approximation:**

$$\frac{\text{pen}_{\text{id}}(m)}{n} \sim 2P_n^{\otimes n} \left(-\ln \frac{\widetilde{s}_m}{\widehat{s}_m} \right) - \underbrace{P_n^{\otimes n} (-\log s_0)}_{\text{independent of } m}.$$

- $P_n^{\otimes n} \left(-\ln \frac{\widetilde{s}_m}{\widehat{s}_m} \right)$ has still **to be estimated!**

Minimal penalty

- If $\text{pen}(m) = \kappa P_n^{\otimes n} \left(-\ln \frac{\tilde{s}_m}{\hat{s}_m} \right)$ then

$$P_n^{\otimes n}(-\ln \hat{s}_m) + \text{pen}(m) = (1 - \kappa) P_n^{\otimes n}(-\ln \hat{s}_m) + \kappa P_n^{\otimes n}(-\ln \tilde{s}_m)$$

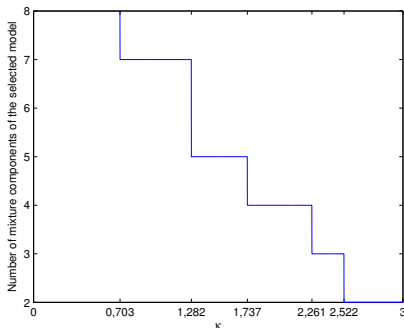
- **No tradeoff is $\kappa < 1$!**
- **Minimal penalty:** $\text{pen}_{\min}(m) = P_n^{\otimes n}(-\ln(\tilde{s}_m/\hat{s}_m))$
- Jump/Slope heuristic strongest assumption: **parametric approximation** of pen_{\min}

$$\text{pen}_{\min}(m) = \text{pen}(\kappa, m)$$

where pen *shape* is given by the theoretical study!

- **Simplest case:** $\text{pen}(\kappa, m) = \kappa \dim S_m$.

Jump heuristic

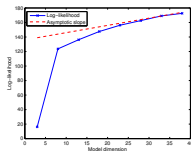


- Minimal penalty for which there is a **tradeoff**:

$$\text{pen}_{\min}(m) = \text{pen}(\kappa, m)$$

- Compute the models selected for several κ and detect a **jump in the model dimensions**.
- Not always a clear single jump...

Slope heuristic



- **Observation:**

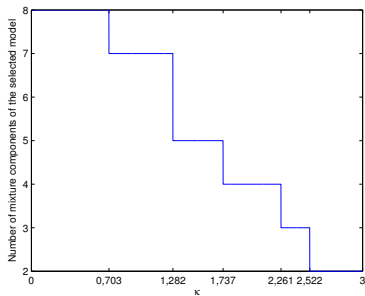
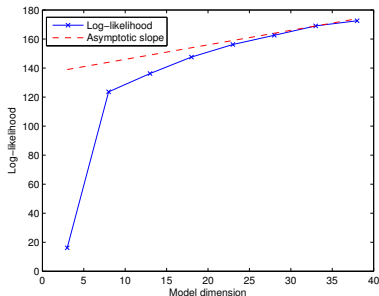
$$\text{pen}_{\text{id}}(m) = P_n^{\otimes n} \left(-\ln \frac{\tilde{S}_m}{\hat{S}_m} \right) = P_n^{\otimes n} (\ln \hat{S}_m) + P_n^{\otimes n} (-\ln \tilde{S}_m)$$

- If the model are more and more complex, one may expect that the **projection bias converges to a constant**: $P_n^{\otimes n} (-\ln \tilde{S}_m) \sim C$
- This implies $\text{pen}_{\text{id}}(m) \sim P_n^{\otimes n} (\ln \hat{S}_m) + C$
- If $\text{pen}_{\text{id}}(\text{ind} m) = \text{pen}(\kappa, m)$ then κ can be estimated by a **regression** as

$$\text{pen}(\kappa, m) - C \sim \underbrace{P_n^{\otimes n} (\ln \hat{S}_m)}_{\text{data driven}}$$

- If $\text{pen}(\kappa, m) = \kappa \dim S_m$, κ measures the **slope** of $P_n^{\otimes n} (\ln \hat{S}_m)$ with respect to $\dim S_m$.

Slope heuristic



● **Slope heuristic** with $\text{pen}(\kappa, m) = \kappa \dim S_m$: $\kappa \sim 1$

● **Resulting penalties:**

● Slope heuristic: $\text{pen}(m) = 2 \dim(S_m)$

● BIC: $\text{pen}(m) = 2.3 \dim(S_m)$

● AIC: $\text{pen}(m) = 2 \dim(S_m)$

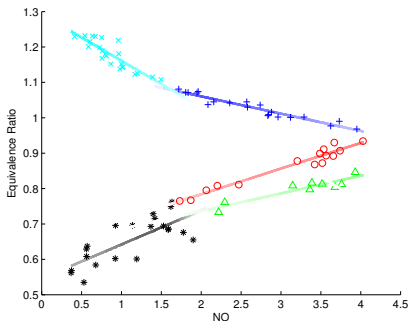
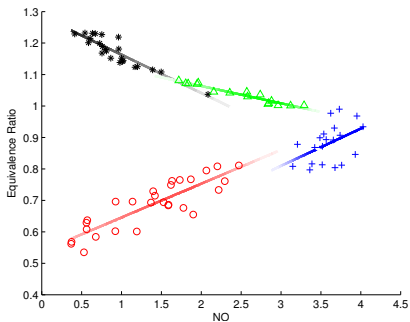
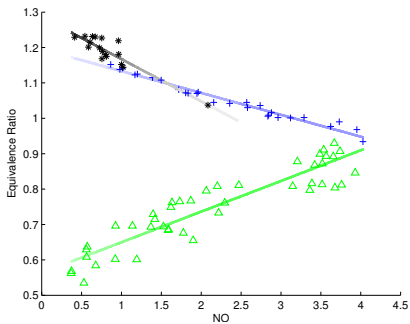
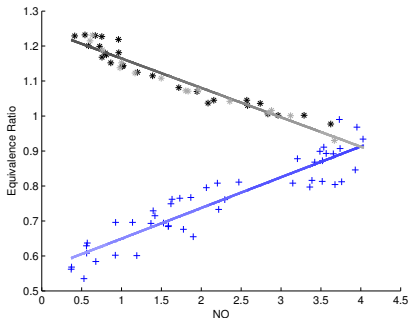
● **Selected number of clusters:**

● Slope heuristic: 4

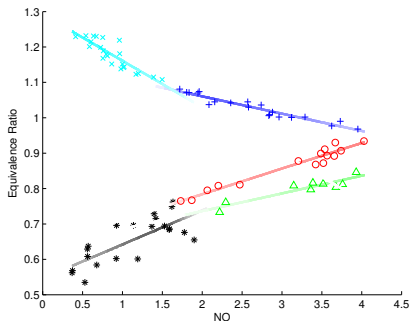
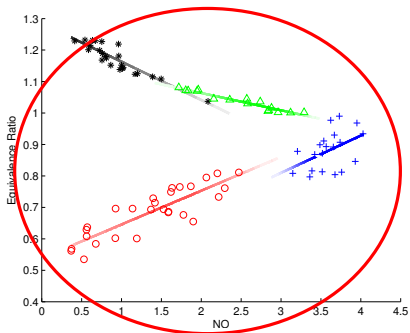
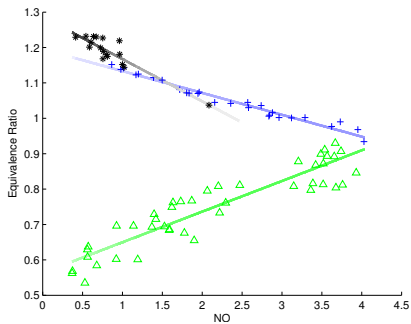
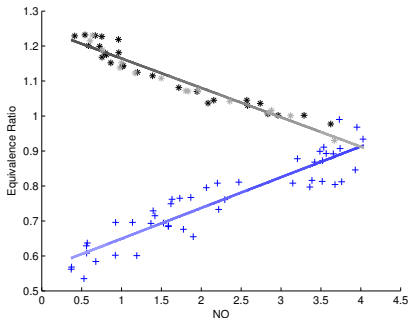
● BIC: 4

● AIC: 7!

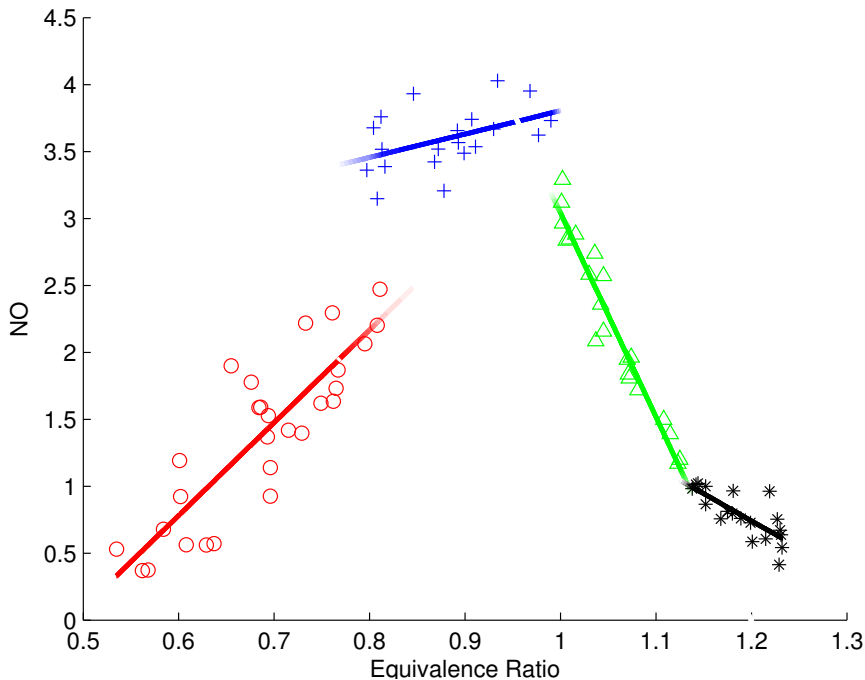
Numerical results



Numerical results



Validation?



Conclusion

● Framework:

- Mixture of regressions.
- Proposed tool: Mixture of Gaussian regressions with logistic mixing weights.
- Penalized maximum likelihood conditional density estimation.

● Contributions:

- Theoretical guarantee for the conditional density estimation problem.
- Efficient minimization algorithm.
- Numerical penalty calibration.

● Perspectives:

- Proof for penalty calibration by slope heuristic.
- Enhanced Spatialized Gaussian Mixture Model with piecewise logistic weights (S. Cohen).