

Model-based Statistical Learning



Pr. Charles BOUVEYRON

Professor of Statistics
Chair of the Institut 3IA Côte d'Azur
Université Côte d'Azur & Inria

charles.bouveyron@univ-cotedazur.fr
[@cbouveyron](https://twitter.com/cbouveyron)

Parsimonious models for GMM

In many situations, it may be useful to consider more constrained models:

- because we have a limited number of obs. and fitting a (full) GMM requires to estimate a lot of parameters \rightarrow parsimonious models are better
- we may have extra informations about the data (e.g. the variables are independent $\rightarrow \Sigma_h = \Sigma_k \Sigma_p$) and we can encode this information as a constraint on the model
- ...

Parsimonious models for GMM

$$\kappa = 4, p = 50$$

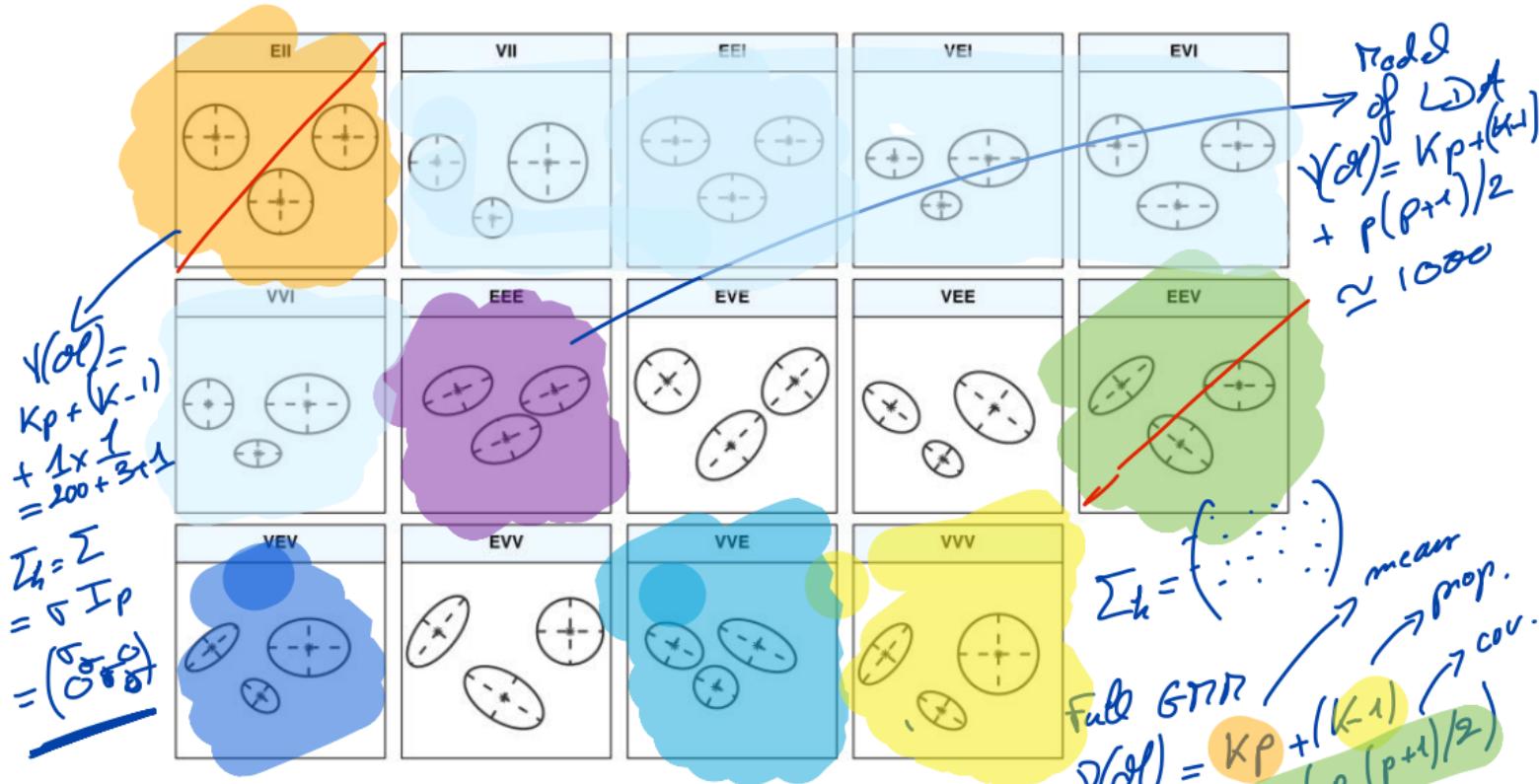


Figure: The parsimonious models of Mclust.

The interest here is to exploit the "catalog" of mixture models to find the most appropriate one and rely on model selection to make the choice



⚠ the list of models that you can test is very large due to the combination of pdf family + constraints + nb of components

Mixture Models



⇒ Model selection is here to pick the most appropriate model for the data.

Parsimonious models for GMM

Such models are available in most softwares:

- **Mclust** package for R
 - ↳ a selection of constrained models
(using the VV \rightarrow EEE nomenclature)
 - ↳ BIC is used for model selection.
- (R)Flexmix (Matlab + C++ + Python)
 - flexmix, ...

How to choose between models?

↳ we rely on model selection criteria allow to pick a good model in a list of candidates.

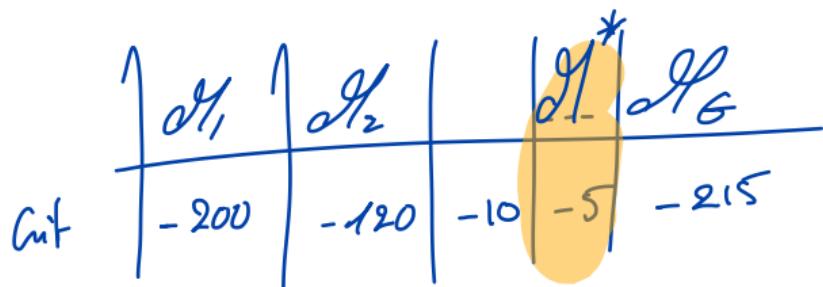
→ AIC

→ Slope heuristic

→ BIC

→ ...

→ ICL



Model selection

The roots of Mod. selection can be found in Bayesian statistical theory.

Let us first consider a list of candidate models $\{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_G\}$ and associated prior probabilities $p(\mathcal{M}_g)$. In practice, we can assume the all prior probabilities are equal. $p(\mathcal{M}_g) = p \forall g$.

The idea of model selection is to evaluate a specific quantity : $p(\mathcal{M}_g | X)$

Model selection

The roots of model selection can be found in Bayesian statistic theory:

Thanks to the Bayes' theorem, we can write :

$$p(\theta_g | X) \propto p(X | \theta_g) p(\theta_g). \quad (1)$$

When the models have unknown parameters, the law of total probability allows to calculate it by integrating out the model parameters:

$$p(X | \theta_g) = \int p(X | \theta_g, \Theta_g) p(\Theta_g | \theta_g) d\Theta_g$$

Ring: this integration is usually very difficult or impossible to do!

Model selection

Note that the quantity $p(X|\mathcal{M}_g)$ is called the integrated likelihood or the marginal likelihood or the evidence.

In the Bayesian framework, the knowledge of the integrated likelihood allows to do model selection by picking the model with the highest posterior proba:

$$\begin{aligned}\mathcal{M}^* &= \operatorname{arg\max}_{\mathcal{M}_g} p(\mathcal{M}_g | X) \\ &= \operatorname{arg\max}_{\mathcal{M}_g} p(X | \mathcal{M}_g).\end{aligned}$$

Note that, when comparing two specific models \mathcal{M}_1 and \mathcal{M}_2 , the ratio $B_{12} = \frac{p(X|\mathcal{M}_1)}{p(X|\mathcal{M}_2)}$ is called the Bayes' factor.

In particular, if $B_{12} > 1$, the model \mathcal{M}_1 should be preferred. And, if $B_{12} > \underline{100}$, the model \mathcal{M}_1 as a strong evidence against \mathcal{M}_2 (Jeffreys, 1961)

In practice, when considering the frequentist case, this framework is also useful because the integrated likelihood can be approximated.

$$\hat{\theta}^* = \underset{\theta_g}{\operatorname{arg\max}} \quad \tilde{p}(x | \theta_g)$$

In particular, the BIC (Bayesian Information Criterion), proposed by Schwarz in 1978, is an approximation of the integrated likelihood:

$$\log p(x | \theta_g) \approx \underbrace{\log p(x | \hat{\theta}_g, \theta_g)}_{\text{BTC}} - \frac{\gamma(\theta_g) \lg(m)}{2}$$

where $\gamma(\theta_g)$ is the number of free parameters in the model θ_g , and m is the number of observations.

BIC is an asymptotic approximation of the integrated likelihood based on a second order Taylor expansion of the logarithm of the integrand, around its maximum $\hat{\theta}_g$.

Ring: it is worth noticing that, unfortunately, the assumptions made about the regularity of the models in BIC approximations are not satisfied for mixture models!

Fortunately, BIC is behaving very well for mixture models in practice!

Model selection: AIC, BIC, ICL

Classical model selection tools can be used easily:

- AIC and BIC are two criteria that can be used in large numbers of situations (regression, mixture models, ...)
- ICL (integrated classification likelihood) was proposed for the specific case of clustering.

The idea of ICL is the same as BIC but instead of approximating $p(X|\theta_g)$ it focuses on the integrated classification likelihood $p(X, Z|\theta_g)$

Model selection: AIC, BIC, ICL

where Z is the latent variable that encodes the group memberships.

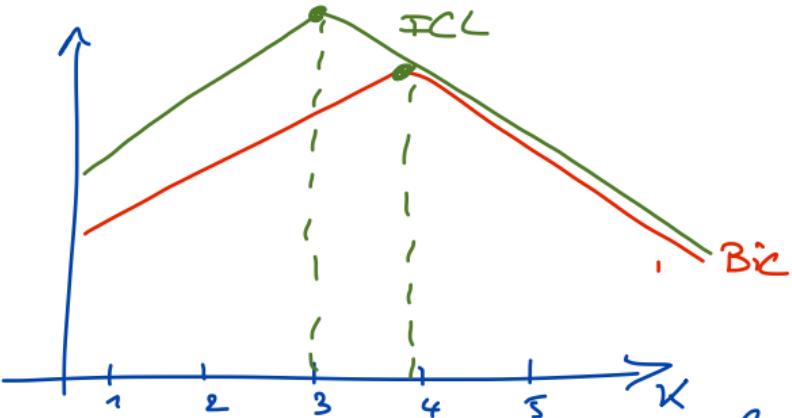
$$p(x|\mathcal{H}_g) \rightarrow p(x, z|\mathcal{H}_g)$$

\uparrow
 Bic
 \uparrow
 ICL.

ICL then uses the same approximations of BIC and we end up with the following approximation:

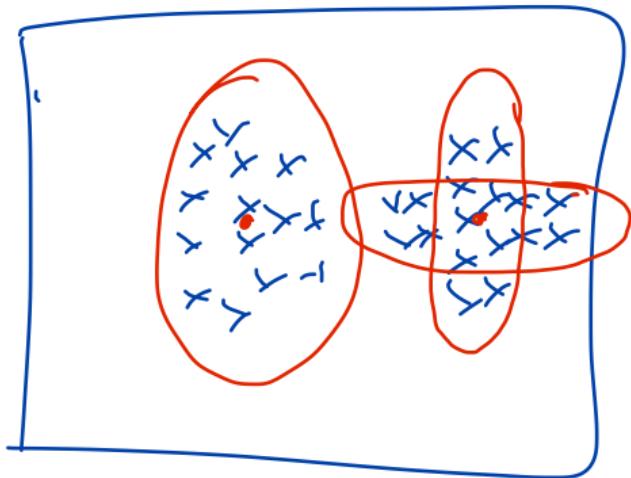
$$\log p(x, z | \theta_g) \approx \left. \begin{aligned} & \log p(x | \hat{\theta}_{g_1}, \hat{\theta}_g) - \frac{v(\theta_g)}{n} \log(n) \\ & - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K \hat{z}_{ik} \log^2(\hat{z}_{ik}) \end{aligned} \right\} \text{ICL.}$$

Rung: it clearly appears that $ICL = BIC - \frac{1}{2} \sum_i^n z_{ik} \log(z_{ik})$ and that ICL penalizes more the likelihood compared to Bic.

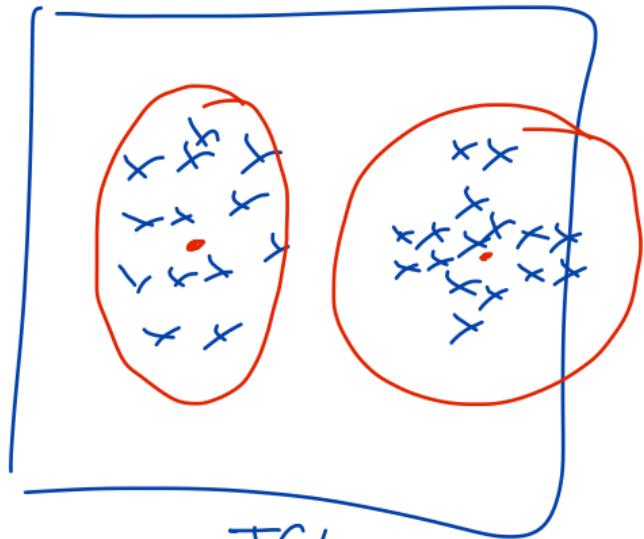


Rung: the objectives of Bic and ICL are slightly different:

- BIC aims to find a good fit of the data
- ICL aims to find a good model to cluster the data.



Bic
 \Rightarrow GMM with 3 comp.



ICL
 \Rightarrow GMM with 2 comp-