# 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

UNIVERSITÉ
CÔTE D'AZUR

*informatiques* *mathématiques*
Inria

# Parsimonious models for GMM

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

$$\text{Full GMM}: \quad \Theta = \{\pi_k, \mu_k, \Sigma_k\}$$

$$\text{Com GMM}: \quad \Theta = \{\pi_k, \mu_k, \Sigma\}$$

$$\text{diag GMM}: \quad \Theta = \{\pi_k, \mu_k, \sigma_k^2 I_p\}$$

$$\text{iso GMM}: \quad \Theta = \{\pi_k, \mu_k, \sigma^2 I_p\}$$

$$\vdots$$

$$k.\text{ means}: \quad \Theta = \{\frac{1}{K}, \mu_k, \sigma^2 I_p\}$$

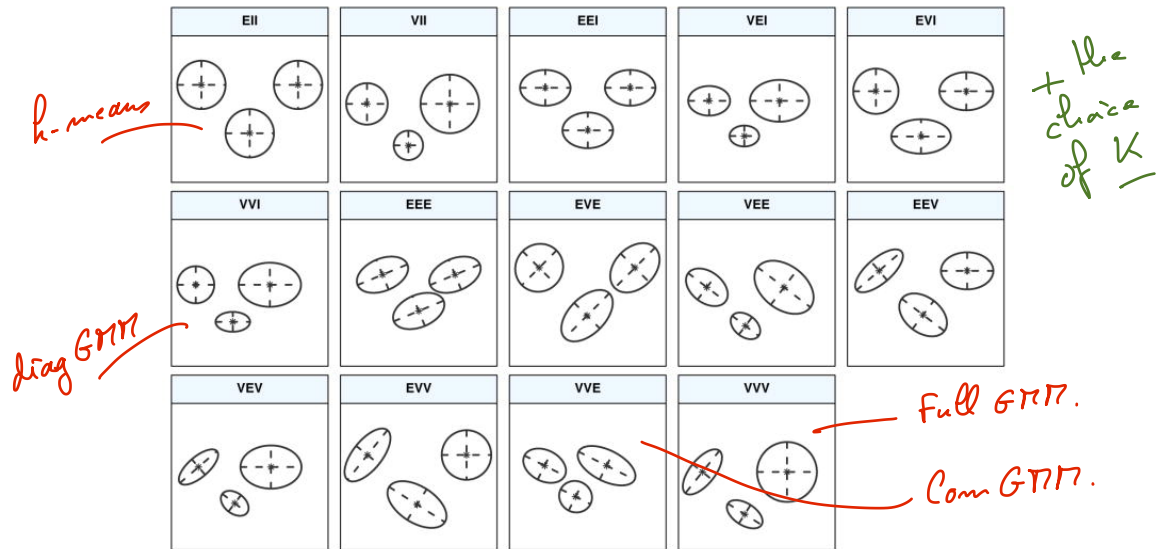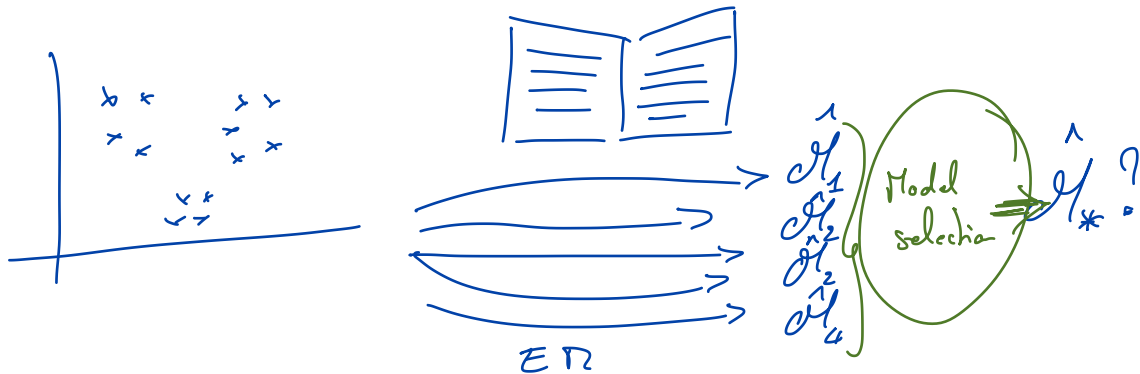# Parsimonious models for GMM



Figure: The parsimonious models of Mclust.

## How to choose between models?

In practice, when facing some data set, we need to pick the appropriate number $K$ of components and the best GMM model to fit our data.



$$\hat{\mathcal{M}}_1$$
$$\hat{\mathcal{M}}_2$$
$$\hat{\mathcal{M}}_3$$
$$\hat{\mathcal{M}}_4$$

Model selection $\Rightarrow \hat{\mathcal{M}}_* \; ?$

$EM$

## Model selection

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

Let's first consider a set of models to test:

$\{ \mathcal{M}_1, \mathcal{M}_2, \ldots, \mathcal{M}_g \}$ and associated to prior probabilities $p(\mathcal{M}_g)$ (in practice we prefer to use
$$p(\mathcal{M}_g) = p \cdot \forall g$$
)

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

Thanks to the Bayes theorem:

$$p(\mathcal{H}_g \mid X) \propto p(X \mid \mathcal{H}_g)\, p(\mathcal{H}_g) \qquad (\iota).$$

To compute this quantity, we have to evaluate it for all possible combinations of the parameters of $\mathcal{H}_g$

$$p(X \mid \mathcal{H}_g) = \int p(X \mid \mathcal{H}_g, \Theta_g) \cdot p(\Theta_g \mid \mathcal{H}_g)\, d\Theta_g.$$

Rmq: due to the high-dimensional integration over $\Theta_g$, this computation is very often intractable!

<u>Note</u>: the quantity $p(X \mid \mathcal{M})$ is called the ==integrated likelihood== or the ==evidence== or the ==marginal likelihood.==

In practice, we prefer to avoid the computational problem by approximating the integrated likelihood. Then:

$$\mathcal{M}^* = \arg\max_{\mathcal{M}} \tilde{p}(X \mid \mathcal{M}).$$

Among the possible approximation of the integrated likelihood, the most popular one is BIC : Bayesian Information Criterion, Schwarz, 1978

$$\log p(X|\mathcal{M}_g) \simeq \underbrace{\log p(X|\mathcal{M}_g, \hat{\theta}_{gML}) - \frac{\nu(\mathcal{M}_g)}{2} \log(n)}_{\text{BIC criterion.}}$$

where $\nu(\mathcal{M}_g)$ is the number of free parameters in the model $\mathcal{M}_g$, and $n$ is the number of observations.

For information, Bic is both an asymptotic approximation of the marginal likelihood $\left( Bic \underset{n \to \infty}{\Longrightarrow} p(X | \mathcal{M}_g) \right)$ and a second order Taylor expansion of the logarithm of the integrand, around it maximum $\hat{\theta}_{g \, ML}$.

In practice :

$$\mathcal{M}^* = \underset{\mathcal{M}_g}{\arg\max} \; Bic(\mathcal{M}_g).$$

<u>Rmk</u>: it is worth mentioning that, unfortunately, the assumptions made about the regularity of the models in BIC are not satisfied for mixture models!
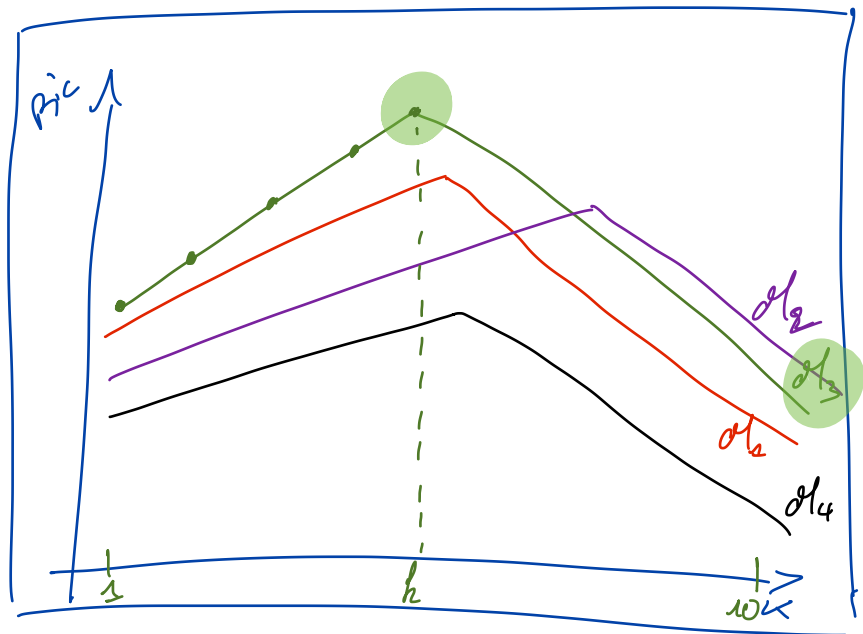
$\Rightarrow$ Fortunately, BIC is behaving very well for mixture models and even when $n < +\infty$.

what we do:

$\mathcal{M}_1$ = Full GMM with $K=2$ $\xrightarrow{EM}$ $\hat{\theta}_{1mc}$ $\xrightarrow{Bic}$ $Bic(\mathcal{M}_1)$

$\mathcal{M}_2$ = diag GMM — $K=2$ $\xrightarrow{EM}$ $\hat{\theta}_2$ $\longrightarrow Bic(\mathcal{M}_2)$

$\mathcal{M}_3$ = Full GMM — $K=3$ $\xrightarrow{EM}$ $\hat{\theta}_3$ $\longrightarrow Bic(\mathcal{M}_3)$

$\mathcal{M}_4$ = diag GMM — $K=3$ $\xrightarrow{EM}$ $\hat{\theta}_4$ $\longrightarrow Bic(\mathcal{M}_4)$

and we choose $\mathcal{M}_3$ because
$Bic(\mathcal{M}_3)$ is the largest one.

with Melot:



pic $\uparrow$

$1$  $h$  $10K$

$M_2$
$M_3$
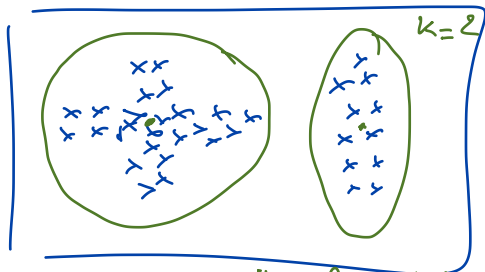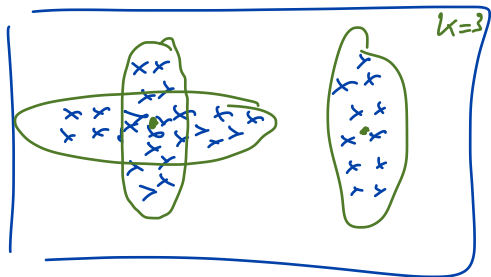$M_4$

Rmk: Bic (and Aic which is a close criterion)
can be used for model selection in all
situations : regression, clustering, ...

# Model selection for clustering:

The notion of mixture components and clusters may differ in some specific situations and it would be great to have a MS criterion that takes this into account.



K=3

K=2

is K=2 better for clustering?

A sound way to answer this it to come back to the model selection theory and consider the integrated complete likelihood instead of the integrated likelihood:

$$\mathcal{M}^* = \arg\max_{\mathcal{M}_g} \; p\left(X, Y \mid \mathcal{M}_g\right)$$

Then, it is possible to make the same approximations.

$$p\left(X, Y \mid \mathcal{M}_g\right) \simeq \log p\left(X \mid \hat{\theta}_g, \mathcal{M}_g\right) - \frac{\nu(\mathcal{M}_g)}{2} \log(n)$$
$$- \frac{1}{2} \sum_{i=1}^{n} \sum_{k=1}^{K} t_{ik} \log\left(t_{ik}\right)^2 = ICL$$