

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

The EM algorithm for GMM

Let's now consider the EM algorithm for GMM

E step: $E[p(x_i) | \theta^*] = \sum_i \sum_k E[z_{ik} | x, \theta] \log (\pi_k \phi(x_i; \theta_k))$



through Bayes' theorem
 $t_{ik} = \frac{\pi_k^* \phi(x_i; \theta_k^*)}{p(x)}$

M step: $\max_{\theta} \prod_{i=1}^n \prod_{k=1}^K t_{ik}^* \pi_k \phi(x_i; \mu_k, \Sigma_k)$
 $= (\pi_h, \mu_h, \Sigma_h)_h$

The EM algorithm for GMM

E step: $\underline{t_{ik}^*} = P[z|x; \theta^*] \stackrel{\text{Bayes}}{=} \frac{\pi_k^* \phi(x_i; \mu_k^*, \Sigma_k^*)}{\sum_h^n \pi_h^* \phi(x_i; \mu_h^*, \Sigma_h^*)}$

π step: $\pi_k^{*(new)} = \frac{\sum_{i=1}^n t_{ik}^*}{n}$

$$\mu_k^{*(new)} = \frac{\sum_{i=1}^n t_{ik}^* x_i}{n}$$

$$\Sigma_k^{*(new)} = \frac{1}{n} \sum_{i=1}^n t_{ik}^* (x_i - \mu_k^*)^T (x_i - \mu_k^*)$$

Deeper in the GMM and EM algorithm:

Let's recall that the EM algorithm introduces a latent variable (z) to encode the group memberships, and then maximizes iteratively the Expectation of the complete data likelihood (of the pair (x, z)). - Start with a given value of θ^*

- E step : compute $E[\ell(x, z | \theta) | \theta^*]$
- M step : maximize $E[\text{---}]$ according to θ to get a new value for θ^* .

Let's first focus on the complete-data log-likelihood:

$$\underline{L}(x, z | \theta) = \sum_{i=1}^n \sum_{h=1}^K z_{ih} \log(\pi_h \phi(x_i | \theta_h))$$

$$\log(L(x, z | \theta)) = \overbrace{\sum_i \log \underbrace{\sum_h}_{\text{one ind. obs.}} p(z_i | x_i, \theta)}_{\text{.}} - \phi(x_i; \theta_h)$$

Note: we can switch the \log and the \sum because for one ind. obs., it can belong to only one group.
 $\log \sum = \sum \log(1)$

$$= \sum_i \sum_h z_{ih} \log \left(p(z_i | x_i, \theta) \times \frac{\phi(x_i; \theta_h)}{\pi_h} \right)$$

$$= \sum_i \sum_h z_{ih} \log \left(\frac{\pi_h}{\pi_h} \phi(x_i; \theta_h) \right)$$

Note: here Z is assumed to be a vector of binary values encoding the group membership:

$$z_i = (0, 0, 1, 0) \Rightarrow x_i \text{ belongs to cluster 3}$$

$$E[\ell(x, z | \theta) | \theta^*] = \sum_i \sum_h [E[z_{ih} | x_i, \theta^*] \log(\pi_h \phi(x_i; \theta_h))]$$

$$t_{ih} = E[z_{ih} | x_i, \theta^*] = P(z_{ih}=1 | x_i, \theta^*)$$

$$\text{Bayes} \quad = \frac{P(z_{ih}=1 | \theta^*) \times P(x_i | z_{ih}=1, \theta^*)}{P(x_i | \theta^*)}$$

This quantity can be easily computed with a table of the distribution ϕ

$$= \frac{\pi_h^* \phi(x_i; \theta_h^*)}{\text{---}}$$

In the Π step, the expected complete data log likelihood is maximised over θ to get the new θ^* :

$$\underset{\theta}{\operatorname{Max}} E[\ell(x, z | \theta) | \theta^*] \Leftrightarrow \underset{\theta}{\operatorname{Max}} \sum_i \sum_h t_{ih} \log(\pi_h \phi(x_i; \theta_h))$$

$$\text{where } \theta = (\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_h)$$

For the Gaussian mixture model:

$$\phi(x_i; \theta_h) = \phi(x_i; \mu_h, \Sigma_h) = \frac{1}{(2\pi)^{p/2}} |\Sigma_h|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (x_i - \mu_h)^T \Sigma_h^{-1} (x_i - \mu_h)\right)$$

$$\text{where } \theta = (\pi_1, \dots, \pi_h, \mu_1, \dots, \mu_h, \Sigma_1, \dots, \Sigma_h)$$

Let's now do the calculations to get update values for θ^* , meaning $\bar{\pi}_h^*$, μ_h^* and Σ_h^* :

Exo 1: Write down the $Q(\theta) = E[\ell(x, z | \theta) | \theta^*]$ for the GMN

$$\begin{aligned} Q(\theta) &= \sum_i \sum_h \text{E}_{\text{th}} \left[\log(\bar{\pi}_h) + \log(\phi(z_i | \theta_h)) \right] \\ &= \sum_i \sum_{+h} \left[\log \bar{\pi}_h - \frac{1}{2} \log(2\pi) - \frac{1}{2} \log \left(|\Sigma_h| \right) \right. \\ &\quad \left. - \frac{1}{2} (z_i - \mu_h)^T \Sigma_h^{-1} (z_i - \mu_h) \right] \end{aligned}$$

Exo 2: maximisation over μ_h

$$\frac{\partial Q(\theta)}{\partial \mu_h} = 0 \Rightarrow \underline{\mu_h^*}$$

$$\frac{\partial}{\partial \mu_h} Q(\theta) = \sum_k \left[\sum_i t_{ih} \left[\log \pi_h - \frac{p}{2} \log(2\pi) - \frac{1}{2} \rho_{\theta} |\Sigma_h|^{-1} \right. \right.$$

$$\left. \left. - \frac{1}{2} \sum_i t_{ih} (x_i - \mu_h)^T \Sigma_h^{-1} (x_i - \mu_h) \right] \right]$$

$$= \frac{\partial}{\partial \mu_h} \left[-\frac{1}{2} \sum_i t_{ih} \underbrace{(x_i - \mu_h)^T \Sigma_h^{-1} (x_i - \mu_h)}_{\cancel{\times \frac{1}{2}}} \right]$$

$$= -\frac{1}{2} \sum_i t_{ih} \cancel{\times \frac{1}{2}} \cancel{\times \sum_h} (x_i - \mu_h)$$

$$\frac{\partial Q(\theta)}{\partial \mu_h} = 0 \Leftrightarrow \sum_i t_{ih} \sum_h (x_i - \mu_h) = 0$$

$$\Leftrightarrow \sum_i t_{ih} (x_i - \mu_h) = 0$$

$$\Leftrightarrow \sum_i t_{ih} x_i = \sum_i t_{ih} \mu_h$$

$$\Leftrightarrow \mu_h^* = \frac{\sum_i t_{ih} \overline{x_i}}{\sum_i t_{ih}} = \frac{1}{m_h} \sum_{i=1}^m t_{ih} \overline{x_i}$$

Exo 3: max over Σ_h :

$$\begin{aligned}
 \frac{\partial Q}{\partial \Sigma_h} &= \frac{\partial}{\partial \Sigma_h} \left[-\frac{1}{2} \sum_{i=1}^m t_{ih} \left(\log |\Sigma_h| + (x_i - \mu_h)^t \Sigma_h^{-1} (x_i - \mu_h) \right) \right] \\
 &= \frac{\partial}{\partial \Sigma_h} \left[-\frac{1}{2} \sum_{i=1}^m t_{ih} \left(\log |\Sigma_h| + \text{tr}((x_i - \mu_h)^t \Sigma_h^{-1} (x_i - \mu_h)) \right) \right] \\
 &\quad = \frac{\partial}{\partial \Sigma_h} \left[\underbrace{-\frac{1}{2} \sum_{i=1}^m t_{ih} \log |\Sigma_h|}_{\text{constant}} + \underbrace{\frac{1}{2} \sum_{i=1}^m t_{ih} \text{tr}((x_i - \mu_h)^t \Sigma_h^{-1} (x_i - \mu_h))}_{\text{differentiation}} \right] \\
 &= \frac{\partial}{\partial \Sigma_h} \left[-\frac{1}{2} \sum_i t_{ih} \log |\Sigma_h| - \frac{1}{2} t_h \left(\Sigma_h^{-1} \sum_i t_{ih} (x_i - \mu_h) (x_i - \mu_h)^t \right) \right] \\
 &= -\frac{1}{2} \sum_i t_{ih} \left(\Sigma_h^{-1} \sum_i t_{ih} (x_i - \mu_h) (x_i - \mu_h)^t \right) - \frac{1}{2} t_h \left(\Sigma_h^{-1} S_h \Sigma_h^{-1} \right) \\
 &= -\frac{1}{2} \left(\sum_i t_{ih} (\Sigma_h^{-1})^t - (\Sigma_h^{-1} S_h \Sigma_h^{-1}) \right) \\
 &= -\frac{1}{2} \left(\sum_{i=1}^m t_{ih} \Sigma_h^{-1} - \Sigma_h^{-1} S_h \Sigma_h^{-1} \right) = \text{?} = 0
 \end{aligned}$$

Some elements about derivative according to matrices (cf. the Matrix Cook Book)

$$\text{tr}(AB) = \text{tr}(BA)$$

$$\text{tr}(\lambda) = \lambda$$

$$\text{tr}(A) = (A^{-1})^t$$

$$\text{trace}(A^t B) = -(A^t B A^{-1})^t$$

$$= \text{tr}(\Sigma_h^{-1} (x_i - \mu_h)(x_i - \mu_h)^t)$$

$$= -\frac{1}{2} \sum_i t_{ih} \log |\Sigma_h| - \frac{1}{2} t_h (\Sigma_h^{-1} S_h)$$

$$\Leftrightarrow m_k \sum_h^{-1} = \sum_h^{-1} S_h \sum_h^{-1}$$

$$\Leftrightarrow m_k \sum_h^{-1} \sum_h = \sum_h^{-1} S_h \sum_h^{-1} \sum_h$$

$$\Leftrightarrow m_k = \sum_h^{-1} S_h \quad (\Rightarrow \sum_h m_k = \sum_h \sum_h^{-1} S_h)$$

$$(\Rightarrow \sum_h m_k = S_k)$$

$$\sum_h^{-1} = \frac{S_h}{m_k} = \frac{\sum_i t_{ih} (x_{ih} - m_k) (x_{ih})^T}{\sum t_{ih}}$$

Exo 4: Max $Q(\theta)$ according to π_h given that $\sum_h \pi_h = 1$

$$L(\theta) = Q(\theta) - \omega \left(\sum_h \pi_h - 1 \right)$$

$$\frac{\partial L}{\partial \pi_h} = \frac{\sum_i t_{ih}}{\pi_h} - \omega \stackrel{(1)}{\Rightarrow}$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = 0 \Leftrightarrow \frac{m_k}{\pi_k} = \omega \Leftrightarrow w \pi_k = m_k$$

and we know that $\sum_k \pi_k = 1$

$$\sum_k w \pi_k = \sum_k m_k \Leftrightarrow w = m$$

Coming back to (1) $\Rightarrow \frac{m_k}{\pi_k} = m$

$$(\Rightarrow) \boxed{\pi_k^* = \frac{m_k}{m}}$$

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

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

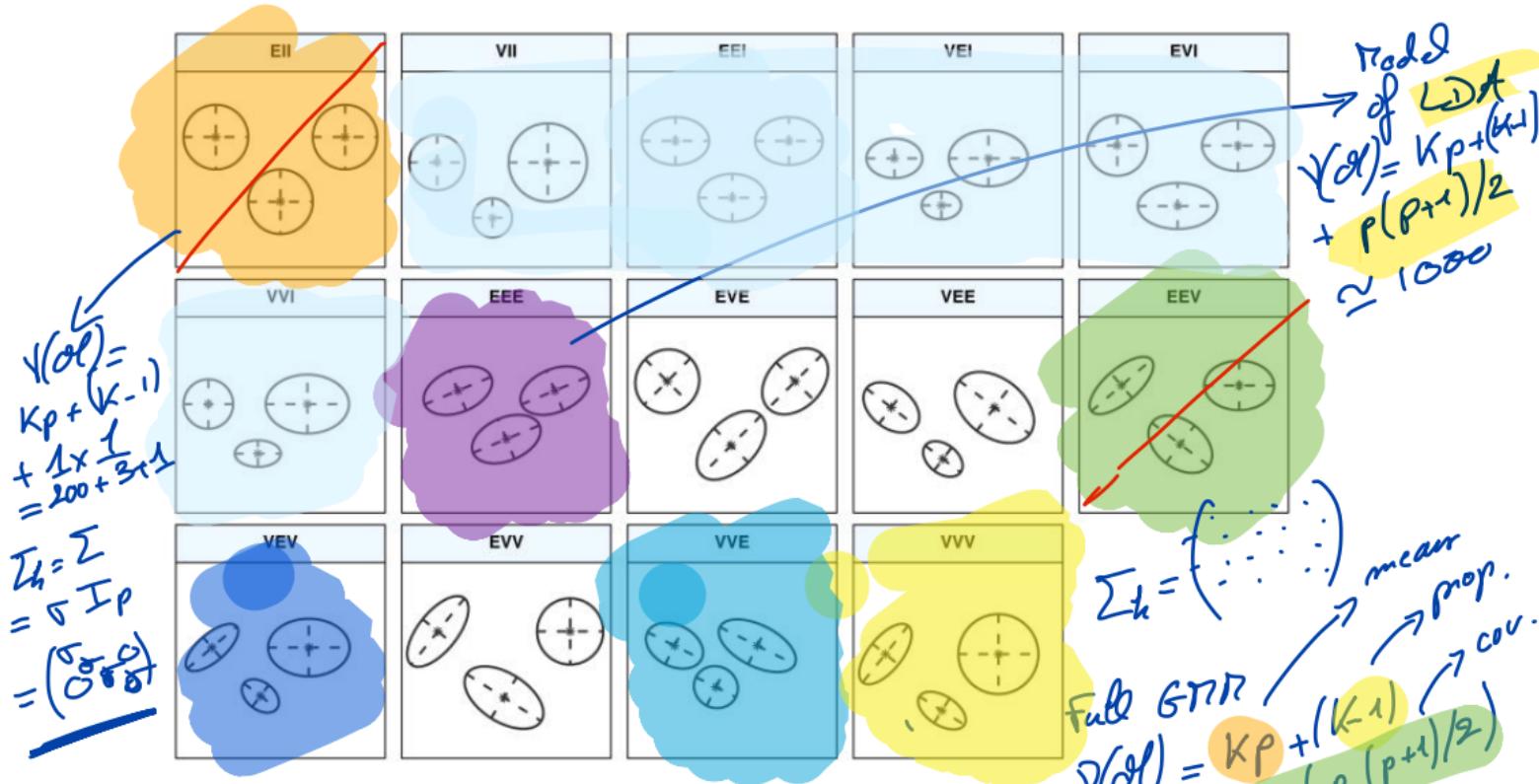


Figure: The parsimonious models of Mclust.

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 $VVV \rightarrow EEE$ nomenclature)
 - ↳ BIC is used for model selection.

- **(R)Flexmix** (R, 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.

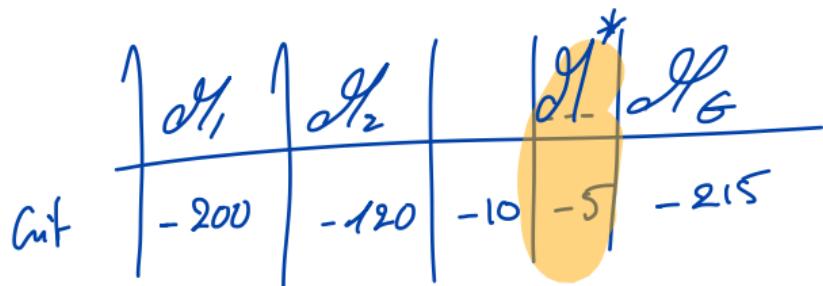
→ AIC

→ 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. $p(\mathcal{M}_g) = p \quad \forall g$.

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:

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

Remark: 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:

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

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 $B_{12} > 100$, 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.

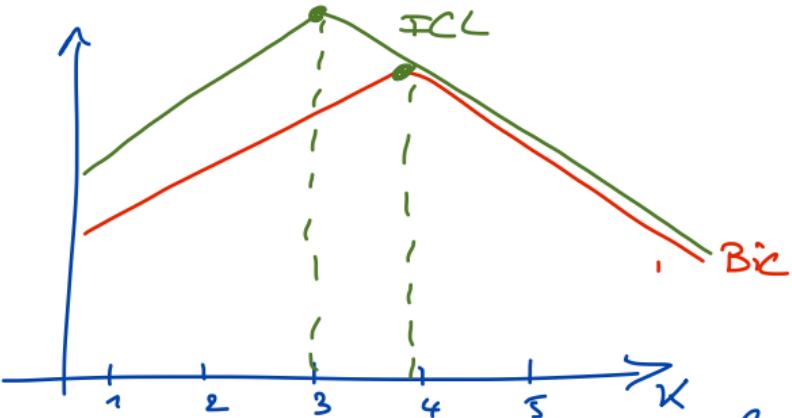
$$p(X|\theta_g) \rightarrow p(X, Z|\theta_g)$$

\uparrow
BIC \uparrow
 ICL.

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

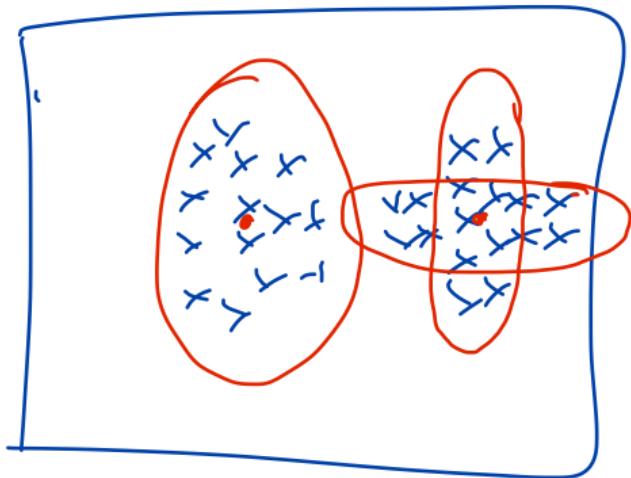
$$\log p(X, Z|\theta_g) \approx \left. \begin{aligned} & \log p(X|\hat{\theta}_g, \theta_g) - \frac{V(\theta_g)}{n} \log(n) \\ & - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K \hat{z}_{ik} \log^2(\hat{z}_{ik}) \end{aligned} \right\} \text{ICL.}$$

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.

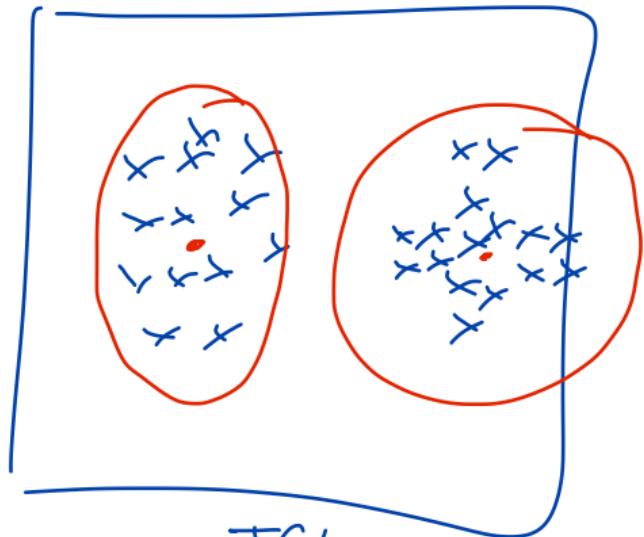


Rung: the objectives of Bic and ICL are slightly different:

- 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-