

Model-based Statistical Learning

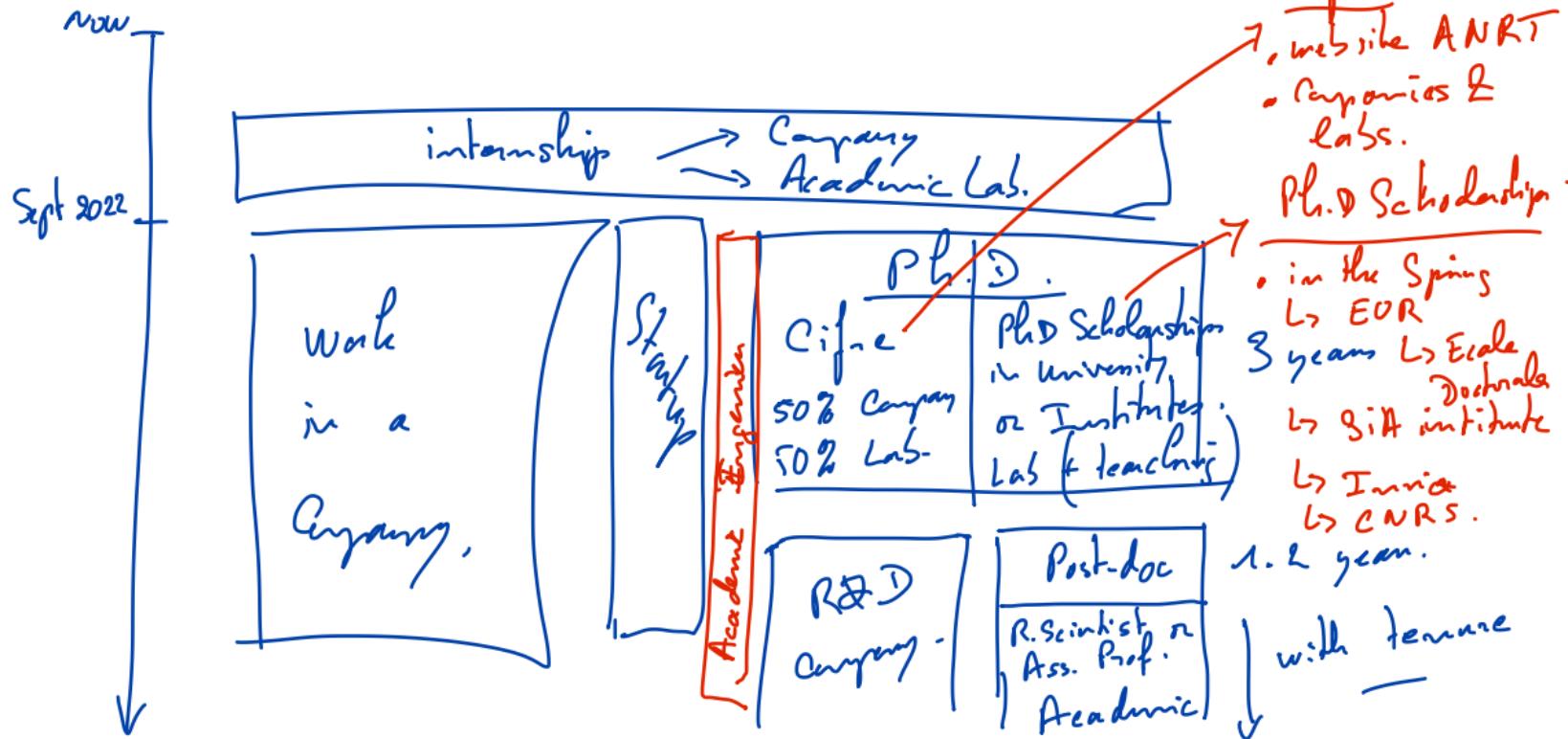


Pr. Charles BOUVEYRON

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)

Aftra the RSc :

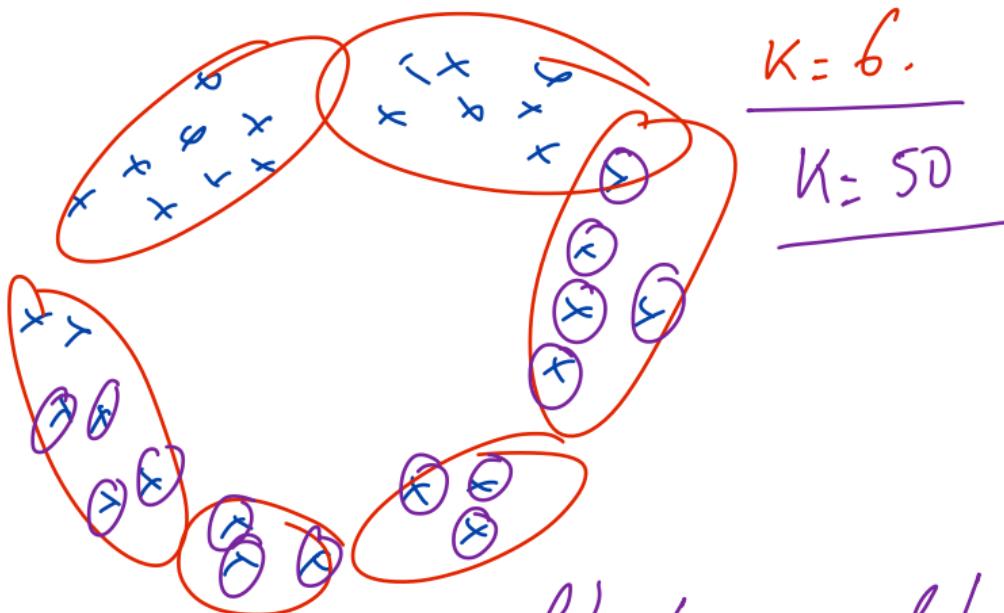


The Gaussian mixture model (GMM)

The GMM is probably the most popular Dist. model for two main reasons:

- a bad one: it is the DP for which the calculations are the simplest.
- a good one: even though it is a simple model it is flexible enough to fit a large variety of situations.

The GMM: $p(x; \Theta) = \sum_{k=1}^K \pi_k \phi(x; \mu_k, \Sigma_k)$



When varying K , we can fit to any data distribution. The limit will be to choose K in an appropriate manner.

Model inference in the Gaussian mixture model
is not easy because of the specific form
of the Log-Likelihood:

$$\begin{aligned}\text{Log } L(x; \theta) &= \text{Log} \left(\prod_{i=1}^n \sum_{h=1}^k \pi_h \phi(x_i; \mu_h, \Sigma_h) \right) \\ &= \sum_{i=1}^n \text{Log} \sum_{h=1}^k \pi_h \phi(\underline{\hspace{2cm}})\end{aligned}$$

\Rightarrow the classical solution is to use the EM algo.

The EM algorithm

The idea of the EM algo. is to first revisit the model by introducing a **latent variable** $z \in \{0,1\}^k$ to encode the class memberships.

$z_{ih} = \begin{cases} 1 & \text{if } x_i \text{ belongs to cluster } h. \\ 0 & \text{otherwise.} \end{cases}$

$\int z | \pi \sim \mathcal{D}(1; \pi)$ integrate over z
 $x | z=h \sim N(x; \mu_h; \Sigma_h)$ $\Rightarrow p(x|\theta) = \sum_{h=1}^k \pi_h \phi(x; \mu_h, \Sigma_h)$

$\rightarrow p(z=h) = \pi_h.$

The EM algorithm

Complete

This allows to write the likelihood of the couple (x, z) :

$$\log \mathcal{L}(x, z; \theta) = \sum_{i=1}^m \left[\log p(z_i | x_i; \theta) + \log p(x_i) \right]$$

$$= \log \mathcal{L}(x; \theta) + \sum_{i=1}^m \log p(z_i | x_i; \theta)$$

$$(\Rightarrow \log \mathcal{L}(x; \theta) = \log \mathcal{L}(x, z; \theta) - \sum_{i=1}^m \log p(z_i | x_i; \theta)$$

there, we demonstrated that the green part is a lower bound of the \log likelihood \Rightarrow if we know z , we can maximize the LB instead of the \log likelihood.

The EM algorithm

The spirit of the EM algorithm is to alternate two steps:

- the Expectation step: knowing a current value of θ^* , we compute the expectation of the lower bound :

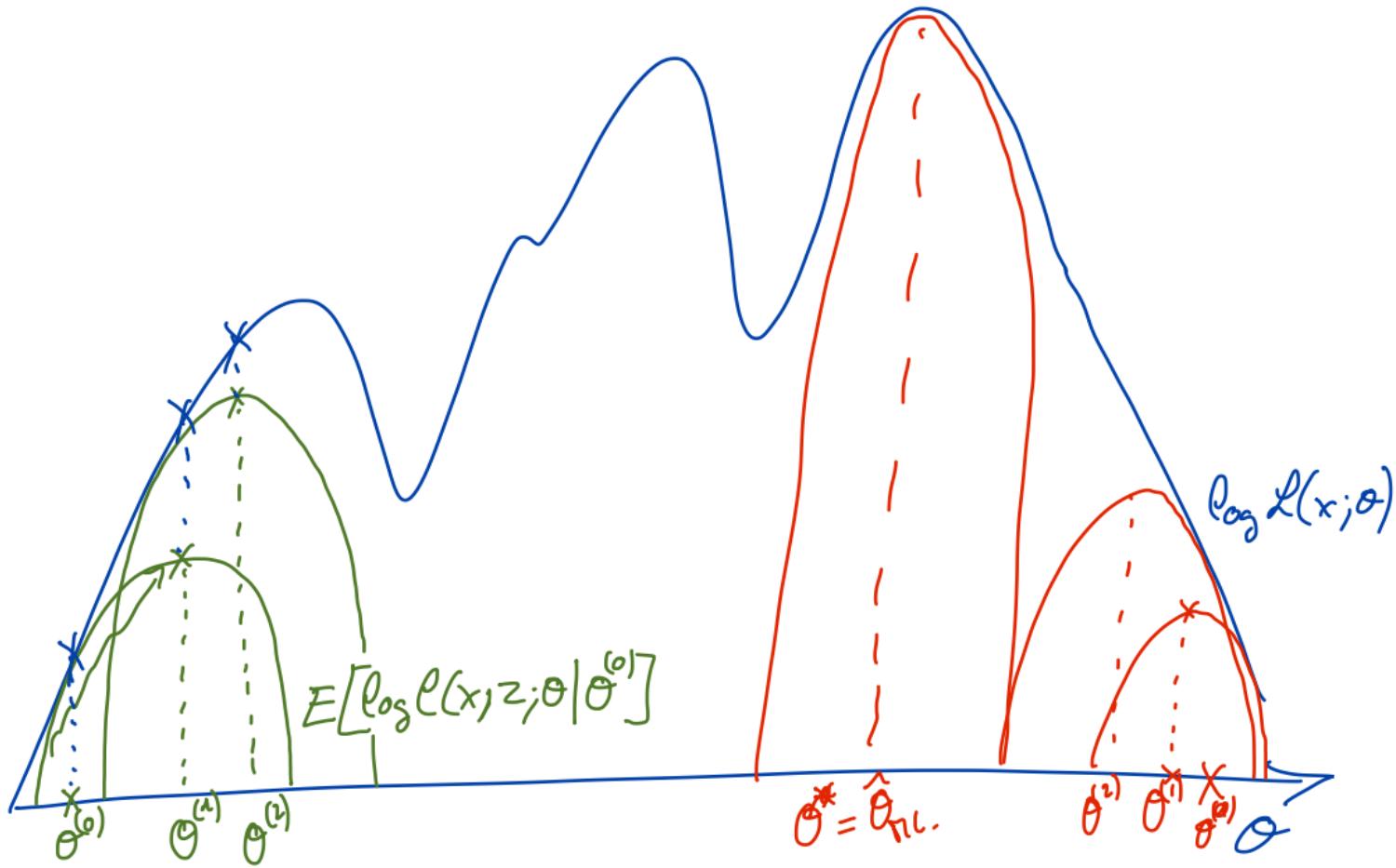
$$E[\ell(x, z | \theta) | \theta^*] = Q(\theta | \theta^*)$$

- the Maximization step : optimize $Q(\theta | \theta^*)$ over θ to obtain a new estimate of θ^{**}

The EM algorithm

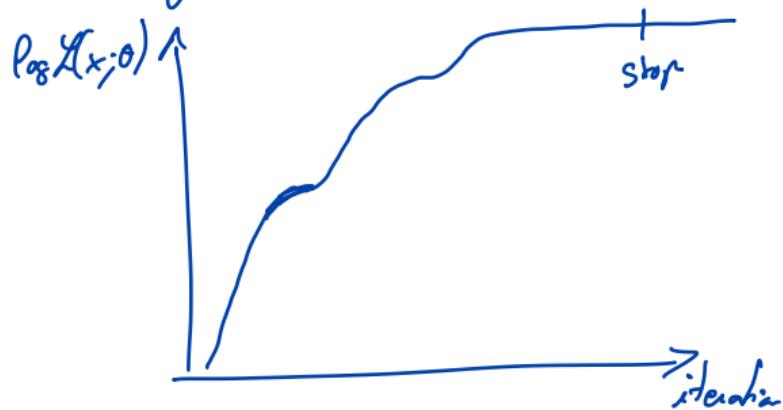
Theorem (Dempster, Laird, Rubin, 1977)
+ Wu, 1981)

The series of parameters $(\theta^*)_q$ generated by
the EM algorithm converges toward a local
maximum of the log-likelihood $\log L(x; \theta)$.



In practice, to avoid a sensitivity to the initialization, we recommend to try several random initializations for $\theta^{(0)}$ and retain at the end the θ^* that leads to the maximum likelihood.

In practice, we also stop the algorithm when a plateau of the likelihood is detected:



The EM algorithm for GMM

The central quantity to calculate in the E step is

$$Q(\theta | \theta^*) = E[\ell(x, z | \theta) | \theta^*]$$

$$\text{and } \ell(x, z | \theta) = \sum_{i=1}^n \sum_{h=1}^k z_{ih} \log(\pi_h \phi(x_i; \mu_h, \Sigma_h))$$

$$\text{and therefore: } Q(\theta | \theta^*) = \sum_i \sum_h E[z_{ih} | \theta^*] \log(\pi_h \phi(x_i; \mu_h, \Sigma_h))$$

So, the E step for the GMM reduces to the computation
of $z_{ih} = E[z_{ih} | \theta^*]$ Bayes $\propto \frac{P(z_{ih}=1 | \theta^*) p(x_i | z_{ih}=1, \theta^*)}{\pi_h^* \cdot \phi(x_i; \mu_h^*, \Sigma_h^*)}$

The EM algorithm for GMM

The E step:

$$\gamma_{ih} \leftarrow \pi_h^* \cdot \phi(x_i; \mu_h^*, \Sigma_h^*)$$

$$\forall i \\ \forall h.$$

The M step: Maximize over π_h, μ_h, Σ_h the function

$$Q(\theta|\theta^*) = \sum_i \sum_h \gamma_{ih} \log (\pi_h \phi(x_i; \mu_h, \Sigma_h))$$

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

where d is the dimensionality of $x_i \in \mathbb{R}^d$

The EM algorithm for GMM

The update equations for π_h , μ_h and Σ_h can be obtained by simply taking the partial derivatives of $Q(\theta|\theta^*)$ regarding π_h , μ_h and Σ_h respectively and equalling to 0

$$\frac{\partial}{\partial \mu_h} Q(\theta|\theta^*) = 0 \iff \mu_h^* = \frac{1}{m_h} \sum_{i=1}^{m_h} \gamma_{ih} x_i \quad \text{where } m_h = \sum \gamma_{ih}.$$

$$\frac{\partial}{\partial \Sigma_h} Q(\theta|\theta^*) = 0 \iff \Sigma_h^* = \frac{1}{m_h} \sum_{i=1}^{m_h} \gamma_{ih} (x_i - \mu_h^*)(x_i - \mu_h^*)^T$$

$$\frac{\partial}{\partial \pi_h} Q(\) \underset{\sum_h \pi_h = 1}{\text{under the constraint}} = 0 \iff \pi_h^* = \frac{m_h}{n}$$

The EM algorithm for GMM

Compute $\frac{\partial}{\partial \mu_k} Q(\theta | \theta^*)$ where

$$Q(\theta | \theta^*) = \sum_i \sum_h \gamma_{ih} \log (\pi_h \phi(x_i; \mu_h, \Sigma_h))$$

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

$$\text{so } Q(\theta | \theta^*) = \sum_i \sum_h \gamma_{ih} \left[\log \pi_h + \log \phi(x_i; \mu_h, \Sigma_h) \right]$$

$$= \sum_i \sum_h \gamma_{ih} \left[\log \pi_h - d \log (2\pi) - \frac{1}{2} \log |\Sigma_h| \right. \\ \left. - \frac{1}{2} (x_i - \mu_h)^T \Sigma_h^{-1} (x_i - \mu_h) \right]$$

$$\frac{\partial}{\partial \mu_k} Q(\theta | \theta^*) = \sum_i \gamma_{ih} \frac{\partial}{\partial \mu_k} \left[-\frac{1}{2} (x_i - \mu_h)^T \Sigma_h^{-1} (x_i - \mu_h) \right]$$

cf. the "book": The Matrix Cook Book.

$$\frac{\partial}{\partial v} v^t B v = 2 B v$$

$$\begin{aligned}\frac{\partial}{\partial \mu_k} Q &= \sum_i \gamma_{ik} \left[-\cancel{\frac{\partial}{\partial x}} \sum_h^{-1} (x_i - \mu_h) \right] \\ &= -\sum_i \gamma_{ik} \sum_h^{-1} (x_i - \mu_h)\end{aligned}$$

and now $\frac{\partial}{\partial \mu_k} Q = 0 \quad (\Rightarrow \sum_i \gamma_{ik} \sum_h^{-1} (x_i - \mu_h) = 0)$

$$(\Rightarrow \sum_i \gamma_{ik} (x_i - \mu_h) = 0$$

$$\Rightarrow \mu_k^* = \frac{\sum_i \gamma_{ik} x_i}{\sum_i \gamma_{ik}} = \frac{1}{n_k} \sum_i \gamma_{ik} x_i$$

Home work: do the calculations for $\sum_h^* \frac{\partial}{\partial A} \text{trace}(A^{-1}B) = -(A^{-1}BA^{-1})^{-1}$

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

$$\text{tr}(\lambda) = \lambda \quad \frac{\partial}{\partial A} \log |A| = (A^{-1})^T$$

Model selection: how to choose K ?

Model selection criteria all penalized the likelihood with a quantity which favors models with a low number of groups.

$$\text{MSCriteria} = \log L(\mathbf{x}; \hat{\boldsymbol{\theta}}) - \text{pen}(K)$$

Two popular criteria:

$$\text{AIC} = \log L(\mathbf{x}; \hat{\boldsymbol{\theta}}) - n(\boldsymbol{\theta})$$

$$\text{BIC} = \log L(\mathbf{x}; \hat{\boldsymbol{\theta}}) - \frac{1}{2} n(\boldsymbol{\theta}) \log(n)$$

where $n(\boldsymbol{\theta})$ is the number of free scalar parameters in the model $\boldsymbol{\theta}$.

In practice, $\gamma(GNN)$ is easy to compute

$$\begin{aligned}\gamma(GNN \text{ with } K) &= \text{mb of } T_h + \text{mb of } \mu_k + \text{mb of } I_k \\ &= (K-1) + Kd + K \frac{d(d+1)}{2}\end{aligned}$$

$$BIC(GNN) = \log L(x; \hat{\theta}) - \frac{(K-1) + Kd + K \frac{d(d+1)}{2}}{2} \log(n)$$

