

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

EM algorithm for GMM:

The goal of the first step is to find the best mixture parameter estimates for the data at hand:

$$X \xrightarrow[\text{with EM also}]{\text{learn}} \hat{\theta} = \left\{ \begin{array}{l} \hat{\pi}_1, \dots, \hat{\pi}_K \\ \hat{\mu}_1, \dots, \hat{\mu}_K \\ \hat{\Sigma}_1, \dots, \hat{\Sigma}_K \end{array} \right\}$$

The EM algorithm will optimize iteratively the likelihood of this model without directly optimize the function.

The EM algorithm

Starting with the GMM model, we can write the log-likelihood of the model:

$$\begin{aligned}\mathcal{L}(x; \theta) &= \log \left(\prod_{i=1}^m p(x_i; \theta) \right) \\ &= \log \left(\prod_{i=1}^m \sum_{h=1}^K \pi_h \mathcal{N}(x_i; \mu_h, \Sigma_h) \right) \\ &= \sum_{i=1}^m \log \left(\sum_{h=1}^K \pi_h \mathcal{N}(x_i; \mu_h, \Sigma_h) \right)\end{aligned}$$

We see here that it is totally possible to evaluate $\mathcal{L}(x; \theta)$ for specific values of x and θ , but the direct optimization is really difficult due to the $\log(\sum)$

The EM algorithm

The idea of the EM algorithm is to introduce an extra and non-observed (latent) variable Y , encoding the group memberships. Estimating both θ and Y from X is finally easier than just estimating θ from X .

$$Y_i \in \{0, 1\}^K \Rightarrow Y_i = (0, 0, 1, 0) \Rightarrow x_i \text{ belongs to the 3rd cluster}$$

$$\begin{cases} Y \sim \mathcal{H}(1; \pi) \\ X|Y=k \sim \mathcal{N}(\mu_k, \Sigma_k) \end{cases}$$

$$\begin{aligned} &\text{integrate over } Y \\ &\Rightarrow p(X) = \sum_{k=1}^K \pi_k \mathcal{N}(x; \mu_k, \Sigma_k) \end{aligned}$$

The EM algorithm

This allows us to write another log-likelihood, called the complete data log-likelihood, of the couple (x, y) :

$$\begin{aligned} &= \sum_{i=1}^n \log p(x_i, y_i; \theta) \\ \mathcal{L}_c(x, y, \theta) &= \sum_{i=1}^n \left[\log p(y_i | x_i; \theta) + \log p(x_i; \theta) \right] \\ &= \mathcal{L}(x; \theta) + \sum_{i=1}^n \log p(y_i | x_i; \theta) \end{aligned}$$

$$\Rightarrow \mathcal{L}(x; \theta) = \underbrace{\mathcal{L}_c(x, y; \theta)} - \sum_{i=1}^n \log p(y_i | x_i; \theta)$$

$$\Rightarrow \mathcal{L}_c \leq \mathcal{L}$$

The EM algorithm

Rule: we see here that L_c is a lower bound of L and optimizing L_c over Θ will automatically lead in the optimization of L .

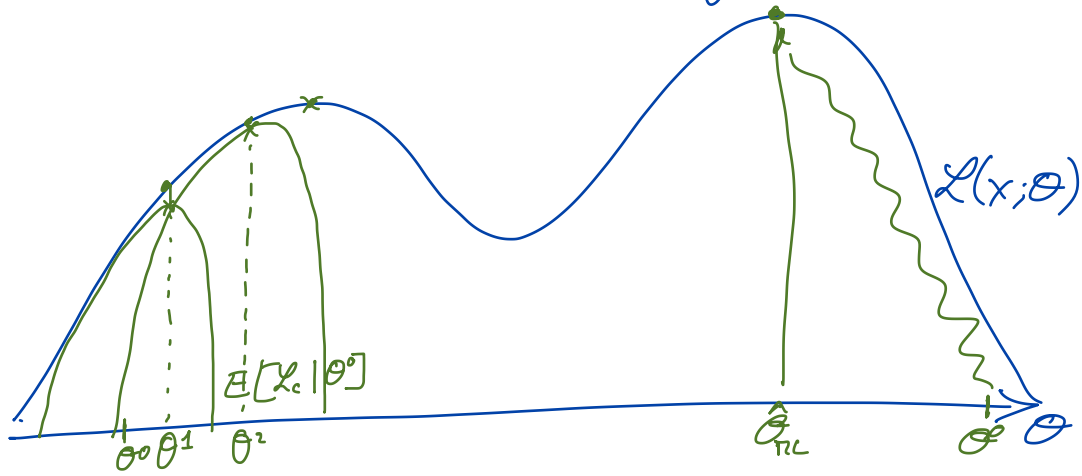
Thanks to this remark, Dempster, Laird and Rubin proposed in 1977 the EM algorithm:

- E step: the E step aims at calculating the expected complete log-likelihood:

$$Q(\theta; \theta^*) = E[\mathcal{L}_c(x, y; \theta) | x; \theta^*]$$

- M step: the M step aims at maximizing this function $Q(\theta; \theta^*)$ over θ to provide a new value for θ^*

Theorem: the sequence of estimates (θ^*) over the iterations of the EM algorithm is converging toward a local maximum of the log-likelihood.



The EM algorithm for GMM

In practice :

- 1) to avoid being trapped in a local maximum, we usually do several (10) different random initializations and we keep afterward the $\hat{\Theta}_{EM}$ with the highest likelihood.
- 2) to stop the algorithm, we just monitor the evolution of the log. likelihood and we stop when a plateau is detected



The EM algorithm for GMM

Starting with the E step, we need to focus on

$$Q(\theta, \theta^*) = E[\mathcal{L}_c(x, y; \theta) | \theta^*, X]$$

$$\begin{aligned} \text{and } \mathcal{L}_c(x, y; \theta) &= \sum_{i=1}^m \log p(x_i, y_i; \theta) \\ &= \sum_{i=1}^m \log \left[\sum_{k=1}^K y_{ik} \pi_k \mathcal{N}(x_i; \mu_k, \Sigma_k) \right] \\ &= \sum_{i=1}^m \sum_{k=1}^K y_{ik} \log(\pi_k \mathcal{N}(x_i; \mu_k, \Sigma_k)) \end{aligned}$$

and therefore:

$$E[\mathcal{L}_c(x, y | \theta) | \theta^*; X] = \sum_{i=1}^m \sum_{h=1}^K \frac{E[y_{ih} | \theta^*, X]}{\pi_h N(x_i; \mu_h, \Sigma_h)} \log \left(\pi_h N(x_i; \mu_h, \Sigma_h) \right)$$

and $\underline{E[y_{ih} | \theta^*, X]} = P(y_{ih}=1 | x, \theta^*)$

$$\begin{aligned} &\stackrel{\text{Bayes}}{=} \frac{P(y_{ih}=1 | \theta^*) P(x_i | y_{ih}=1; \theta^*)}{p(x)} \\ &\propto \frac{\pi_h^* \times N(x_i; \mu_h^*, \Sigma_h^*)}{\pi_h^*} \end{aligned}$$

Let's call $t_{ih} = E[y_{ih} | \theta^*, X]$

$$\Rightarrow E[\mathcal{L}_c(x, y; \theta) | \theta^*, X] = \sum_{i=1}^m \sum_{h=1}^K t_{ih} \log \pi_h N(x_i; \mu_h, \Sigma_h)$$

The EM algorithm for GMM

In the π step, we just have to optimize in θ the $E[\mathcal{L}_c(x, y; \theta) | \theta^*; X]$:

$$\max_{\theta} E[\mathcal{L}_c(x, y; \theta) | \theta^*; X] = \sum_{i=1}^n \sum_{k=1}^K \text{tikh} \rho_{\theta} \pi_k \mathcal{N}(x_i; \mu_k, \Sigma_k)$$

$Q(\theta; \theta^*) =$

• Finding the update for π_k :

$$\frac{\partial}{\partial \pi_k} Q(\theta; \theta^*) = \frac{\partial}{\partial \pi_k} \left[\sum_i \sum_k \text{tikh} \left[\rho_{\theta} \pi_k + \rho_{\theta} \mathcal{N}(x_i, -) \right] \right]$$