

The Expectation-Maximization algorithm

Masmoudi Afif

Laboratoire de probabilités et statistique
Université de Sfax
Journées de Statistique Mathématique et Data Science
(JSMDS 2019)

Plan

- 1 Introduction
- 2 A Finite Gaussian mixture models
 - Parameters estimation
- 3 The EM Algorithm
 - The EM algorithm properties

Plan

- 1 Introduction
- 2 A Finite Gaussian mixture models
 - Parameters estimation
- 3 The EM Algorithm
 - The EM algorithm properties

Plan

- 1 Introduction
- 2 A Finite Gaussian mixture models
 - Parameters estimation
- 3 The EM Algorithm
 - The EM algorithm properties

Plan

- 1 Introduction
- 2 A Finite Gaussian mixture models
 - Parameters estimation
- 3 The EM Algorithm
 - The EM algorithm properties

Introduction

- A finite mixture models is a combination of two or more probability density functions.
- It is a probabilistic model used to represent a subpopulation within an overall population.
- It represents a natural extension of the finite Gaussian mixture of distributions.
- It has been adopted in such diverse areas as ecology, bioinformatics, computer science and biostatistics..

Plan

- 1 Introduction
- 2 A Finite Gaussian mixture models
 - Parameters estimation
- 3 The EM Algorithm
 - The EM algorithm properties

A Finite Gaussian mixture model

Gaussian mixture distribution

The distribution f is a mixture of K components Gaussian distributions defined by

$$f(x) = f(x; \Theta) = \sum_{k=1}^K \pi_k N_k(x; \mu_k, \sigma_k) \quad (1)$$

where,

- π_k is the mixing weights, verify $0 < \pi_k < 1$, and $\sum_{k=1}^K \pi_k = 1$.

A Finite Gaussian mixture models

- $\Theta = (\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \sigma_1, \dots, \sigma_K)$ is the parameter vector of the Gaussian mixture models.

Plan

- 1 Introduction
- 2 A Finite Gaussian mixture models
 - Parameters estimation
- 3 The EM Algorithm
 - The EM algorithm properties

Parameters estimation

Let (X_1, \dots, X_N) be N independent random variables with the same Gaussian mixture density $f(\cdot; \Theta)$ and let (x_1, \dots, x_N) be N associated observations.

The incomplete likelihood function

The incomplete likelihood function l is given by

$$\begin{aligned} l(x_1, x_2, \dots, x_N; \Theta) &= \prod_{i=1}^N f(x_i; \Theta) \\ &= \prod_{i=1}^N \left(\sum_{k=1}^K \pi_k \frac{1}{\sigma_k \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - \mu_k}{\sigma_k} \right)^2} \right) \end{aligned}$$

Parameters estimation

The incomplete log -likelihood

the log-likelihood L is given by

$$\begin{aligned}L(x_1, x_2, \dots, x_N; \Theta) &= \log l(x_1, x_2, \dots, x_N; \Theta) \\ &= \sum_{i=1}^N \log \left(\sum_{k=1}^K \pi_k N_k(x_i; \sigma_k, \mu_k) \right) \\ &= \sum_{i=1}^N \log \left(\sum_{k=1}^K \pi_k \frac{1}{\sigma_k \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - \mu_k}{\sigma_k} \right)^2} \right)\end{aligned}$$

Parameters estimation

For each random variable X_i , we associate a random vector $Z_i = (Z_{i1}, \dots, Z_{iK})$ following a multivariate Bernoulli distribution with vector parameters (π_1, \dots, π_K) . That is

$$\mathbb{P}(Z_i = z_i) = \prod_{k=1}^K \pi_k^{z_{ik}}, \quad (2)$$

where

- $z_i = (z_{i1}, \dots, z_{iK}) \in \{0, 1\}^K$, and $\sum_{k=1}^K z_{ik} = 1$.

Parameters estimation

The maximum likelihood function from complete data

The maximum likelihood function from complete data l_C is given by

$$\begin{aligned}
 l_C(x_1, \dots, x_N, z_1, \dots, z_N; \Theta) &= \prod_{i=1}^N \mathbb{P}(Z_i = z_i) \prod_{k=1}^K N_k^{z_{ik}}(x_i; \sigma_k, \mu_k) \\
 &= \prod_{i=1}^N \left(\prod_{k=1}^K (\pi_k N_k(x_i; \sigma_k, \mu_k))^{z_{ik}} \right) \\
 &= \prod_{i=1}^N \left(\prod_{k=1}^K \left(\pi_k \frac{1}{\sigma_k \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - \mu_k}{\sigma_k} \right)^2} \right)^{z_{ik}} \right).
 \end{aligned}$$

The log-likelihood function from complete data

The log-likelihood function from complete data L_C is given by

$$\begin{aligned}
 L_C(x_1, \dots, x_N, z_1, \dots, z_N; \Theta) &= \sum_{i=1}^N \sum_{k=1}^K z_{ik} \log(\pi_k N_k(x_i; \sigma_k, \mu_k)) \\
 &= \sum_{i=1}^N \sum_{k=1}^K z_{ik} \log\left(\pi_k \frac{1}{\sigma_k \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - \mu_k}{\sigma_k}\right)^2}\right)
 \end{aligned}$$

where $\sum_{k=1}^K z_{ik} = 1$.

The conditional expectation of the complete-data log likelihood

$$\begin{aligned}
 Q(\theta || \theta^{(l)}) &= \mathbb{E}(L_C(x_1, \dots, x_N, z_1, \dots, z_N; \theta^{(l)}) | x_1, \dots, x_N) \\
 &= \sum_{i=1}^N \sum_{k=1}^K \tau_{ik}^{(l)} \log(\pi_k N_k(x_i; \mu_k, \sigma_k)) \\
 &= \sum_{i=1}^N \sum_{k=1}^K \tau_{ik}^{(l)} \log\left(\pi_k \frac{1}{\sigma_k \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - \mu_k}{\sigma_k}\right)^2}\right)
 \end{aligned}$$

Where

The conditional expectation of the complete-data log likelihood

The posterior probability

$$\begin{aligned}\tau_{ik}^{(l)} &= \mathbb{E}\left(\mathbf{Z}_{ik} \mid \Theta^{(l)}, \mathbf{X}_1 = x_1, \dots, \mathbf{X}_N = x_N\right) \\ &= \mathbb{E}\left(\mathbf{Z}_{ik} = z_{ik} \mid \mathbf{X}_i = x_i, \Theta^{(l)}\right) \\ &= \frac{\pi_k^{(l)} \mathcal{N}_k(x_i; \mu_k^{(l)}, \sigma_k^{(l)})}{\sum_{k=1}^K \pi_k^{(l)} \mathcal{N}_k(x_i; \mu_k^{(l)}, \sigma_k^{(l)})},\end{aligned}$$

The conditional expectation of the complete-data log likelihood

and

$\Theta = (\mu_1, \dots, \mu_K, \pi_1, \dots, \pi_K, \sigma_1, \dots, \sigma_K)$ is the parameter vector of the Gaussian mixture models in the l^{th} iteration.

The conditional expectation of the complete-data log likelihood

Challenge

Maximization of the conditional expectation log likelihood from complete data $Q(\Theta || \Theta^{(l)})$.

Plan

- 1 Introduction
- 2 A Finite Gaussian mixture models
 - Parameters estimation
- 3 The EM Algorithm
 - The EM algorithm properties

The EM Algorithm

Definition

The EM algorithm is an iterative method for finding maximum likelihood or maximum a posteriori (MAP) estimates of parameters in statistical models involving incomplete data.

The EM algorithm

The EM algorithm steps

- **Step 1 : Initialization:** $\Theta^{(0)}$
- **Step 2 : Expectation(E)**

We calculate the posterior probability

$$\tau_{ik}^{(l)} = \frac{\pi_k^{(l)} \mathcal{N}_k(x_i; \mu_k^{(l)}, \sigma_k^{(l)})}{\sum_{k=1}^K \pi_k^{(l)} \mathcal{N}_k(x_i; \mu_k^{(l)}, \sigma_k^{(l)})} \quad (3)$$

The EM algorithm

The EM algorithm steps

and we calculate

$$Q(\Theta || \Theta^{(l)}) = \sum_{i=1}^N \sum_{k=1}^K \tau_{ik}^{(l)} \log (\pi_k N_k(x_i; \mu_k, \sigma_k)). \quad (4)$$

- **Step 3 : Maximization (M)**

$$\Theta^{(l+1)} = \arg \max_{\Theta} Q(\Theta || \Theta^{(l)})$$

The EM algorithm

The EM algorithm steps

The maximum likelihood estimators of μ_k and π_k are given respectively by

$$\mu_k^{(l+1)} = \frac{\sum_{i=1}^N \tau_{ik}^{(l)} x_i}{\sum_{i=1}^N \tau_{ik}^{(l)}} \quad (5)$$

$$\pi_k^{(l+1)} = \frac{1}{N} \sum_{i=1}^N \tau_{ik}^{(l)} \quad (6)$$

The EM algorithm

The EM algorithm steps

The maximum likelihood estimator of σ_k is given by

$$\sigma_k^{(l+1)} = \sqrt{\frac{\sum_{i=1}^N \tau_{ik}^{(l)} (x_i - \mu_k^{(l+1)})^2}{\sum_{i=1}^N \tau_{ik}^{(l)}}} \quad (7)$$

Plan

- 1 Introduction
- 2 A Finite Gaussian mixture models
 - Parameters estimation
- 3 The EM Algorithm
 - The EM algorithm properties

The EM algorithm properties

Proposition 1

At each iteration of the EM algorithm,

$$L(x_1, x_2, \dots, x_N; \Theta^{(l+1)}) \geq L(x_1, x_2, \dots, x_N; \Theta^{(l)}), \quad (8)$$

where, $\Theta^{(l+1)}$ is the parameter vector in the $(l + 1)$ iteration.

Then,

$$Q(\Theta^{(l+1)} || \Theta^{(l)}) \geq Q(\Theta^{(l)} || \Theta^{(l)}) \quad (9)$$

Proof of Proposition 1

Let $\mathbb{Z} = (Z_1 = z_1, \dots, Z_N = z_N)$ the missing data vector and $\mathbb{X} = (X_1 = x_1, \dots, X_N = x_N)$ the observed data vector.

We can see that

$$\begin{aligned}
 \log l(\mathbb{X}|\Theta^{(l)}) &= \sum_{\mathbb{Z}} l(\mathbb{Z}|\mathbb{X}, \Theta^{(l)}) \log \frac{l_c(\mathbb{Z}, \mathbb{X}; \Theta^{(l)})}{l(\mathbb{Z}|\mathbb{X}, \Theta^{(l)})} \\
 &= \sum_{\mathbb{Z}} l(\mathbb{Z}|\mathbb{X}, \Theta^{(l)}) \log l_c(\mathbb{Z}, \mathbb{X}|\Theta^{(l)}) \\
 &\quad - \sum_{\mathbb{Z}} l(\mathbb{Z}|\mathbb{X}, \Theta^{(l)}) \log l(\mathbb{Z}|\mathbb{X}, \Theta^{(l)}) \\
 &= \mathbb{E}(\log l_c(\mathbb{Z}, \mathbb{X}|\Theta^{(l)})|\mathbb{X}, \Theta^{(l)}) \\
 &\quad - \mathbb{E}(\log l(\mathbb{Z}|\mathbb{X}, \Theta^{(l)})|\mathbb{X}, \Theta^{(l)}).
 \end{aligned}$$

Proof of Proposition 1

So, we have

$$\begin{aligned} \mathbb{Q}(\Theta^{(l)} || \Theta^{(l)}) &= \mathbb{E}(\log l_c(\mathbb{Z}, \mathbb{X} | \Theta^{(l)}) | \mathbb{X}, \Theta^{(l)}) \\ &= \mathbb{H}(\Theta^{(l)} || \Theta^{(l)}) + \log l(\mathbb{X} | \Theta^{(l)}), \end{aligned}$$

Proof of Proposition 1

where

$$\mathbb{H}(\Theta^{(l)} || \Theta^{(l)}) = \mathbb{E}(\log l(\mathbb{Z} | \mathbb{X}, \Theta^{(l)}) | \mathbb{X}, \Theta^{(l)}). \quad (10)$$

Proof of Proposition 1

and further we have

$$\begin{aligned}\log l(\mathbb{X}; \Theta) &= \log \left(\sum_{\mathbb{Z}} l_c(\mathbb{Z}, \mathbb{X} | \Theta) \right) \\ &= \log \left(\sum_{\mathbb{Z}} l_c(\mathbb{Z}, \mathbb{X} | \Theta) \frac{l(\mathbb{Z} | \mathbb{X}, \Theta^{(l)})}{l(\mathbb{Z} | \mathbb{X}, \Theta^{(l)})} \right) \\ &= \log \mathbb{E} \left[\frac{l_c(\mathbb{Z}, \mathbb{X} | \Theta)}{l(\mathbb{Z} | \mathbb{X}, \Theta^{(l)})} | \mathbb{X}, \Theta^{(l)} \right].\end{aligned}$$

Proof of Proposition 1

As though, the log function is concave so, by Jensen's inequality, we obtain

$$\begin{aligned}\log l(\mathbb{X}|\Theta) = \log l(\Theta) &\geq \mathbb{E}\left[\log\left(\frac{l_c(\mathbb{Z}, \mathbb{X}|\Theta)}{l(\mathbb{Z}|\mathbb{X}, \Theta^{(l)})}\right) \middle| \mathbb{X}, \Theta^{(l)}\right] \\ &\geq \mathbb{E}\left[\log l_c(\mathbb{Z}, \mathbb{X}|\Theta) \middle| \mathbb{X}, \Theta^{(l)}\right] \\ &\quad - \mathbb{E}\left[\log(l(\mathbb{Z}|\mathbb{X}, \Theta^{(l)})) \middle| \mathbb{X}, \Theta^{(l)}\right].\end{aligned}$$

Proof Proposition 1

Hence, we get the following result

$$\log l(\mathbb{X}|\Theta) \geq Q(\Theta|\Theta^{(l)}) - \mathbb{H}(\Theta^{(l)}|\Theta^{(l)}). \quad (11)$$

Since, we have

$$\Theta^{(l+1)} = \arg \max_{\Theta} Q(\Theta|\Theta^{(l)}).$$

Then, as result

$$Q(\Theta^{(l+1)}|\Theta^{(l)}) \geq Q(\Theta^{(l)}|\Theta^{(l)}) \quad (12)$$

Proof Proposition 1

we can conclude that

$$\begin{aligned} & \log l(\mathbb{X}|\Theta^{(l+1)}) - \log l(\mathbb{X}|\Theta^{(l)}) \\ \geq & \mathbb{Q}(\Theta^{(l+1)}|\Theta^{(l)}) - \mathbb{H}(\Theta^{(l)}|\Theta^{(l)}) - \mathbb{Q}(\Theta^{(l)}|\Theta^{(l)}) + \mathbb{H}(\Theta^{(l)}|\Theta^{(l)}) \\ \geq & \mathbb{Q}(\Theta^{(l+1)}|\Theta^{(l)}) - \mathbb{Q}(\Theta^{(l)}|\Theta^{(l)}) \geq 0 \end{aligned}$$

which completes this proof.

The EM algorithm properties

Theorem 1

Suppose that, we have

$$1) \lim_{l \rightarrow +\infty} \Theta^{(l)} = \hat{\Theta}.$$

$$2) \frac{\partial Q(\Theta^{(l+1)} || \Theta^{(l)})}{\partial \Theta} = 0.$$

and

3) $\frac{\partial^2 Q(\Theta^{(l+1)} || \Theta^{(l)})}{\partial \Theta^2}$ is negative definite with eigenvalues bounded away from zero.

Then,

$\frac{\partial l(\hat{\Theta})}{\partial \Theta} = 0$ and $\frac{\partial^2 Q(\hat{\Theta} || \hat{\Theta})}{\partial \Theta^2}$ is negative definite.

The EM algorithm properties

According to Theorem 1

$\hat{\Theta}$ is the maximum likelihood estimator parameter vector of \mathcal{Q} function.

and

The EM likelihood sequence $l(\Theta^{(l)})$ converges to some local maximum $l(\hat{\Theta})$.

The EM algorithm properties

Theorem 2

If $L(\Theta)$ is unimodal and satisfies the same differentiability conditions in Theorem 1 then, $\Theta^{(l)}$ converges to the unique maximum likelihood of Θ .

The EM algorithm properties

Dempster et al. (1977) defined an EM mapping

$$M : \Theta \longrightarrow \Theta \quad (13)$$

which verified at each iteration (l),

$$\Theta^{(l+1)} = M(\Theta^{(l)}). \quad (14)$$

Theorem 3

If $\Theta^{(l)}$ converges to some point $\hat{\Theta}$, then $\hat{\Theta}$ is a fixed point of the EM algorithm.