

Expectation- Maximization Algorithm

Mohammad Emtiyaz Khan
EPFL

Nov 5, 2015



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

©Mohammad Emtiyaz Khan 2015

Motivation

Computing maximum likelihood for Gaussian mixture model is difficult due to the log outside the sum.

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) := \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\boldsymbol{\theta} = \{ \boldsymbol{\mu}_{1:K}, \boldsymbol{\Sigma}_{1:K}, \boldsymbol{\pi}_{1:K} \}$$

$$\sum_{n=1}^N (y_n - \beta \bar{x}_n)^2$$

$$\sum_n \log \sum_k \pi_k e^{-\frac{1}{2} \frac{(x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)}{2}}$$

$$\sum_n \sum_k \left(\log \pi_k - \frac{1}{2} \frac{(x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)}{2} \right)$$

Expectation-Maximization (EM) algorithm provides an elegant and general method to optimize such optimization problems. It uses an iterative two-step procedure where individual steps usually involve problems that are easy to optimize.

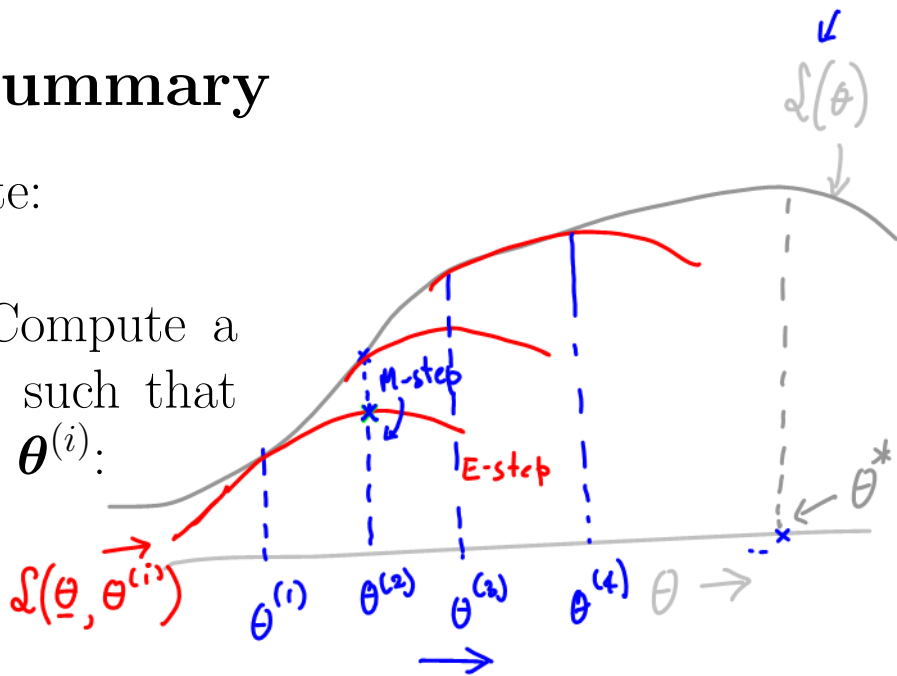
EM algorithm: Summary

Start with $\boldsymbol{\theta}^{(1)}$ and iterate:

- (Expectation step) Compute a lower bound to the cost such that it is tight at the previous $\boldsymbol{\theta}^{(i)}$:

$$\mathcal{L}(\boldsymbol{\theta}) \geq \underline{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i)}) \text{ and}$$

$$\underline{\mathcal{L}}(\boldsymbol{\theta}^{(i)}, \boldsymbol{\theta}^{(i)}) = \mathcal{L}(\boldsymbol{\theta}^{(i)}, \boldsymbol{\theta}^{(i)}).$$



- (Maximization step) Update $\boldsymbol{\theta}$:

$$\boldsymbol{\theta}^{(i+1)} = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i)}).$$

Concavity of log

Given a vector $\mathbf{p} = [p_1, p_2, \dots, p_K]^T$
 s.t. $0 < p_k < 1, \forall k, \sum_k p_k = 1$, the
 following holds for any $t_k > 0$:

$$\log \left(\sum_{k=1}^K p_k t_k \right) \geq \sum_{k=1}^K p_k \log t_k$$

Handwritten diagram illustrating the concavity of log:

$$\log \sum_{k=1}^K \pi_k \mathcal{N} \left(\frac{\mathbf{x}_n}{\cdot} \mid \underline{\mu}_k, \underline{\Sigma}_k \right)$$

The expression is enclosed in a green box. A red line underlines the sum, with p_{kn} written below it. A red arrow points from t_{kn} to the sum. Below the box, it says: s.t. $0 < p_{kn} < 1, \sum_{k=1}^K p_{kn} = 1$.

The expectation step

(A)
$$\log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n \mid \underline{\mu}_k, \underline{\Sigma}_k) \geq \sum_{k=1}^K p_{kn} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n \mid \underline{\mu}_k, \underline{\Sigma}_k)}{p_{kn}}$$

Annotations: $\mathcal{L}(\theta)$ points to the log of the sum. θ points to the parameters of the Gaussian. $\theta^{(i)}$ points to the parameters in the denominator. $p_{kn}^{(i)}$ points to the denominator.

with equality when,

(B)
$$p_{kn}^{(i)} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n \mid \underline{\mu}_k, \underline{\Sigma}_k)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n \mid \underline{\mu}_k, \underline{\Sigma}_k)}$$

This is not a coincidence.

$$\underline{\theta}^{(i)} = \{ \underline{\pi}^{(i)}, \underline{\mu}^{(i)}, \underline{\Sigma}^{(i)} \}$$

Annotations: $\underline{\pi}^{(i)}$ is labeled "Prior", $\underline{\mu}^{(i)}$ and $\underline{\Sigma}^{(i)}$ are labeled "Likelihood".

To show equality at $\theta^{(i)}$

$$\mathcal{L}(\theta^{(i)}) = \underline{\mathcal{L}}(\theta^{(i)}, \theta^{(i)})$$

We substitute $p_{kn}^{(i)}$ from (B) in (A)
 to show equality $\mathcal{L}(\theta^{(i)}) = \underline{\mathcal{L}}(\theta^{(i)}, \theta^{(i)})$

Let
$$\tilde{p}_{kn}^{(i)} = \pi_k \mathcal{N} \left(\frac{\mathbf{x}_n}{\cdot} \mid \underline{\mu}_k, \underline{\Sigma}_k \right)$$

and
$$p_{kn}^{(i)} = \frac{\tilde{p}_{kn}^{(i)}}{\sum_j \tilde{p}_{jn}^{(i)}}$$

$p_{kn}^{(i)}$ is the posterior (defined later)

$$\begin{aligned} \sum_k p_{kn}^{(i)} \log \frac{\pi_k \mathcal{N}(\cdot)}{p_{kn}^{(i)}} &= \sum_k \frac{\pi_k \mathcal{N}(\cdot)}{\sum_j \pi_j \mathcal{N}(\cdot)} \log \sum_j \pi_j \mathcal{N}(\cdot) \\ &= \frac{\sum_k \pi_k \mathcal{N}(\cdot)}{\sum_j \pi_j \mathcal{N}(\cdot)} \log \sum_j \pi_j \mathcal{N}(\cdot) \\ &= \log \sum_j \pi_j \mathcal{N}(\cdot) \end{aligned}$$

The lower bound is

The maximization step

$$= \sum_n \sum_k p_{kn}^{(i)} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n)}{p_{kn}^{(i)}}$$

Maximize the lower bound w.r.t. θ .

~~$p_{kn}^{(i)}$~~ is constant so we can ignore it.

$$\max_{\theta} \sum_{n=1}^N \sum_{k=1}^K p_{kn}^{(i)} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

$\left[\log \pi_k \quad -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad -\frac{1}{2} \log |\boldsymbol{\Sigma}_k| \right]$

Differentiating w.r.t. $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k^{-1}$, we can get the updates for $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$.

$$\boldsymbol{\mu}_k^{(i+1)} = \frac{\sum_{n=1}^N p_{kn}^{(i)} \mathbf{x}_n}{\sum_n p_{kn}^{(i)}} = \frac{\sum_n r_{kn} \mathbf{x}_n}{\sum_n r_{kn}}$$

$$\boldsymbol{\Sigma}_k^{(i+1)} = \frac{\sum_{n=1}^N p_{kn}^{(i)} (\mathbf{x}_n - \boldsymbol{\mu}_k^{(i+1)}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{(i+1)})^T}{\sum_n p_{kn}^{(i)}}$$

For π_k , we use the fact that they sum to 1. Therefore, we add a Lagrangian term, differentiate w.r.t. π_k and set to 0, to get the following update:

Homework

$$\pi_k^{(i+1)} = \frac{1}{N} \sum_{n=1}^N p_{kn}^{(i)}$$

Summary of EM for GMM

Initialize $\underline{\boldsymbol{\mu}}^{(1)}, \underline{\boldsymbol{\Sigma}}^{(1)}, \underline{\boldsymbol{\pi}}^{(1)}$ and iterate between the E and M step, until $\mathcal{L}(\boldsymbol{\theta})$ stabilizes.

1. E-step: Compute responsibilities $p_{kn}^{(i)}$:

$$p_{kn}^{(i)} = \frac{\pi_k^{(i)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{(i)}, \boldsymbol{\Sigma}_k^{(i)})}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{(i)}, \boldsymbol{\Sigma}_k^{(i)})} \quad \leftarrow$$

Prior x likelihood

2. Compute the marginal likelihood (cost).

$$\mathcal{L}(\boldsymbol{\theta}^{(i)}) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k^{(i)} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k^{(i)}, \boldsymbol{\Sigma}_k^{(i)}) \quad \leftarrow$$

3. M-step: Update $\boldsymbol{\mu}_k^{(i+1)}, \boldsymbol{\Sigma}_k^{(i+1)}, \pi_k^{(i+1)}$.

$$\boldsymbol{\mu}_k^{(i+1)} = \frac{\sum_{n=1}^N p_{kn}^{(i)} \mathbf{x}_n}{p_{kn}^{(i)}}$$
$$\boldsymbol{\Sigma}_k^{(i+1)} = \frac{\sum_{n=1}^N p_{kn}^{(i)} (\mathbf{x}_n - \boldsymbol{\mu}_k^{(i+1)}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{(i+1)})^T}{p_{kn}^{(i)}}$$
$$\pi_k^{(i+1)} = \frac{1}{N} \sum_{n=1}^N p_{kn}^{(i)}$$

If we let, covariance be diagonal i.e. $\boldsymbol{\Sigma}_k := \sigma^2 \mathbf{I}$, then EM algorithm is same as K-means as $\sigma^2 \rightarrow 0$.

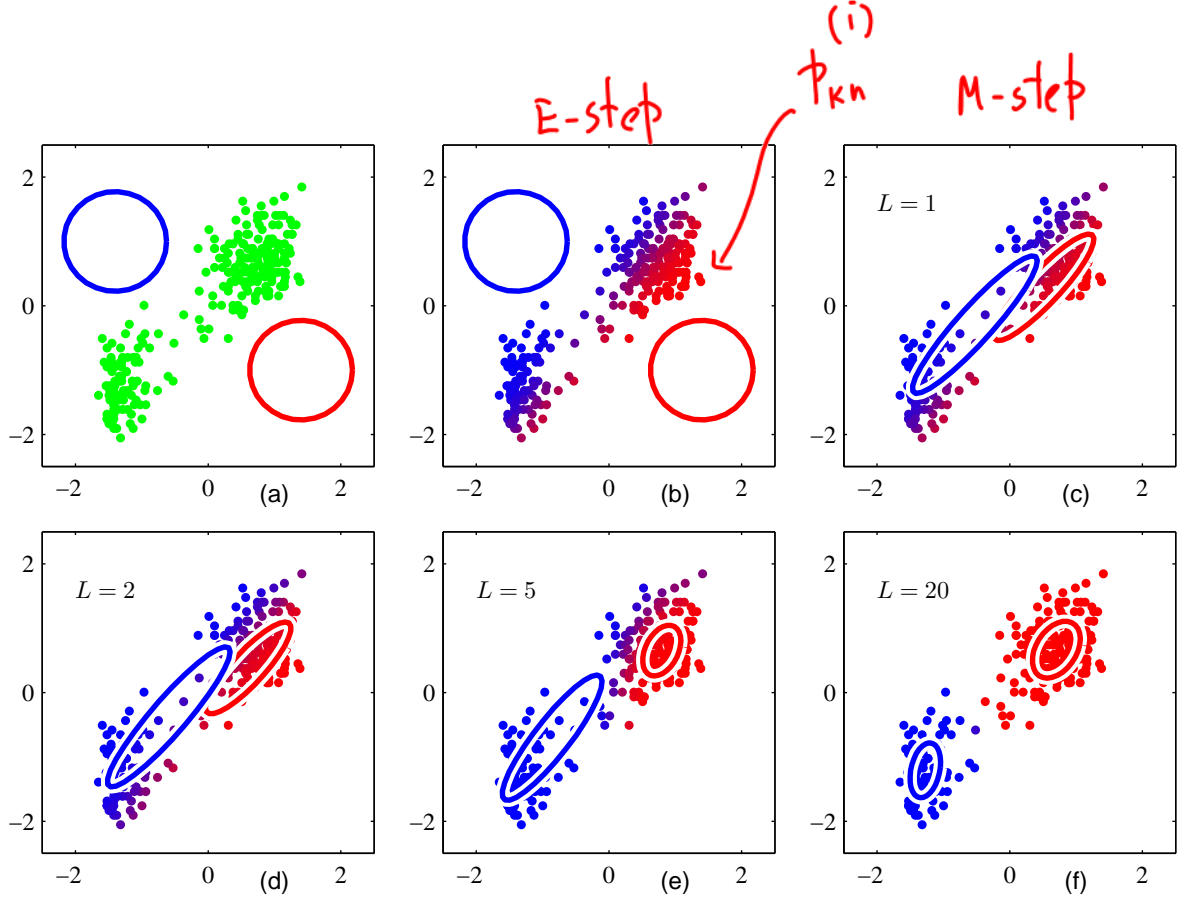


Figure 1: EM algorithm for GMM

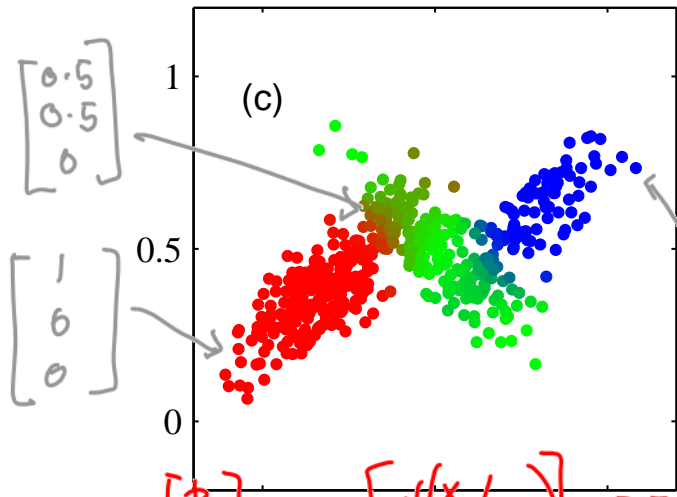
Posterior distribution

We now show that $p_{kn}^{(i)}$ is the posterior distribution of the latent variable, i.e. $p_{kn}^{(i)} = p(r_n = k | \mathbf{x}_n, \boldsymbol{\theta}^{(i)})$

$$p(\mathbf{x}_n, r_n | \boldsymbol{\theta}) = p(\mathbf{x}_n | r_n, \boldsymbol{\theta}) p(r_n | \boldsymbol{\theta}) = p(r_n | \mathbf{x}_n, \boldsymbol{\theta}) p(\mathbf{x}_n | \boldsymbol{\theta})$$

Joint Likelihood x Prior Posterior x Marginal likelihood
 $P(A, B) = P(A|B) P(B) = P(B|A) P(A)$ Bayes' rule

Bayes rule enables probabilistic inversion:
 Forward model: $A\underline{x} = \underline{b}$
 Inverse model: $\underline{x} = A^{-1}\underline{b}$



$$\underline{p} = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix} = \begin{bmatrix} \mathcal{N}(\underline{x}_n / \mu_1, \Sigma_1) \\ \mathcal{N}(\underline{x}_n / \mu_2, \Sigma_2) \\ \mathcal{N}(\underline{x}_n / \mu_3, \Sigma_3) \end{bmatrix} \cdot \begin{bmatrix} 0.5 \\ 0.4 \\ 0.1 \end{bmatrix}$$

$$p(r_n | \underline{x}_n, \boldsymbol{\theta}) = \frac{p(\underline{x}_n / r_n, \boldsymbol{\theta}) p(r_n / \boldsymbol{\theta})}{\sum_{j=1}^K p(\underline{x}_n, r_n = j / \boldsymbol{\theta})}$$

$$= \frac{p(\underline{x}_n / r_n, \boldsymbol{\theta}) p(r_n / \boldsymbol{\theta})}{\sum_{j=1}^K p(\underline{x}_n / r_n = j, \boldsymbol{\theta}) p(r_n = j / \boldsymbol{\theta})}$$

$$\underline{p} \propto \underline{p} \cdot \text{sum}(\underline{p})$$

EM in general $\log p(X|\theta)$ ^{Data} \leftarrow ^{Likelihood \times Prior} $= \log \sum p(X, r|\theta) \times p(r|\theta, X)$

Given a generic joint distribution $p(\mathbf{x}_n, r_n|\theta)$, marginal likelihood can be lower bounded in a similar way.

$\hookrightarrow p(r|\theta^{(i)}, X)$

EM algorithm can be compactly written as follows:

$\geq \sum [\log p(X, r|\theta)] p(r|\theta^{(i)}, X)$
 $+ H(p(r|\theta^{(i)}, X))$

$$\theta^{(i+1)} = \arg \max_{\theta} \sum_{n=1}^N \mathbb{E}_{p(r_n|\mathbf{x}_n, \theta^{(i)})} [\log p(\mathbf{x}_n, r_n|\theta)]$$

Another interpretation is that part of the data is missing, i.e. (\mathbf{x}_n, r_n) is the “complete” data and r_n is missing. EM algorithm averages over the “unobserved” part of the data.

To do

1. Identify the joint, likelihood, prior, and marginal distributions respectively. Understand the use of Bayes rule that relates all these distributions together.
2. Derive the posterior distribution for GMM.
3. Understand the relation between EM and K-means.
4. Relate the lower bound to EM for probabilistic models in general.
5. Read the Wikipedia page on how to find a good K.
6. Read Bishop Section 14.5 to learn about conditional mixture models and mixture of experts.
7. Read about other mixture models in KPM book.