

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.

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

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:

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

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

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

2. (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 $t_k > 0$:

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

The expectation step

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

with equality when,

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

This is not a coincidence.

The maximization step

Maximize the lower bound w.r.t. θ .

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

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

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

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:

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

Summary of EM for GMM

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

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

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

2. Compute the marginal likelihood (cost).

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

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 $\sigma^2 \rightarrow 0$.

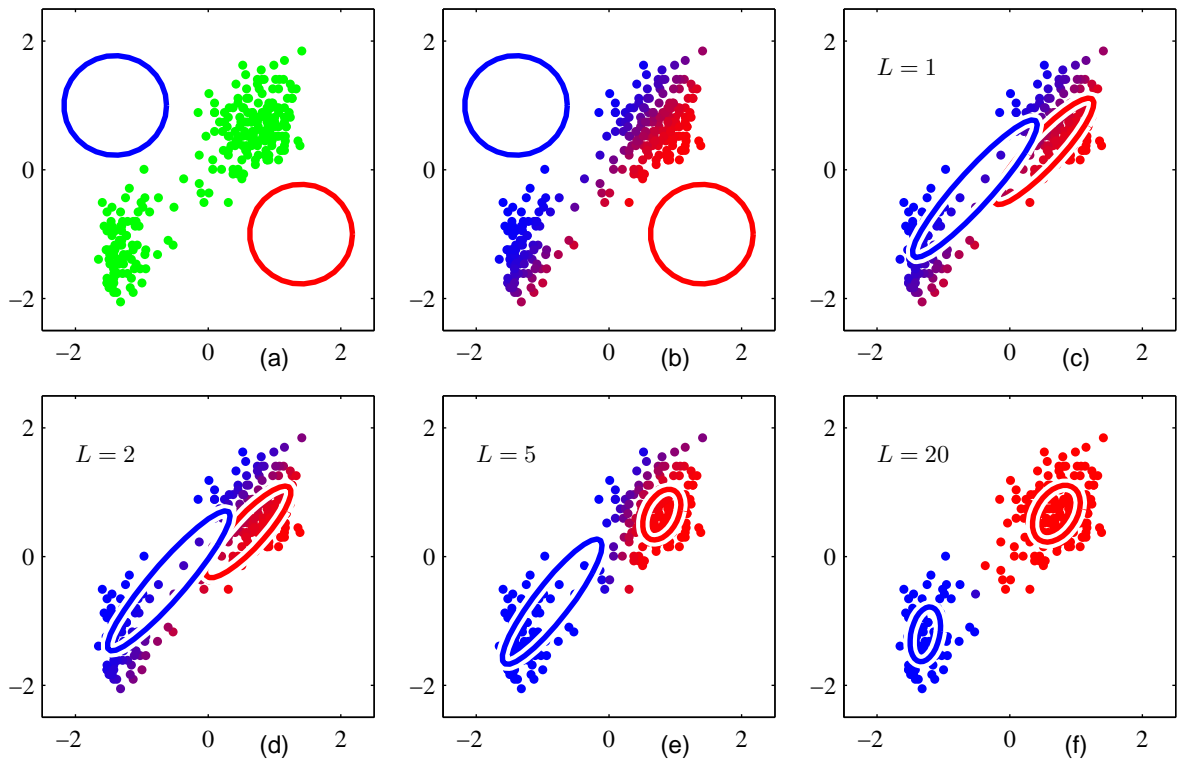
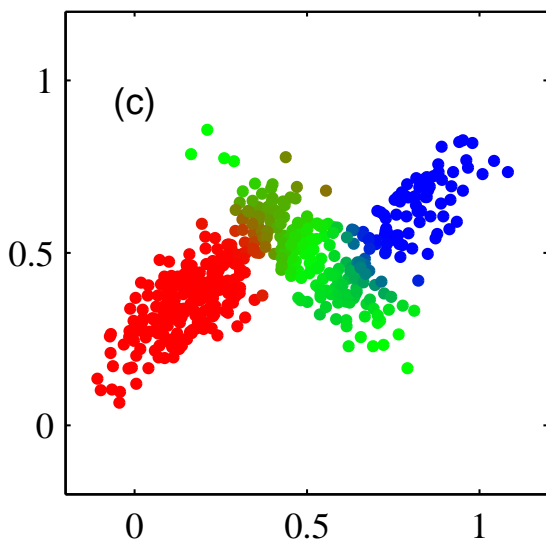


Figure 1: EM algorithm for GMM

Posterior distribution

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

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



EM in general

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

EM algorithm can be compactly written as follows:

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

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.