

annotated  
version

Machine Learning Course - CS-433

# Expectation- Maximization Algorithm

find  $\theta$

Nov 17, 2016

©Mohammad Emtiyaz Khan 2015

minor changes by Martin Jaggi 2016



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

# Motivation

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

$$\theta = (\pi_{1:K}, \mu_{1:K}, \Sigma_{1:K})$$

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

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

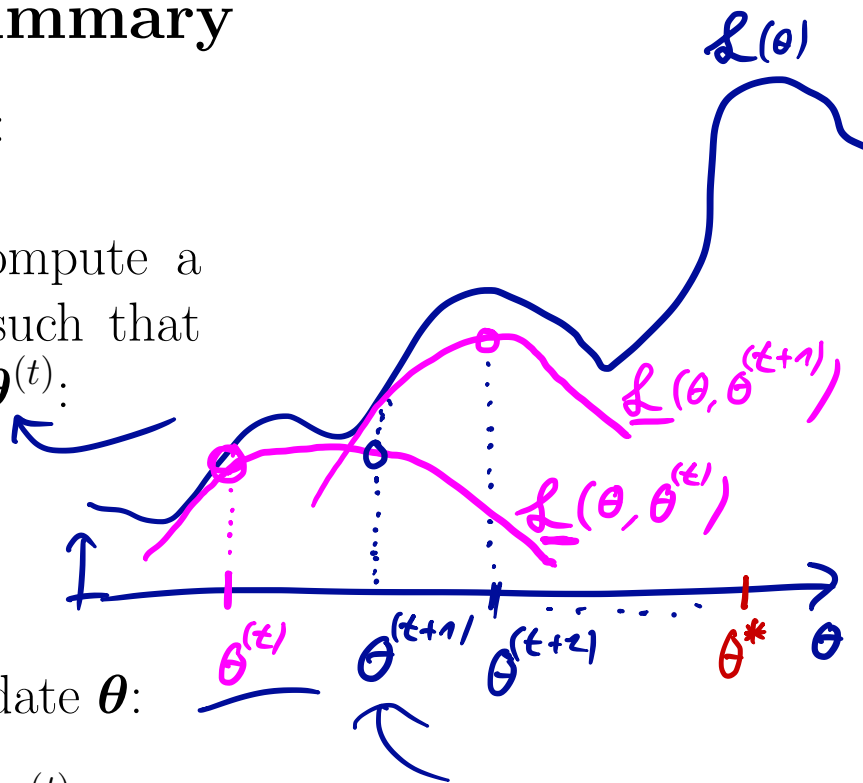
Start with  $\theta^{(1)}$  and iterate:

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

$$\mathcal{L}(\theta) \geq \underline{\mathcal{L}}(\theta, \theta^{(t)}) \text{ and } \mathcal{L}(\theta^{(t)}) = \underline{\mathcal{L}}(\theta^{(t)}, \theta^{(t)}).$$

2. **Maximization step:** Update  $\theta$ :

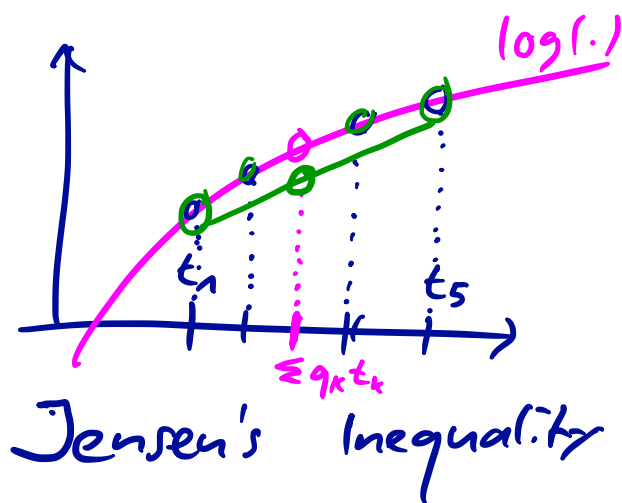
$$\theta^{(t+1)} = \arg \max_{\theta} \underline{\mathcal{L}}(\theta, \theta^{(t)}).$$



## Concavity of log

Given non-negative weights  $q$  s.t.  $\sum_k q_k = 1$ , the following holds for any  $t_k > 0$ :

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



## The expectation step

(per data example)

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

with equality when,

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

This is not a coincidence.

$$\mathcal{L}(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^{(t)}) \stackrel{?}{=} \mathcal{L}(\boldsymbol{\theta}^{(t)})$$

$$= \sum_k q_{kn} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}^{(t)})}{q_{kn}}$$

↓

$$= \sum_k \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}^{(t)})}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}^{(t)})} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}^{(t)})}{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}^{(t)})}$$

$$= \log \sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\theta}^{(t)})$$

$$= \mathcal{L}(\boldsymbol{\theta}^{(t)})$$

$$\rightarrow = \log \sum_k \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{kn}} q_{kn} = \log \sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \mathcal{L}(\boldsymbol{\theta})$$

for fixed  $q_{kn}^{(t)}$

## The maximization step

Maximize the lower bound w.r.t.  $\theta$ .

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

$$\max_{\theta} \sum_{n=1}^N \sum_{k=1}^K q_{kn}^{(t)} [\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)] - \lambda$$

$$q \log \frac{\pi \mathcal{N}(\cdot)}{q}$$

Differentiating w.r.t.  $\mu_k, \Sigma_k^{-1}$ , we can get the updates for  $\mu_k$  and  $\Sigma_k$ .

$$\mu_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} \mathbf{x}_n}{\sum_n q_{kn}^{(t)}}$$

$$\nabla_{\mu} \underline{\mathcal{L}}(\theta, \theta^{(t)}) \stackrel{!}{=} 0$$

$$\nabla_{\Sigma} \underline{\mathcal{L}}(\theta, \theta^{(t)}) \stackrel{!}{=} 0$$

$$\Sigma_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} (\mathbf{x}_n - \mu_k^{(t+1)}) (\mathbf{x}_n - \mu_k^{(t+1)})^{\top}}{\sum_n q_{kn}^{(t)}}$$

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

$$\dots + \beta (\sum_k \pi_k - 1)$$

$$\hookrightarrow \nabla_{\pi} \stackrel{!}{=} 0$$

$$\hookrightarrow \pi$$

$$\pi_k^{(t+1)} := \frac{1}{N} \sum_{n=1}^N q_{kn}^{(t)}$$

# 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 assignments  $q_{kn}^{(t)}$ :

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

for  $\boldsymbol{\Sigma}_k = \sigma^2 \mathbf{I}_D$   
 if  $\sigma \rightarrow 0$   
 $\rightarrow$  get k-means

2. Compute the marginal likelihood (cost).

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

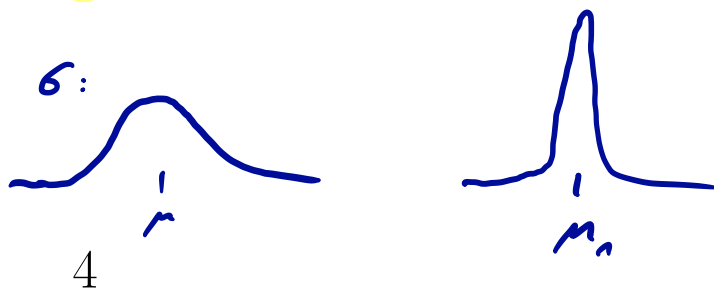
3. **M-step:** Update  $(\boldsymbol{\mu}_k^{(t+1)}, \boldsymbol{\Sigma}_k^{(t+1)}, \pi_k^{(t+1)}) = \boldsymbol{\theta}^{(t+1)}$

$$\boldsymbol{\mu}_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} \mathbf{x}_n}{\sum_n q_{kn}^{(t)}}$$

$$\boldsymbol{\Sigma}_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t+1)}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{(t+1)})^\top}{\sum_n q_{kn}^{(t)}}$$

$$\pi_k^{(t+1)} := \frac{1}{N} \sum_n q_{kn}^{(t)}$$

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



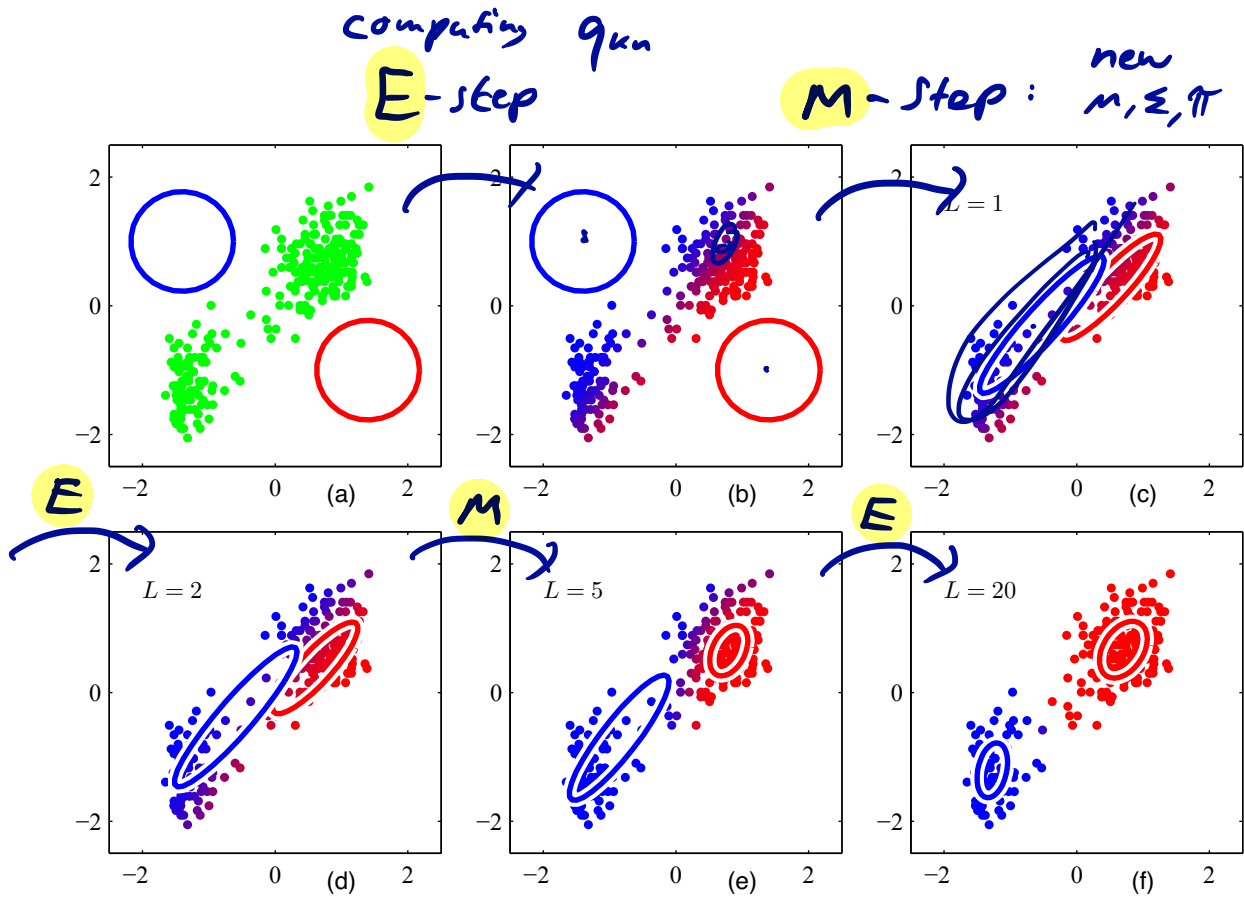
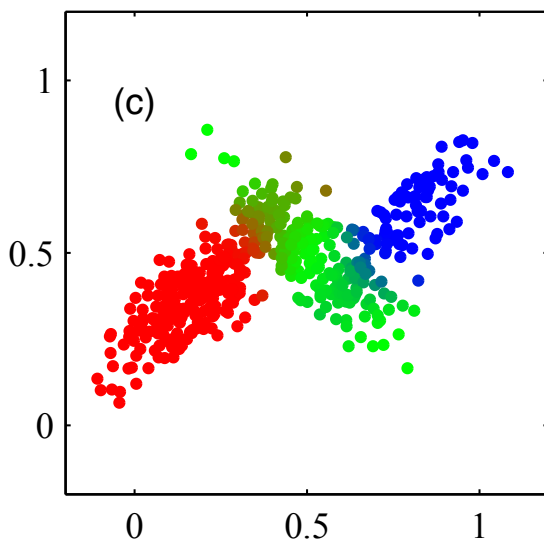


Figure 1: EM algorithm for GMM

## Posterior distribution

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

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



# EM in general

Given a general joint distribution  $p(\mathbf{x}_n, z_n | \boldsymbol{\theta})$ , the marginal likelihood can be lower bounded similarly:

The EM algorithm can be compactly written as follows:

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

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

## ToDo

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 about other mixture models in the KPM book.