

*Annotated  
version*

Machine Learning Course - CS-433

# Expectation-Maximization Algorithm

Nov 29, 2023

Martin Jaggi

Last updated on: November 28, 2023

credits to Mohammad Emtiyaz Khan & Rüdiger Urbanke



# Motivation

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

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

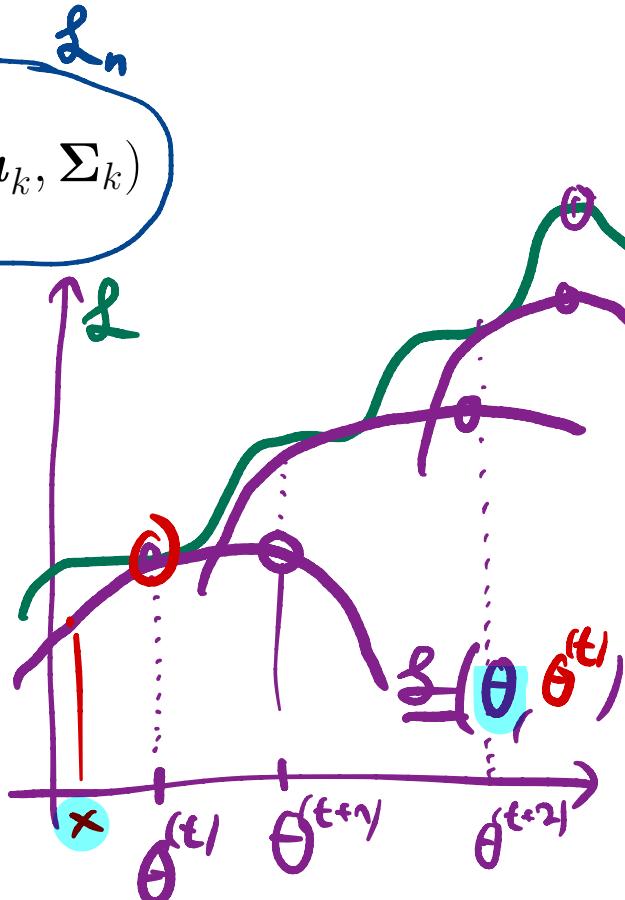
Expectation-Maximization (EM) algorithm provides an elegant and general method to optimize such optimization problems. It uses an iterative two-step procedure where individual steps usually involve problems that are easy to optimize.

## EM algorithm: Summary

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

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

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



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

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

2. Maximization step: Update  $\theta$ :

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

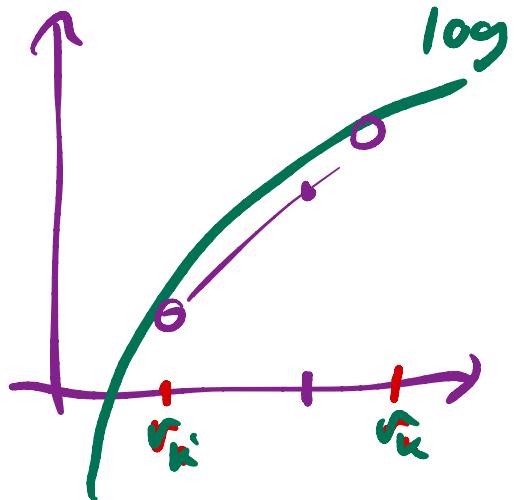
$-\log$  is convex

Concavity of  $\log$   
Jensen's inequality

Given non-negative weights  $q$  s.t.

$\sum_k q_k = 1$ , the following holds for any  $r_k > 0$ :

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



## The expectation step

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

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

$$\geq \sum_{k=1}^K q_{kn} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{r_k}$$

with equality when,

$$q_{kn} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}$$

- $\sum_k q_{kn} = 1$  ✓
- $\underline{\mathcal{L}}_n(\theta^{(t)}, \theta^{(t)}) = \frac{1}{2} \underline{\mathcal{L}}_n(\theta^{(t)})$

This is not a coincidence.

$$\begin{aligned} & \sum_k q_{kn}^{(t)} \log(r_k) \\ &= \sum_k \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)} \\ &= \log \sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) = \underline{\mathcal{L}}_n(\theta^{(t)}) \end{aligned}$$

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \dots e^{-(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})}$$

## The maximization step

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

$$\underline{\mathcal{L}}_n(\theta, \theta^{(t)}) \quad \text{indep of } \theta$$

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

$$= \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \dots)}{q_{kn}}$$

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

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

$$\nabla_{\boldsymbol{\mu}} \underline{\mathcal{L}}_n(\theta, \theta^{(t)}) \stackrel{!}{=} 0 \quad \text{solve for } \boldsymbol{\mu}$$

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

$$\nabla_{\boldsymbol{\Sigma}} \underline{\mathcal{L}}_n(\theta, \theta^{(t)}) \stackrel{!}{=} 0 \quad \text{solve for } \boldsymbol{\Sigma}^{-1}$$

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:

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

~~$$\nabla_{\pi} \underline{\mathcal{L}}_n \stackrel{!}{=} 0 \quad ??$$~~

must respect  $\sum_k \pi_k = 1$

$$\nabla_{\pi} \tilde{\underline{\mathcal{L}}}_n \stackrel{!}{=} 0$$

$$\tilde{\underline{\mathcal{L}}} := \underline{\mathcal{L}} + \beta (\sum_k \pi_k - 1)$$

**enforce  $\Sigma = \sigma^2 \mathbf{I}$**

## Summary of EM for GMM

k-means as  
a special case

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

1. E-step: Compute assignments  $q_{kn}^{(t)}$ .

$$q_{kn}^{(t)} := \frac{\pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{k=1}^K \pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t)}, \Sigma_k^{(t)})}$$

2. Compute the marginal likelihood (cost).

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

= k-means  
assignment

3. M-step: Update  $\mu_k^{(t+1)}, \Sigma_k^{(t+1)}, \pi_k^{(t+1)}$ .

$$\theta^{(t+1)} = \begin{cases} \mu_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} \mathbf{x}_n}{\sum_n q_{kn}^{(t)}} \\ \Sigma_k^{(t+1)} := \frac{\sum_n q_{kn}^{(t)} (\mathbf{x}_n - \mu_k^{(t+1)}) (\mathbf{x}_n - \mu_k^{(t+1)})^\top}{\sum_n q_{kn}^{(t)}} \\ \pi_k^{(t+1)} := \frac{1}{N} \sum_n q_{kn}^{(t)} \end{cases}$$

$q_{kn} = z_{kn}$   
 $= k\text{-means update}$

$\cancel{\Sigma = \sigma^2 \mathbf{I}}$

$\cancel{\# points assigned to cluster k}$

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

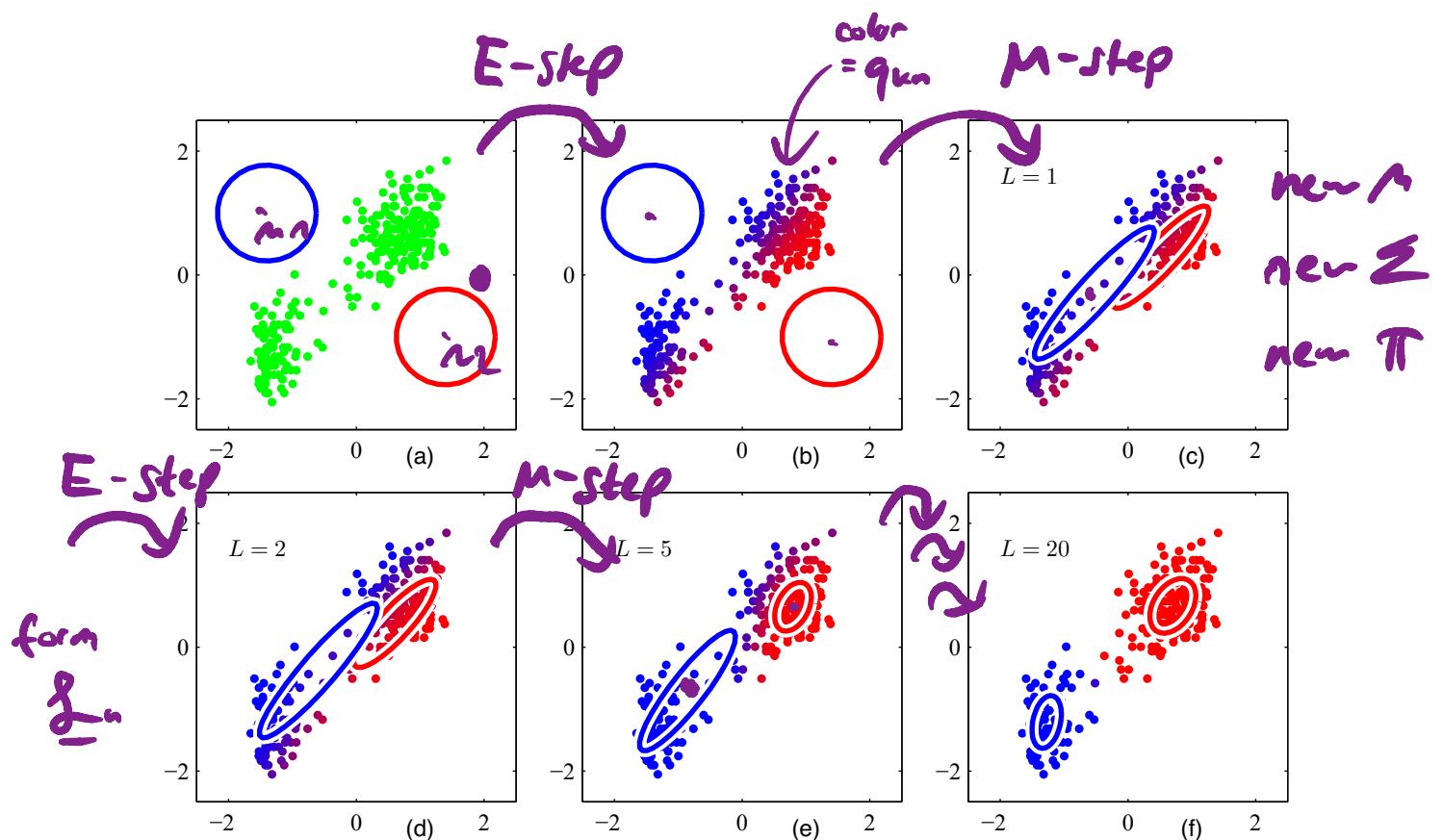


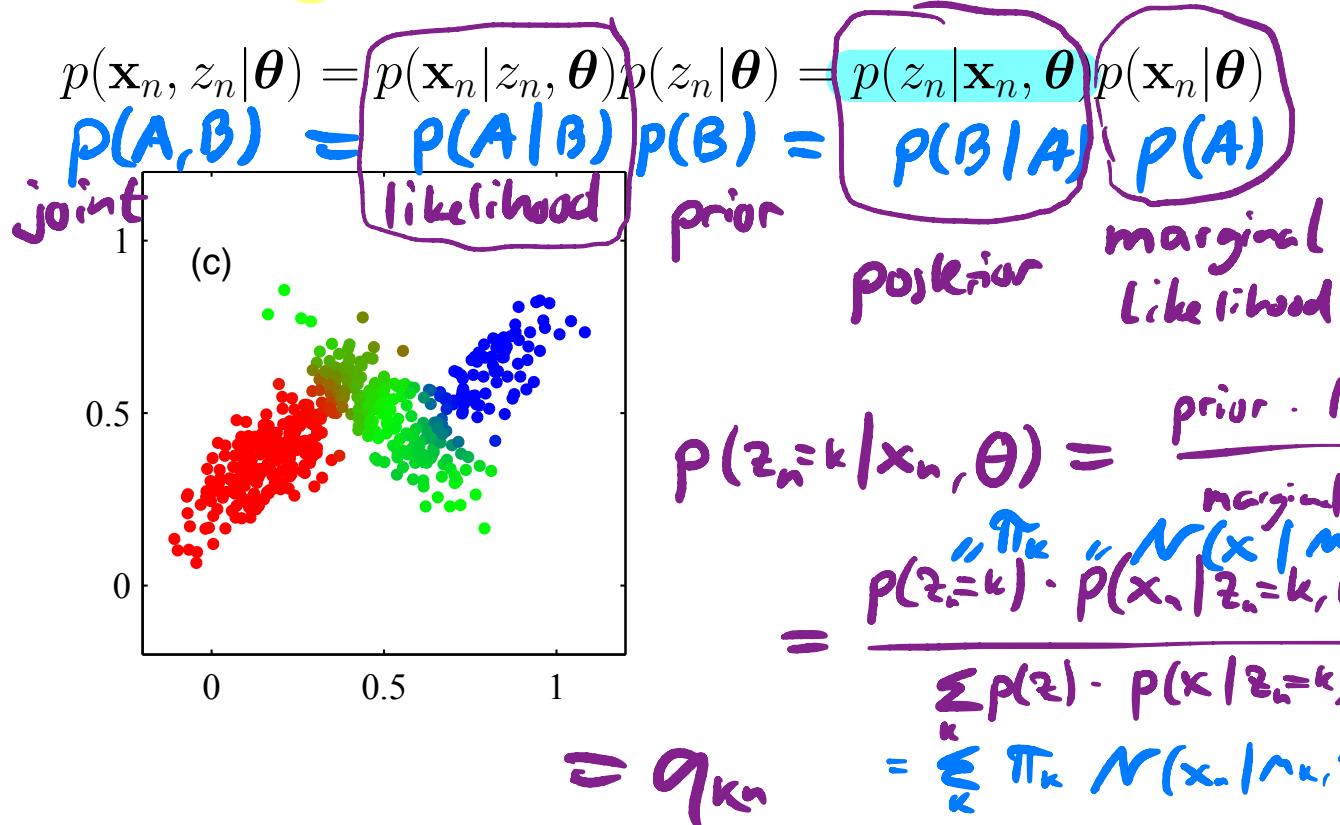
Figure 1: EM algorithm for GMM

## Posterior distribution

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

$$\sum_{k=1}^K p(\mathbf{x}_n | z_n) p(z_n)$$

!!



## EM in general

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

*Maximization step*

The EM algorithm can be compactly written as follows:

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

*$\mathcal{L}_n$*

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.

*expectation over  $z$   
= taking marginal  
= expectation step*