

*Annotated
Version*

Machine Learning Course - CS-433

Expectation-Maximization Algorithm

Dec 2, 2021

changes by Martin Jaggi 2021, 2020, 2019, changes by Rüdiger Urbanke 2018, changes by
Martin Jaggi 2017, 2016 ©Mohammad Emtiyaz Khan 2015

Last updated on: November 30, 2021

EPFL

Motivation

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

$$\theta = (\mu, \Sigma, \pi)$$

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

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

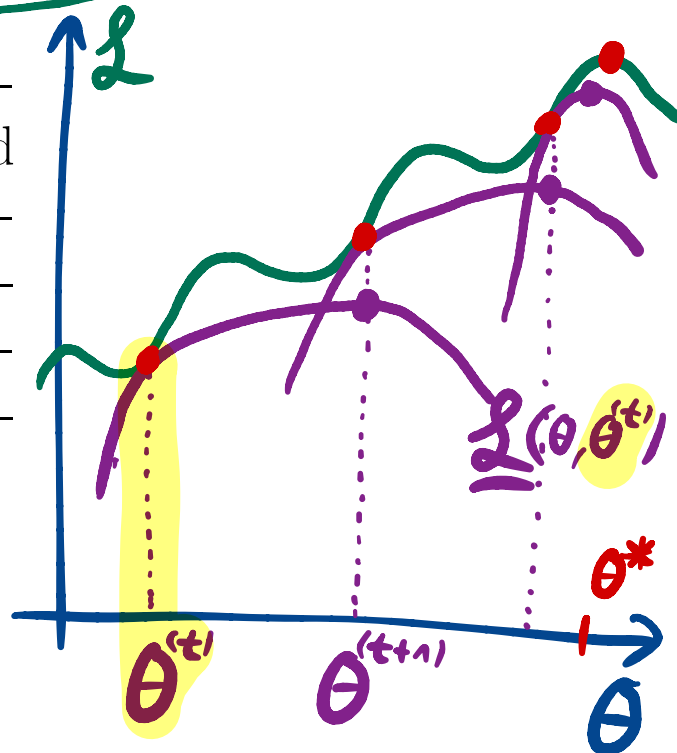
Start with $\theta^{(1)}$ and iterate:

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 θ :

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

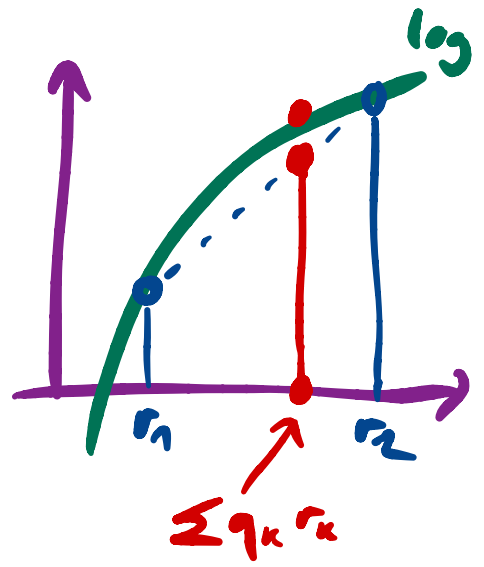


\Rightarrow Convexity of $-\log$

Concavity of \log

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

\Rightarrow Jensen's Inequality



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

The expectation step

lower bound on \mathcal{L}

$$\log \sum_{k=1}^K \underbrace{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}_{q_k \cdot r_k} \geq \sum_{k=1}^K q_{kn} \log \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{q_{kn}}}_{=: r_k} =: \mathcal{L}_n(\theta, \theta^{(t)})$$

with equality when,

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

This is not a coincidence.

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

$$= \sum_{k=1}^K q_{kn}^{(t)} \log \square$$

$$= \sum_{k=1}^K \frac{\pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{k'=1}^K \pi_{k'}^{(t)} \mathcal{N}(\mathbf{x}_n | \mu_{k'}^{(t)}, \Sigma_{k'}^{(t)})} \log \frac{\pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t)}, \Sigma_k^{(t)})}{\pi_k^{(t)} \mathcal{N}(\mathbf{x}_n | \mu_k^{(t)}, \Sigma_k^{(t)})}$$

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

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

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

The maximization step

Maximize the lower bound w.r.t. θ .

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

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

constant in θ

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

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

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

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

$$\nabla_{\boldsymbol{\Sigma}_k^{-1}} \underline{\mathcal{L}}(\theta, \theta^{(t)}) \stackrel{!}{=} 0$$

$\mathbf{v} = \mathbf{x} - \boldsymbol{\mu}$

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:

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

!!

$$\underline{\mathcal{L}}_n(\theta, \theta^{(t)}) + \beta (\sum \pi_k - 1)$$

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

Summary of EM for GMM

and k-mean as 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 $q_{kn}^{(t)}$:

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

$$\approx \exp \left(-\frac{\|\mathbf{x}_n - \mu_k\|^2}{2\sigma^2} \right)$$

2. Compute the marginal likelihood (cost).

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

$\begin{cases} 1 & \text{k closest to } \mathbf{x}_n \\ 0 & \text{otherwise} \end{cases}$
= k-means assignment

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

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

← mean of cluster

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

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

← # 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 $\sigma^2 \rightarrow 0$.

σ width



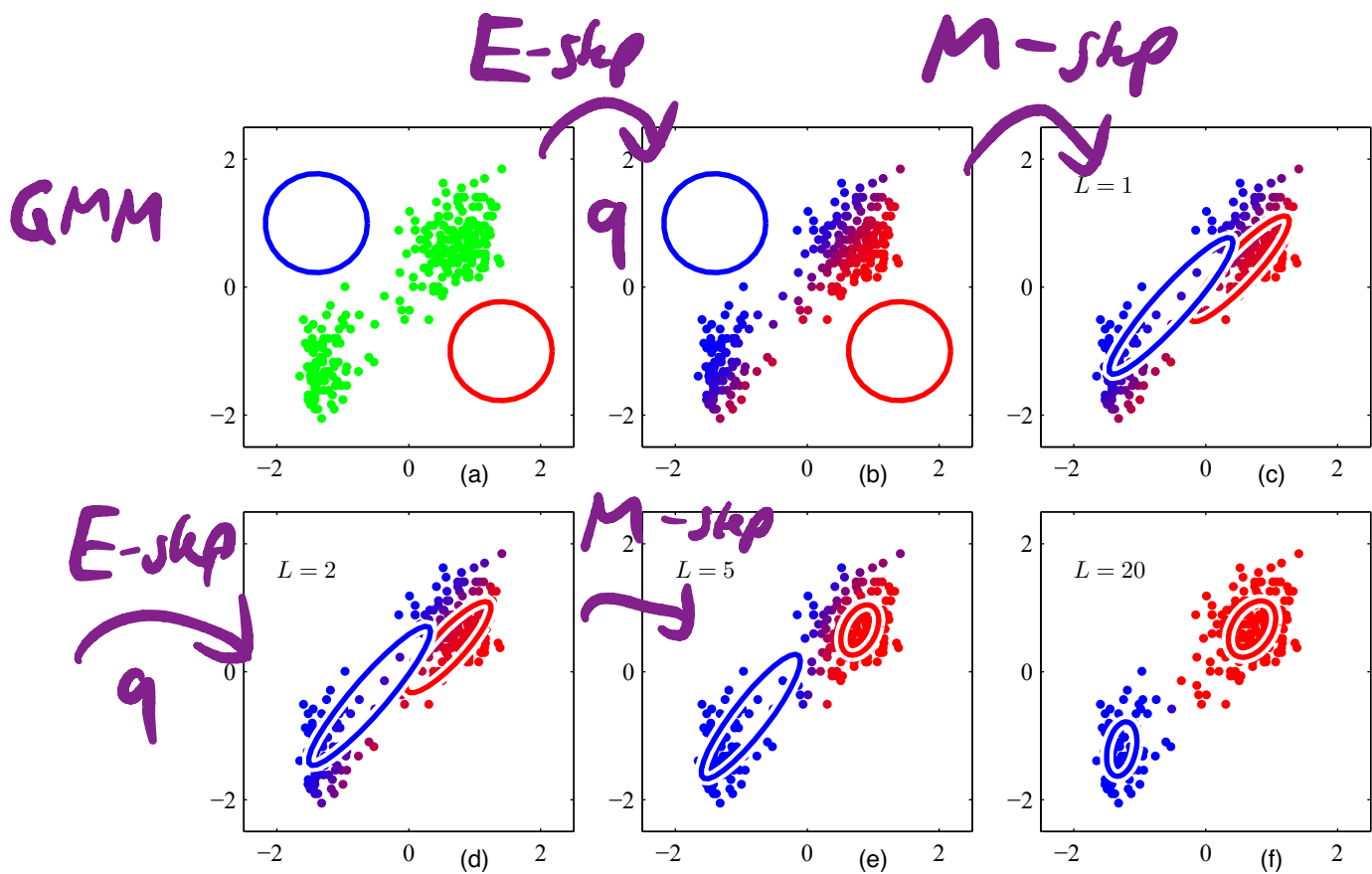


Figure 1: EM algorithm for GMM

Posterior distribution

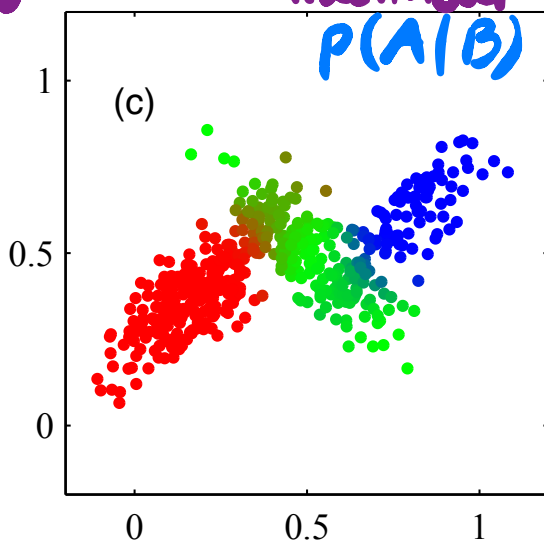
We now show that $q_{kn}^{(t)}$ is the posterior distribution of the latent variable, i.e. $q_{kn}^{(t)} = p(z_n = k | \mathbf{x}_n, \boldsymbol{\theta}^{(t)})$

$$p(\mathbf{x}_n, z_n | \boldsymbol{\theta}) = \underbrace{p(\mathbf{x}_n | z_n, \boldsymbol{\theta})}_{\text{likelihood}} \underbrace{p(z_n | \boldsymbol{\theta})}_{\text{prior}} = \underbrace{p(z_n | \mathbf{x}_n, \boldsymbol{\theta})}_{\text{posterior}} \underbrace{p(\mathbf{x}_n | \boldsymbol{\theta})}_{\text{marginal likelihood}}$$

joint **likelihood** **prior** **posterior** **marginal likelihood**

$P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$

Bayes Rule



$$p(z_n = k | \mathbf{x}_n, \boldsymbol{\theta}) = \frac{\text{prior} \cdot \text{likelihood}}{\text{ML}}$$

$$= \frac{\underbrace{\pi_k}_{p(z_n = k | \boldsymbol{\theta})} \cdot \underbrace{\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{p(\mathbf{x}_n | z_n = k, \boldsymbol{\theta})}}{\sum_k \underbrace{\pi_k}_{p(z_n = k)} \cdot \underbrace{\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{p(\mathbf{x}_n | z_n = k, \boldsymbol{\theta})}}$$

$$= q_{nk}$$

EM in general

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

The EM algorithm can be compactly written as follows:

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

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.

Expectation
over z

Maximization