

Annotated
Version

Machine Learning Course - CS-433

Expectation-Maximization Algorithm

Nov 26, 2020

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

Last updated on: November 24, 2020

EPFL

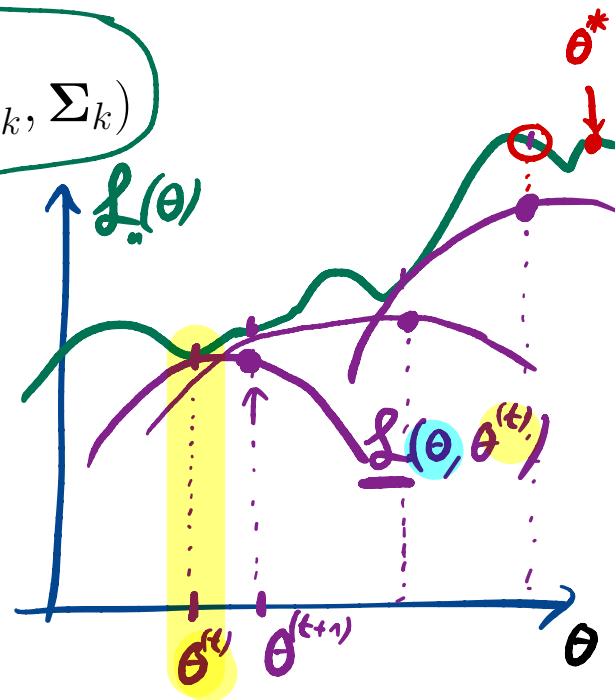
Motivation

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

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

Handwritten notes: $\theta = (\mu, \Sigma, \pi)$ (in purple); $\mathcal{L}_n(\theta)$ (in green, above the sum); $\mathcal{L}(\theta)$ (in green, above the log)

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 $\theta^{(t)}$:

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

$\forall \theta$

- lower bound
- equality if $\theta = \theta^{(t)}$

2. **Maximization step:** Update θ :

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

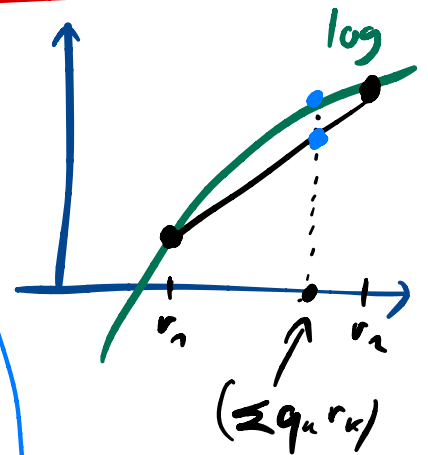
⇒ Convexity of $-\log$

Concavity of log

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

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

⇒ Jensen's Inequality



The expectation step

lower bound to $\mathcal{L}_n(\theta)$

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

with equality when,

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

This is not a coincidence.

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

• lower bound ✓

• coincides at $\theta = \theta^{(t)}$ ✓

$$\begin{aligned} &= \sum_{k=1}^K \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{k'=1}^K \pi_{k'} \mathcal{N}(\mathbf{x}_n | \mu_{k'}, \Sigma_{k'})}}_{=1, q_{kn}} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)} \\ &= \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \\ &= \mathcal{L}_n(\theta^{(t)}) \end{aligned}$$

The maximization step

Maximize the lower bound w.r.t. θ .

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

Handwritten notes: $\mathcal{L}_n(\theta, \theta^{(t)})$, $\log \frac{\pi_k \mathcal{N}(\mathbf{x}_n, \dots)}{q_{kn}^{(t)}}$, $e^{-(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)}$, independent of θ

Differentiating w.r.t. μ_k, Σ_k^{-1} , we can get the updates for μ_k and Σ_k .

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

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

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

$$\nabla_{\Sigma_k^{-1}} \mathcal{L}(\theta, \theta^{(t)}) \stackrel{!}{=} 0$$

$$\frac{1}{\sqrt{T}} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} = \text{rank one matrix}$$

$V = \mathbf{x}_n - \mu_n$

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:

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

want: $\sum_k \pi_k = 1$ (constraint)

maximize $\mathcal{L}_n + \beta (\sum_k \pi_k - 1)$

(unconstrained)

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

Summary of EM for GMM

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

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

Handwritten notes: $\approx \frac{\exp(-\|\mathbf{x}_n - \mu_k\|^2 / \sigma^2)}{\sum_{k=1}^K \exp(-\|\mathbf{x}_n - \mu_k\|^2 / \sigma^2)}$

2. Compute the marginal likelihood (cost).

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

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

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

$$\pi_k^{(t+1)} := \frac{1}{N} \sum_n q_{kn}^{(t)} \quad \leftarrow \text{\# points assigned to } k$$

Handwritten notes: $\sigma^2 \rightarrow 0$, $\{ \dots \}$, \rightarrow clock \leftarrow , $= k\text{-mean assignment}$, $q_{kn} \approx z_{kn}$

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



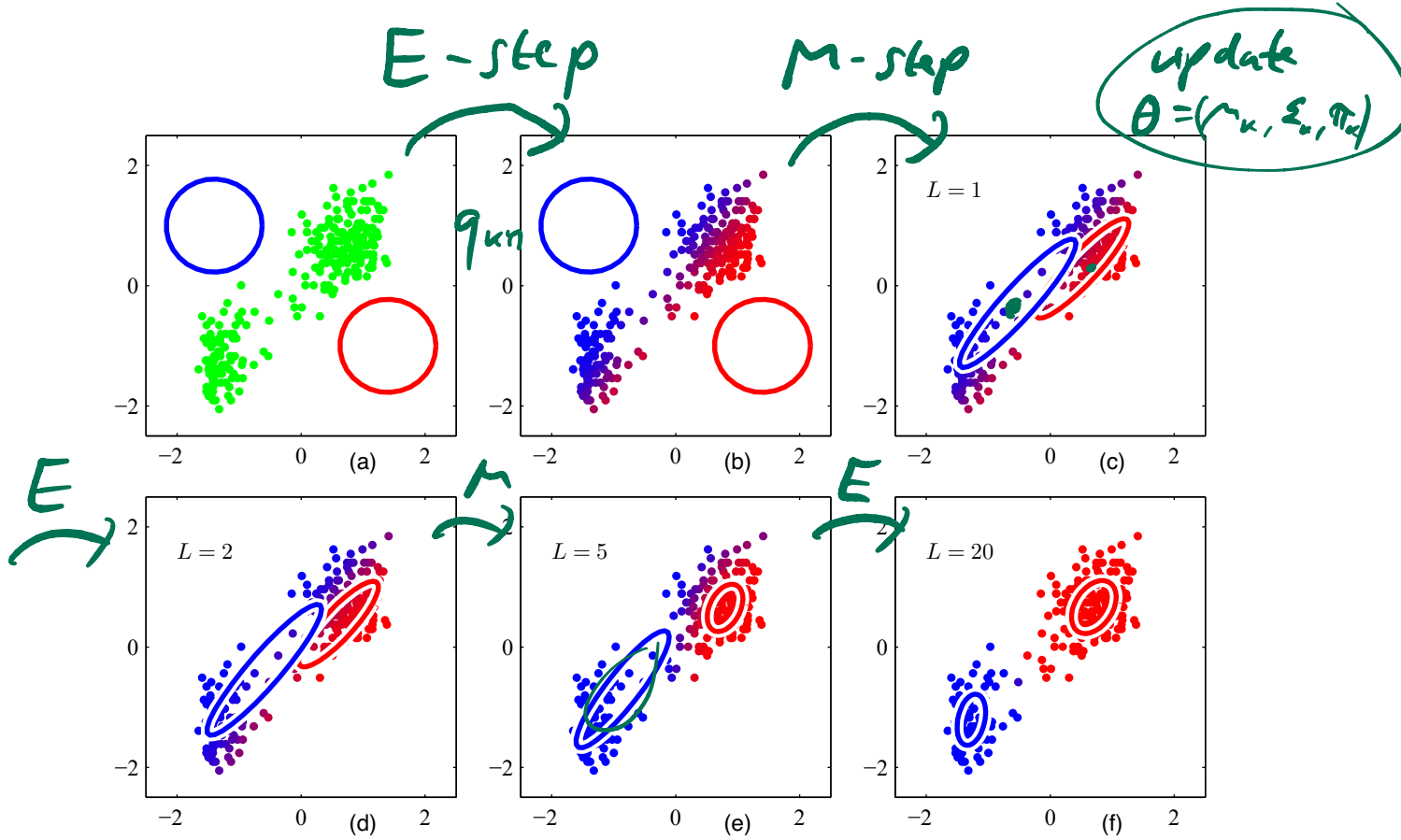


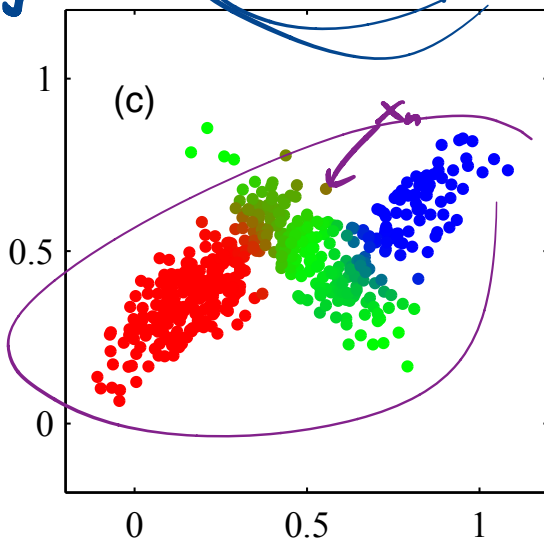
Figure 1: EM algorithm for GMM

Posterior distribution

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

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

joint = likelihood · prior = posterior · marginal likelihood



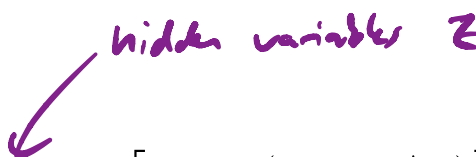
$$p(z_n = k | \mathbf{x}_n, \theta) = \frac{\text{prior} \cdot \text{likelihood}}{\sum_{k=1}^K \text{prior} \cdot \text{likelihood}}$$

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

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:

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


hidden variables z

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.