Annotated Version

**Machine Learning Course - CS-433**

# Expectation-Maximization Algorithm

Nov 7, 2019

EPFL

$$\boldsymbol{\theta} = \left( \{\pi_k\}_{k=1}^{K}, \{\boldsymbol{\mu}_k\}_{k=1}^{K}, \{\boldsymbol{\Sigma}_k\}_{k=1}^{K} \right)$$
$$\quad \uparrow \mathbb{R} \qquad\quad \uparrow \mathbb{R}^D \qquad\quad \uparrow \mathbb{R}^{D \times D}$$

# Motivation

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

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

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 $\boldsymbol{\theta}^{(t)}$:

- $\mathcal{L}(\boldsymbol{\theta}) \geq \underline{\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)})}$ and $\quad \forall \boldsymbol{\theta}$  — lower bound
- $\mathcal{L}(\boldsymbol{\theta}^{(t)}) = \underline{\mathcal{L}(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^{(t)})}$.  — equality at $\boldsymbol{\theta} = \boldsymbol{\theta}^{(t)}$
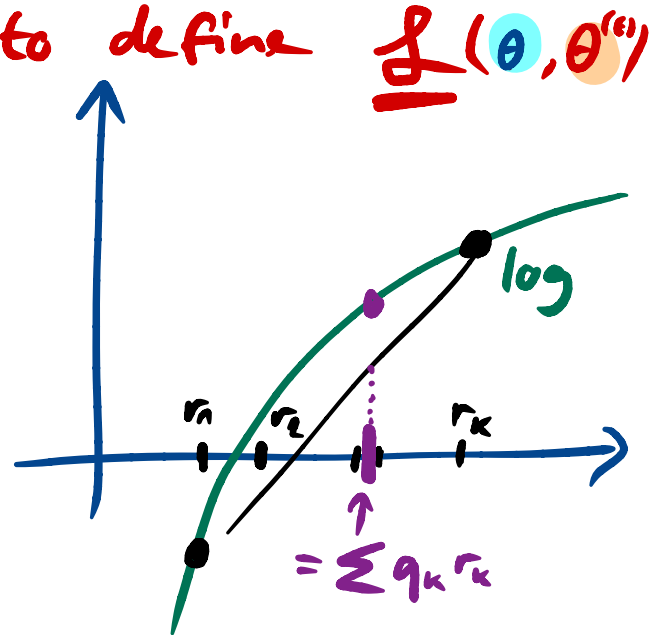
2. Maximization step: Update $\boldsymbol{\theta}$:

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

**How to define** $\underline{\mathcal{L}}(\theta, \theta^{(t)})$

# Concavity of log

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

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



$$= \sum q_k r_k$$

$q_k \cdot r_k$

Jensen's inequality
$\Leftrightarrow -\log$ is convex

# The expectation step

$\mathcal{L}_n(\theta)$

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

$r_k$

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

with equality when,

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

This is not a coincidence.

- lower bound ✓
- coincides with $\mathcal{L}$ at $\theta^{(t)}$

$\theta'' = \theta^{(t)}$

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

$$\sum_{k=1}^{k} \underbrace{\left(\frac{\pi_k \, \mathcal{N}(.)}{\sum_{k'} \pi_{k'} \mathcal{N}(.)}\right)}_{q_{kn}} \log \underbrace{\frac{\pi_k \, \mathcal{N}(.)}{\frac{\pi_k \, \mathcal{N}(.)}{\sum_{k'} \pi_{k'} \mathcal{N}(.)}}}_{q_{kn}}$$

$$= \log \sum_{k=1}^{k} \pi_k \mathcal{N}(.)$$

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

# The maximization step

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

$$\log\left(\frac{\pi_k \, \mathcal{N}(x_n | ..)}{q_{kn}}\right)$$

Maximize the lower bound w.r.t. $\boldsymbol{\theta}$.

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

$$\exp^{-(x_n - \mu)\Sigma^{-1}(x_n - \mu)}$$

independent of $\theta$

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

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

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

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

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

$$v = x_n - \mu_k$$

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:

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

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

want $\sum_k \pi_k = 1$

$$\mathcal{L}_n + \beta\left(\sum_{k=1}^{K} \pi_k - 1\right)$$

modified   unconstrained objective.

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

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

$$\approx \exp\left(-\frac{\|x_n - m_k\|^2}{\sigma^2}\right)$$
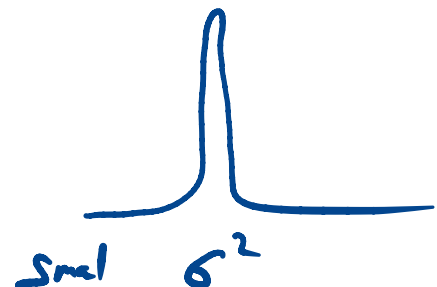
$$\approx \sum_{k=1}^{k} \exp(\cdots/\sigma^2)$$

$$\sigma^2 \to 0 \implies \begin{cases} 1 & \text{clost } k \\ 0 & \text{other} \end{cases}$$

k-means assignment

$q_{kn} \approx z_{kn}$

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 $\boldsymbol{\mu}_k^{(t+1)}, \boldsymbol{\Sigma}_k^{(t+1)}, \pi_k^{(t+1)}$.

mean

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

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

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

\#member of cluster k

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

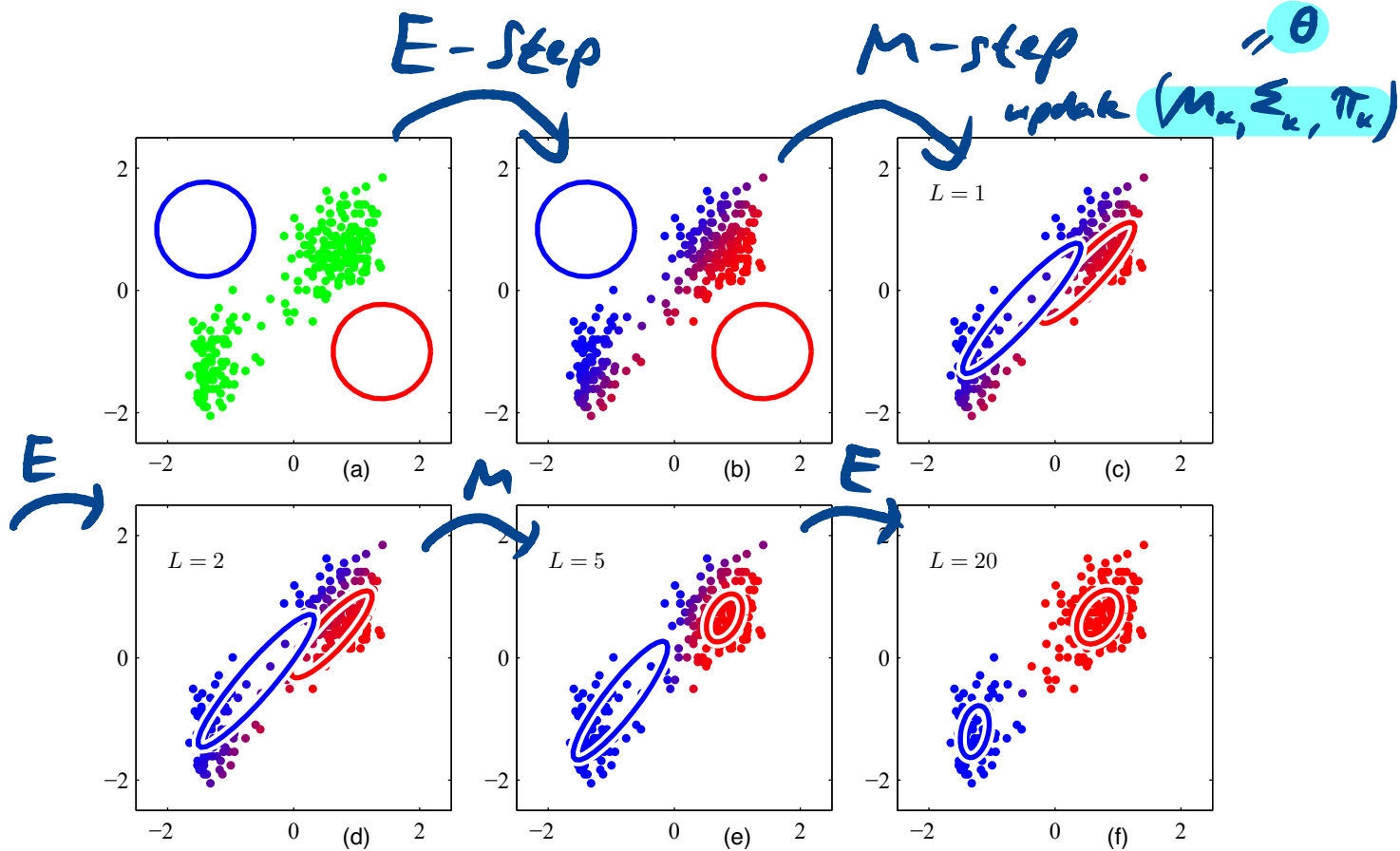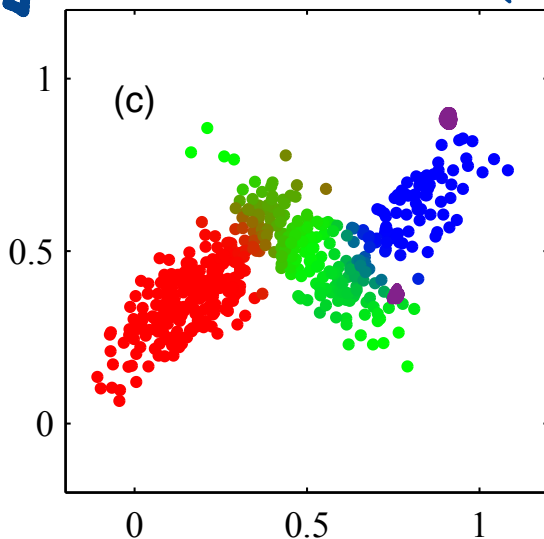large $\sigma^2$

smal $\sigma^2$

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. $q_{kn}^{(t)} = p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\theta}^{(t)})$

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

*(handwritten annotations:)*

$= \theta$

M-step update $(M_k, \Sigma_k, \pi_k)$

E-step   M-step

E   M   E

posterior   Bayes

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

joint   likelihood · prior   posterior · marginal likelihood



$$p(z_n = k \mid x_n, \theta) = \frac{prior \cdot likelihood}{ML}$$

$$= \frac{p(z_n = k \mid \theta)\, p(x_n \mid z_n = k, \theta)}{\sum_{k=1}^{K} p(z = k)\, p(x_n \mid z_n = k, \theta)}$$

$$= \frac{\pi_k\, \mathcal{N}(x_n \mid M_k, \Sigma_k)}{\sum_k \pi_k\, \mathcal{N}(x_n \mid \cdot \; \cdot)} =: q_{kn}$$

$\pi_k$   $p(x) =$

# 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:

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

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.