

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 Emamiyaz Khan 2015

Last updated on: November 24, 2020

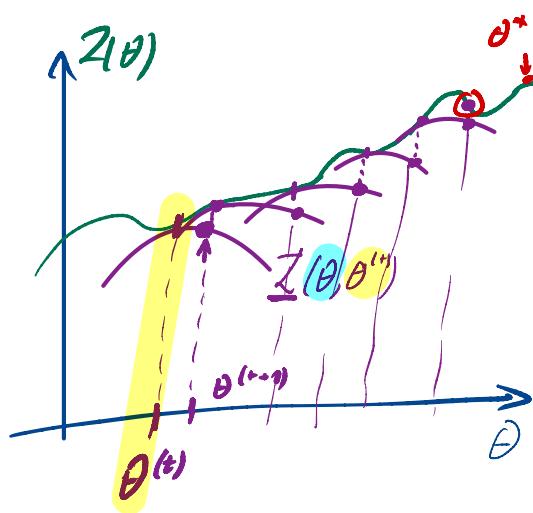


Motivation

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

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

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

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

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &\geq \underline{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(t)}) \text{ and } \\ \mathcal{L}(\boldsymbol{\theta}^{(t)}) &= \underline{\mathcal{L}}(\boldsymbol{\theta}^{(t)}, \boldsymbol{\theta}^{(t)}). \end{aligned}$$

$\forall \boldsymbol{\theta}$ *lower bound equality if $\boldsymbol{\theta} = \boldsymbol{\theta}^{(t)}$*

- ② **Maximization step:** Update $\boldsymbol{\theta}$:

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

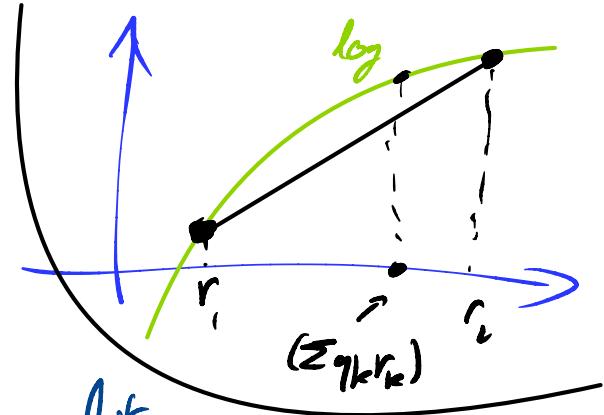
\Leftrightarrow convexity of $-\log$

Concavity of log

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

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

\Leftrightarrow Jensen's inequality



The expectation step

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

lower bound for Z_n

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

$= Z_n(\theta, \theta^{(t)})$

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.

$$Z_n(\theta, \theta^{(t)}) = Z_n(\theta^{(t)})$$

- lower bound ✓
- coincide at $\theta = \theta^{(t)}$ ✓

$$= \sum_{k=1}^K \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_n | \mu_{k'}, \Sigma_{k'})}$$

$q_{kn} = 1$

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

The maximization step

Maximize the lower bound w.r.t. θ .

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

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

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

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

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:

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

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

We want: $\sum_k \pi_k = 1$ constraint

$$\sum_k + \beta \left(\sum_k \pi_k - 1 \right)$$

(unconstrained)

$$F = \frac{1}{V^T} \sum_k \pi_k = \frac{1}{V^T} \sum_k q_{kn}^{(t)}$$

$$V = \mathbf{x}_n - \boldsymbol{\mu}_k$$

$$\mathcal{L}(\theta, \theta^{(t)}) = \sum_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

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

$$\nabla_{\boldsymbol{\Sigma}_k} \mathcal{L}(\theta, \theta^{(t)}) = 0$$

Summary of EM for GMM

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

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

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

$$(if \Sigma_k = \mathbf{I})$$

$$-||\mathbf{x}_n - \boldsymbol{\mu}_k||^2 / \sigma^2$$

$$\approx \sum_{k=1}^K e^{-||\mathbf{x}_n - \boldsymbol{\mu}_k||^2 / \sigma^2}$$

2. Compute the marginal likelihood (cost).

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

$q_{kn} \rightarrow z_k$
k-means assignment

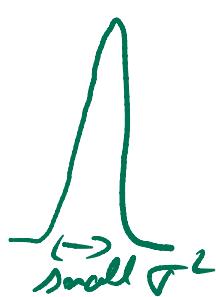
3. **M-step:** Update $\boldsymbol{\mu}_k^{(t+1)}, \Sigma_k^{(t+1)}, \pi_k^{(t+1)}$.

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

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

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

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



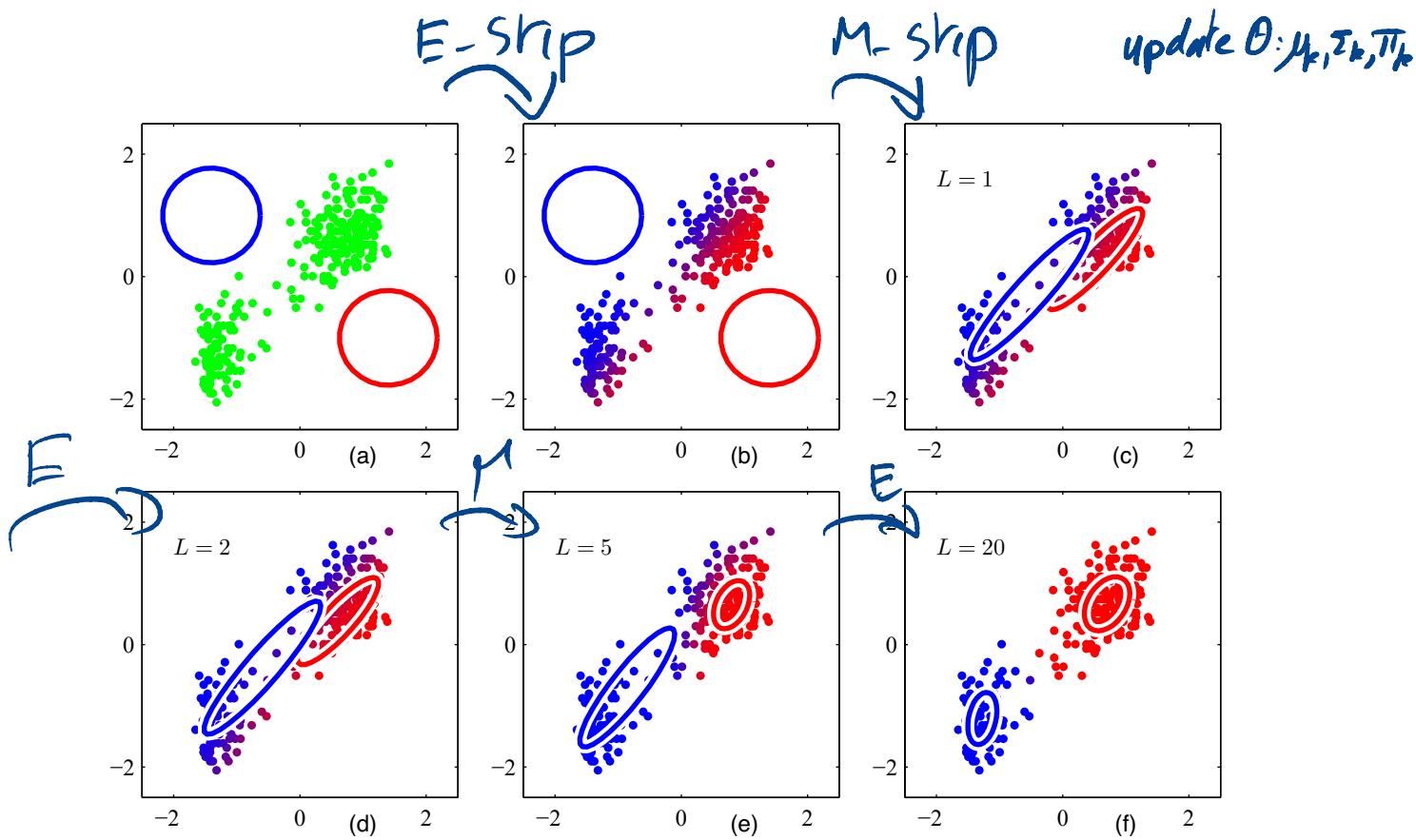
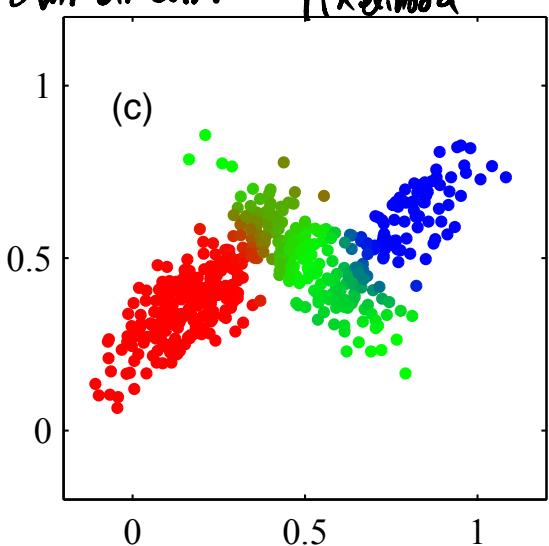


Figure 1: EM algorithm for GMM

Posterior distribution

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

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



$$P(Z_n = k | \mathbf{x}_n, \boldsymbol{\theta}) = \frac{\text{prior likelihood}}{\text{M.L.}}$$

$$\begin{aligned} &= \frac{\pi_k P(Z_n = k | \boldsymbol{\theta}) P(\mathbf{x}_n | Z_n, \boldsymbol{\theta})}{\sum_{k=1}^K \pi_k P(Z_n = k) P(\mathbf{x}_n | Z_n, \boldsymbol{\theta})} \\ &= \frac{\pi_k N(\mathbf{x}_n | Z_n, \boldsymbol{\theta})}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | Z_n, \boldsymbol{\theta})} = q_{kn} \end{aligned}$$

EM in general

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

(maximizing posterior)

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.

hidden vars z