

Expectation Maximization = probabilistic variant of K-means.

- Make probabilistic calculations, instead of hard assignments.
- More powerful algorithm of doing the same task.

Clustering:

Inputs: n objects (data points) $\{x_i\}_{i=1}^n$ and K clusters.

Goal: group the points into several groups.

group ID: $z_i \in \{1, \dots, K\}$, for each x_i .

Intro to EM (Idea)

K-Means

Cluster:

Assignment:

$$z_i = \underset{k=1, \dots, K}{\operatorname{argmin}} \|x_i - \mu_k\|^2$$

cluster i

// iterate the centroids and find the one with smallest squared distance to x_i , and set it to z_i .

Centroid:

$$\mu_k = \frac{\sum_{i=1}^n \mathbb{I}(z_i=k) x_i}{\sum_{i=1}^n \mathbb{I}(z_i=k)}$$

Centroid

$$\mu_k = \frac{\sum_{i=1}^n \mathbb{I}(z_i=k)}{\sum_{i=1}^n \mathbb{I}(z_i=k)}$$

where $\mathbb{I}(z_i=k) = \begin{cases} 1, & z_i = k \\ 0, & z_i \neq k \end{cases}$

// calculate the average location of those belong to cluster $z_i = k$.

EM

$$p(z_i = k) = \frac{\exp(-\|x_i - \mu_k\|^2/\lambda)}{\sum_{i=1}^k \exp(-\|x_i - \mu_k\|^2/\lambda)}$$

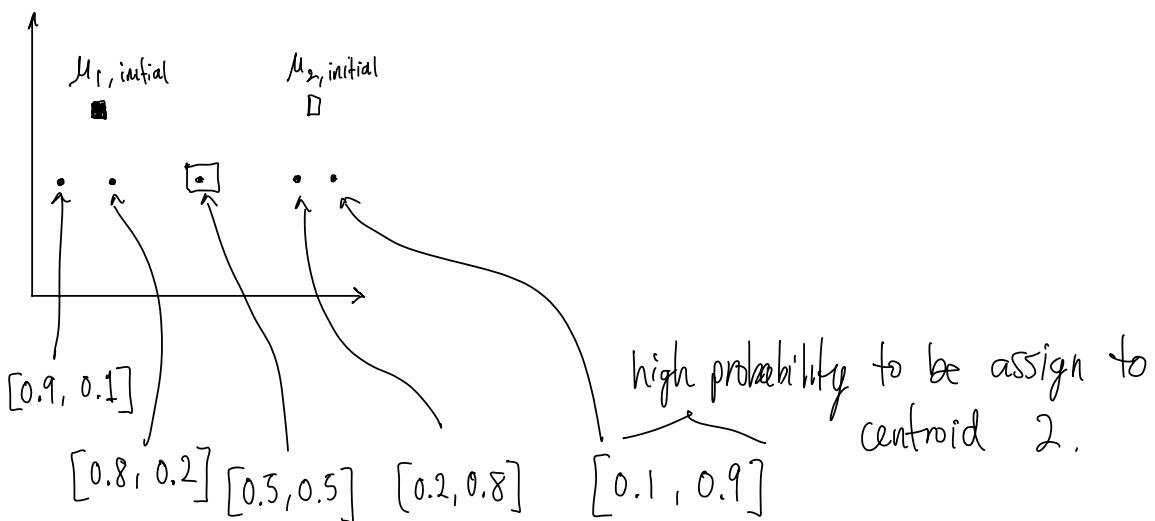
// sigmoid function.

const

replaced the indicator

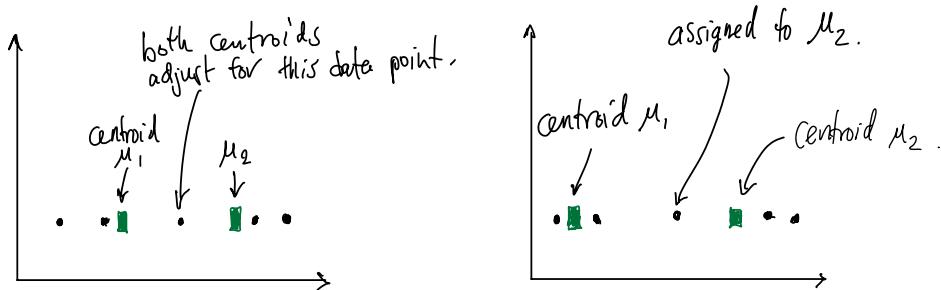
$$\mu_k = \frac{\sum_{i=1}^n p(z_i = k) x_i}{\sum_{i=1}^n p(z_i = k)}$$

// weighted average. Instead of indicator function, it uses a probability.



For $\mu_1 = \frac{0.9X_1 + 0.8X_2 + 0.5X_3 + 0.2X_4 + 0.1X_5}{0.9 + 0.8 + 0.5 + 0.2 + 0.1}$ // weighted average,
 contributes the most.

$$\mu_2 = \frac{0.1X_1 + 0.2X_2 + 0.5X_3 + 0.8X_4 + 0.9X_5}{0.1 + 0.2 + 0.5 + 0.8 + 0.9}$$

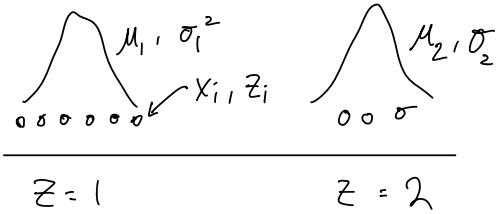


- More generally, EM algorithm is a special optimization algorithm for maximum likelihood estimation of Gaussian mixtures (or latent variable models).

Probabilistic Modeling of Clustering: General

Probabilistic approach:

- Assume a joint distribution $p(x, z | \theta)$ that generates data and labels (x, z) . $x = \text{data}$, $z = \text{label}$.
- Estimate parameter $\theta = \{\pi_k, \mu_k, \sigma_k^2\}_{k=1}^K$
- Infer the labels from the posterior distribution $p(z | x, \theta)$.



data points | label | parameters μ, σ, π

$$p(x, z) = p(z)p(x|z) \quad // \text{multiplication rule.}$$

joint distribution of data and label

①

$$p(z=k) = \pi_k$$

$\left\{ \begin{array}{l} \pi_k \geq 0 \\ \sum_k \pi_k = 1 \end{array} \right.$ given the axioms of probability

probability of $z=k$.

②

$$p(x | z=k) = N(x | \mu_k, \sigma_k^2) \quad \text{Assume Gaussian,}$$

data given μ_k, σ_k^2 (from $z=k$).

③

$$\underbrace{p(x, z=k)}_{\substack{\text{joint} \\ x \cap z=k}} = \pi_k \frac{N(x | \mu_k, \sigma_k^2)}{p(z=k)}$$

$$p(z=k) p(x | z=k) \quad // \text{law of conditional probability}$$

unknown parameters $\theta = \{\pi_k, \mu_k, \sigma_k^2\}_{k=1}^K$

$$p(x, z | \theta)$$

we don't observe.

MLE with Complete Information

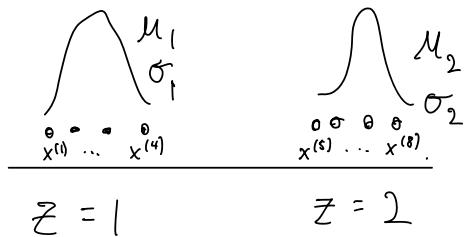
If we observe both data and labels

$\{x_i, z_i\}_{i=1}^n$ (complete information), we can estimate

the parameter θ by maximizing the joint probability:

$$\max_{\theta} \sum_{i=1}^n \log p(\underbrace{x_i, z_i}_{\text{Joint distribution}} | \theta) .$$

Example:



$$\pi_1 = p(z=1) = \frac{\#(z=1)}{n} = \frac{4}{8} = \frac{1}{2} .$$

$$\mu_1 = \frac{\sum_{i=1}^n \mathbb{I}(z_i=1) x_i}{\sum_{i=1}^n \mathbb{I}(z_i=1)} = \frac{\sum \text{values of } x^{(1)} \text{ to } x^{(4)}}{4}$$

$$\sigma_1^2 = \text{Var}\left(\{x_i | z_i=1\}\right) \quad // \text{variance of the data points with } z_i=1.$$

$$\Rightarrow P(x_i, z_i=k | \theta) = \pi_k \cdot N(x | \mu_k, \sigma_k^2)$$

Joint Log Likelihood Function :

$$\underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(z_i = k) \log P(x_i | z_i = k | \theta)$$

only fires when
 $z_i = k$

$$= \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(z_i = k) \log (\pi_k N(x_i | \mu_k, \sigma_k^2))$$

$$= \sum_{i=1}^n \sum_{k=1}^K \mathbb{I}(z_i = k) (\log \pi_k + \log N(x_i | \mu_k, \sigma_k^2))$$

If we max, then (without further derivation) :

$$\pi_k = \frac{\sum_{i=1}^n \mathbb{I}(z_i = k)}{n}$$

$$\mu_k = \text{empirical mean of } S_k : \mu_k = \frac{\sum_{i=1}^n \mathbb{I}(z_i = k) x_i}{\sum_{i=1}^n \mathbb{I}(z_i = k)}$$

$$\sigma_k^2 = \text{empirical variance of } S_k : \text{var}(\{x_i | z_i = k\})$$

If there is complete information, estimating θ is easy.
use the empirical mean and variance.

Gaussian Mixture Models: Incomplete Information

Clustering can be formulated as parameter estimation in GMM.
If we only observe data $\{x_i\}_{i=1}^n$, but missing $\{z_i\}_{i=1}^n$,
(incomplete information), we shall estimate θ by maximizing
the marginal probability of x_i , (instead of joint $\{x_i, z_i\}$).

Marginal Log-Likelihood Function:

$$\sum_{i=1}^n \log p(x_i | \theta) = \sum_{i=1}^n \log \sum_{z_i} p(x_i, z_i | \theta) \quad // \text{LLTP}$$

\curvearrowleft sum over z_i (from 1 to K)

we only observe x_i . z_i = missing information.

$$\theta := \{\pi_k, \mu_k, \sigma_k^2\}_{k=1}^K, \text{ where:}$$

μ_k and σ_k^2 = mean, variance of the k^{th} cluster,

π_k = its percent in the data.

Need numerical method, as there is no closed form solution.

In other words:

$$\begin{aligned}
 p(x|\theta) &= \sum_k p(x, z=k | \theta) \\
 &= \sum_k \pi_k N(x | \mu_k, \sigma_k^2)
 \end{aligned}
 \quad \left. \begin{array}{l} \text{known} \\ \text{unknown} \end{array} \right\} \quad \left. \begin{array}{l} \text{unknown} \\ \text{unknown} \end{array} \right\} \quad \left. \begin{array}{l} \text{mixture of} \\ \text{Gaussian} \\ \text{distributions} \end{array} \right\}$$

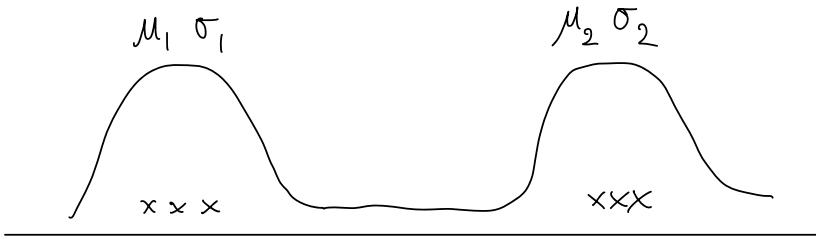
$$\Rightarrow \underbrace{\sum_{i=1}^n \log p(x_i | \theta)}_{\text{want to find}} = \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k N(x | \mu_k, \sigma_k^2) \right)$$

want to find
 θ s.t. this function
 is maximized.

objective function that we want to
 maximize.

$$\theta = \{\pi_k, \mu_k, \sigma_k^2\}_{k=1}^K$$

Example:



$$\begin{aligned}
 p(x|\theta) &= \frac{1}{2} \exp \left(-\frac{(x-\mu_1)^2}{2\sigma_1^2} \right) \left(\frac{1}{\sqrt{2\pi}\sigma_1} \right) \quad // \text{first mixture} \\
 &\quad + \frac{1}{2} \exp \left(-\frac{(x-\mu_2)^2}{2\sigma_2^2} \right) \left(\frac{1}{\sqrt{2\pi}\sigma_2} \right) \quad // \text{2nd mixture}
 \end{aligned}$$

No closed form solution, Use EM Algorithm!

SUMMARY: Gaussian Mixture Models:

- probabilistic model for clustering, in which each cluster is represented by a Gaussian distribution,
- Density function of a GMM is a weighted linear combination of several Gaussian density functions:

$$p(x|\theta) = \sum_{k=1}^K \pi_k N(x; \mu_k, \sigma_k^2), \text{ where :}$$

1) $\theta = \{\pi_k, \mu_k, \sigma_k\}_{k=1}^K$

2) $\{\pi_k\}_{k=1}^K$ = set of mixture weights that satisfy the axioms of probability: $\sum_{k=1}^K \pi_k = 1, \pi_k \geq 0 \forall k.$

3) $N(x; \mu_k, \sigma_k^2) = \underbrace{\frac{1}{\sqrt{2\pi}\sigma_k}}_{\text{each of these is called "component" of } p(x|\theta)} \exp\left(-\frac{1}{2\sigma_k^2}(x-\mu_k)^2\right)$.

each of these is called "component" of $p(x|\theta)$, which corresponds to a cluster

Draw latent label Z : $p(Z=k|\theta) = \pi_k$ // draw $Z=k$ with probability π_k .

Draw observation X : $p(X=x | Z=k, \theta) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{1}{2\sigma_k^2}(x-\mu_k)^2\right)$
 // If $Z=k$, X is a Gaussian with mean μ_k , variance σ_k^2 .

Joint Probability of X and Z : What is the probability of observation x belonging to cluster k ? used for complete problems with known X and Z (data and labels).

$$\begin{aligned} p(X=x, Z=k | \theta) &= p(Z=k | \theta) \cdot p(X=x | Z=k, \theta) \\ &= \pi_k \times N(x; \mu_k, \sigma_k^2). \end{aligned}$$

Marginal Distribution of X : can calculate by summing over Z .
Used for problems without known Z (labels).

$$p(X=x | \theta) = \underbrace{\sum_{k=1}^K p(X=x, Z=k | \theta)}_{\text{summing over all values of } k \text{ in } Z=k} = \sum_{k=1}^K w_k \times N(x; \mu_k, \sigma_k^2)$$

\therefore This is the density function of the GMM.

How to infer cluster ID of an observation X , based on the posterior distribution:

$$\begin{aligned} p(Z=k | X=x, \theta) &= \frac{p(X=x, Z=k | \theta)}{p(X=x | \theta)} \\ &= \frac{w_k \times N(x; \mu_k, \sigma_k^2)}{\sum_{l=1}^K w_l \times N(x; \mu_l, \sigma_l^2)} \end{aligned}$$

EM Algorithm: iterative method for maximizing the marginal likelihood (since we don't observe the latent variable z):

$$\max_{\theta} \sum_{i=1}^n \log P(x_i | \theta).$$

1) Initialize $\theta_0 = \{w_{k,0}, \mu_{k,0}, \sigma_{k,0}^2\}_{k=1}^K$. (random guess)

2) Denote $\gamma_{ik} = P(z^{(i)} = k | X = x^{(i)}, \theta) = \frac{\text{posterior probability for the } i^{\text{th}} \text{ data point belonging to cluster } k}{\text{sum of posterior probabilities}}$

3) Perform until convergence:

i) E-step: Fixing θ , update $\{\gamma_{ik}\}$

Given θ_t , "impute" the missing labels by drawn samples from the posterior distribution,

$$z_i \sim \underbrace{P(z | x_i; \theta_t)}_{\text{posterior distribution}}$$

$$\begin{aligned} \gamma_{ik, t+1} &= P(z = k | X = x^{(i)}; \theta_t) \\ &= \frac{w_{k,t} \cdot N(x^{(i)}; \mu_{k,t}, \sigma_{k,t}^2)}{\sum_{l=1}^K w_{l,t} N(x^{(i)}; \mu_{l,t}, \sigma_{l,t}^2)} // k^{\text{th}} \text{ cluster for one data point } x^{(i)}. \\ \text{where} \end{aligned}$$

$$\theta_t = \{\omega_{k,t}, \mu_{k,t}, \sigma_{k,t}\}_{k=1}^K = \text{values at } t^{\text{th}} \text{ iteration}$$

$\gamma_{ik,t}$ = posterior probability @ t^{th} iteration.

This step calculates the posterior probability of the cluster IDs of each point, assuming

θ_t is the free parameter. For each data point $x^{(i)}$, the sum of γ_{ik} over all clusters = 1.

ii) M-Step : Fixing $\{\gamma_{ik}\}$, update θ :

update θ by maximizing the expected joint likelihood.

$$l(\theta) \triangleq \theta^{t+1} = \arg \max_{\theta} \sum_{i=1}^n E_{z_i \sim p(\cdot | x_i, \theta_t)} [\log p(x_i, z_i | \theta)]$$

$$\Rightarrow l(\theta) = \sum_{i=1}^n \log p(x_i | \theta) = \underbrace{\sum_{i=1}^n \sum_{k=1}^K \gamma_{ik}^t \log p(x_i, z_i = k | \theta)}_{\text{Max this.}}$$

parameter updates (provided without derivation):

updating mean : $\mu_{k,t+1} = \frac{\sum_{i=1}^n \gamma_{ik,t+1} x^{(i)}}{\sum_{i=1}^n \gamma_{ik,t+1}}$

(for cluster k)

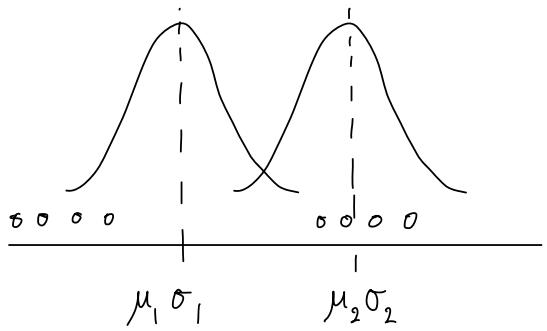
updating variance : $\sigma_{k,t+1}^2 = \frac{\sum_{i=1}^n \gamma_{ik,t+1} (x^{(i)} - \mu_{k,t+1})^2}{\sum_{i=1}^n \gamma_{ik,t+1}}$

updating weights : $w_{k,t+1} = \frac{\sum_{i=1}^n \gamma_{ik,t+1}}{\sum_{i=1}^n \sum_{k=1}^K \gamma_{ik,t+1}}$

where mean μ_k and variance σ_k^2 of the k -th cluster equals the empirical mean and variance of the data, when each data point is "weighted" by posterior probability of belonging to the k -th cluster, based on estimation $\gamma_{ik,t+1}$ from the current iteration.

converge to local optimum.

Implementation using previous example:



$$\theta_t = [\pi_k^t, \mu_k^t, \sigma_k^t]_{k=1}^K$$

After t iterations, these are the values.

$$P(\underbrace{Z^{(i)}=k}_{A} \mid \underbrace{X^{(i)}, \theta_t}_{B}) = \frac{P(\underbrace{X^{(i)}, Z^{(i)}=k}_{A \cap B} \mid \theta_t)}{\sum_l P(\underbrace{X^{(i)}, Z^{(i)}=l}_{B} \mid \theta_t)}$$

$$= \frac{\pi_k N(x^{(i)} \mid \mu_k^t, \sigma_k^t)}{\sum_l \pi_l N(x^{(i)} \mid \mu_l^t, \sigma_l^t)}$$

// law of conditional probability $\approx P(A \mid B) = \frac{P(A \cap B)}{P(B)}$

So, for example:

$$P(Z^{(i)}=1 \mid X^{(i)}, \theta_t) = \frac{\pi_1 \exp\left(-\frac{(X^{(i)} - \mu_1)^2}{2\sigma_1^2}\right) \frac{1}{\sqrt{2\pi}\sigma_1}}{\sum_l \left\{ \pi_l \exp\left(-\frac{(X^{(i)} - \mu_l)^2}{2\sigma_l^2}\right) \frac{1}{\sqrt{2\pi}\sigma_l} \right\}}$$

Derivation of EM algorithm : ↗ arbitrary distribution?

- Consider general "imputation distributions" $f(z|x)$
- we can construct a tight lower bound $LB(\theta, p)$ of the marginal log likelihood function:

$$l(\theta) \geq LB(\theta, p), \quad \forall p \text{ and } l(\theta) = \max_p LB(\theta, p)$$

- Maximizing $l(\theta)$ is then equivalent to maximizing $LB(\theta, p)$:

$$\max_{\theta} l(\theta) = \max_{\theta, p} LB(\theta, p)$$

- Optimizing θ and p alternatively (coordinate descent) yields EM algorithm.

②

Dive In:

$$l(\theta) \triangleq \max_{\theta} \sum_{i=1}^n \log P(x_i | \theta)$$

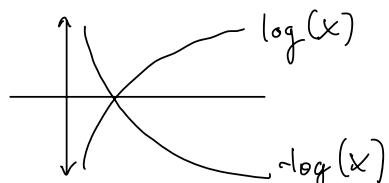
$$= \max_{\theta} \sum_{i=1}^n \log \sum_{z_i} P(x_i, z_i | \theta) \quad // \text{Marginal distribution} = \text{sum of joint distribution}$$

want to move the summation outside of $\log \Rightarrow$
use Jensen's Inequality:

for convex function f : $f(E[x]) \leq E[f(x)]$.



\log is concave function $\rightarrow -\log$ is a convex function.



$$\Rightarrow \sum_{i=1}^n \log \sum_{z_i} \frac{p(x_i, z_i | \theta)}{f(z_i | x_i)} p(z_i | x_i)$$

↗ multiply both num and den (equivalent)
 ↗ by arbitrary distribution f .

$$f(E[X])$$

$$= \sum_{i=1}^n \overbrace{\log E_{z_i \sim f(\cdot | x_i)} \left[\frac{p(x_i, z_i | \theta)}{f(z_i | x_i)} \right]}^{\text{rewrite summation as expectation over } f.}$$

$$\geq \sum_{i=1}^n E_{z_i \sim f(\cdot | x_i)} \underbrace{\log \left(\frac{p(x_i, z_i | \theta)}{f(z_i | x_i)} \right)}_{E[f(x)]} \triangleq LB(\theta, f).$$

If $p(z | x, \theta) = p(z | x, \theta)$ then

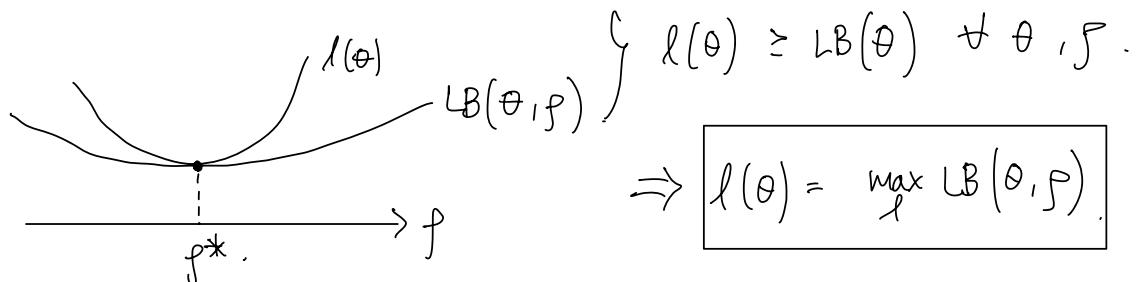
$$\frac{p(x_i, z_i | \theta)}{f(z_i | x_i)} = \frac{p(x_i, z_i | \theta)}{\underbrace{p(z_i | x_i, \theta)}_{\text{posterior}}} = \underbrace{p(x_i | \theta)}_{\text{marginal}}$$

$$\Rightarrow LB(\theta, \ell^*) = \sum_{i=1}^n E_{z \sim f(\cdot | x_i)} \log p(x_i | \theta).$$

the expectation does not depend on z_i anymore,
thus can be removed.

$$= \sum_{i=1}^n \log p(x_i | \theta) = L(\theta).$$

By Jensen Inequality, we showed that LB is always a Lower Bound of $\ell(\theta)$ for any f . If we choose f to be to be the posterior distribution, then the lower bound matches the marginal likelihood.



$\Rightarrow \max_{\theta, f} LB(\theta, f)$: Use coordinate descent:

1) Initialize θ_0 .

$$p(z|x, \theta_t)$$

2) Iterate:

Fix θ_t : $f_{t+1} = \overbrace{\arg \max_f LB(\theta_t, f)}^{\text{expectation step}}$

Fix f_{t+1} : $\theta_{t+1} = \arg \max_{\theta} LB(\theta, f_{t+1})$ // maximization step.

To see that:

$$\begin{aligned} LB(\theta, f_{t+1}) &= \sum_{i=1}^n E_{z_i \sim p(\cdot | x_i)} \left[\log p(x_i, z_i | \theta) - \underbrace{\log p(z_i | x_i)}_{\text{constant}} \right] \\ &= \underbrace{\sum_{i=1}^n E_{z_i \sim p(\cdot | x_i)} [\log (x_i, z_i | \theta)]}_{\text{Expected joint likelihood}} + \text{constant}, \end{aligned}$$

Properties: this procedure:

- ① Monotonically decrease $LB(\theta, f)$, and $\ell(\theta)$.
- ② converge to local optima of $\ell(\theta)$.