

# Mixture Models and EM

(Part 1)

18-698 / 42-632

Neural Signal Processing  
Prof. Byron Yu

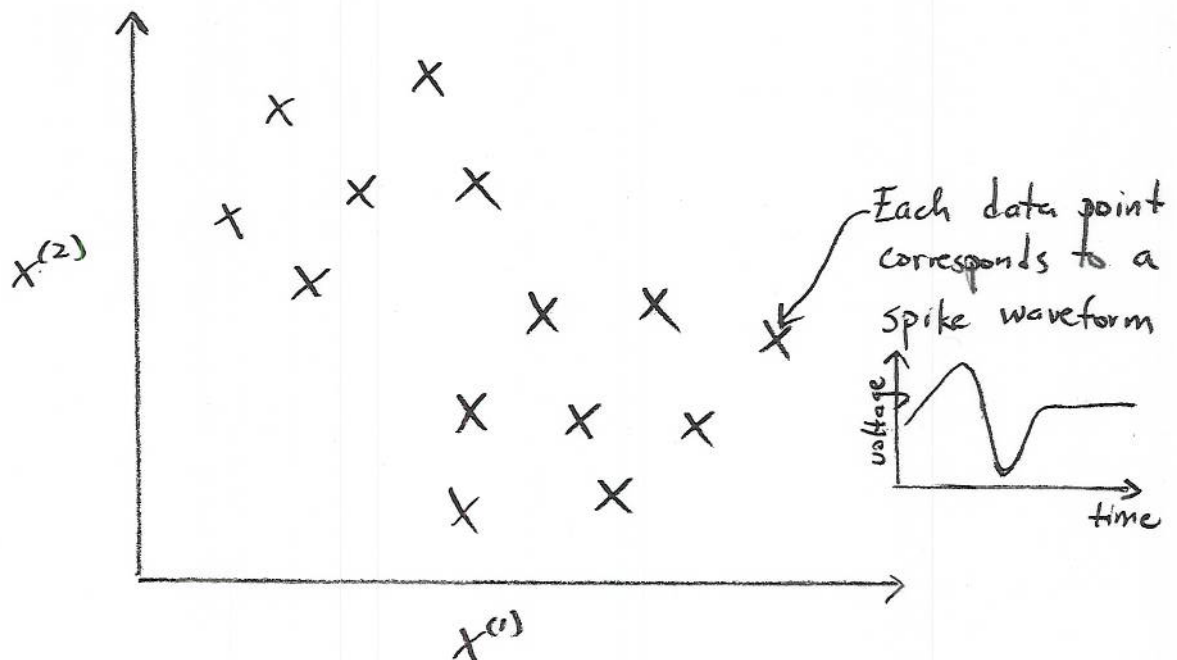
## A) K-means Clustering

Suppose we have a data set  $x_n \in \mathbb{R}^D$   
where  $n=1, \dots, N$ .

Goal: Partition the data set into some  
number  $K$  of clusters.

For now, assume  $K$  is given.

Picture to have in mind:



How would you partition this dataset into  
 $K=2$  clusters?

## Intuitive definition of a cluster:

A group of data points whose inter-point distances are small compared with the distances to points outside the cluster.

## Let's formalize this notion:

Let  $\mu_k \in \mathbb{R}^D$  where  $k=1, \dots, K$  be the "prototype" associated with the  $k$ th cluster.

$$r_{nk} = \begin{cases} 1 & \text{if } x_n \text{ belongs to the } k\text{th cluster} \\ 0 & \text{else} \end{cases}$$

objective  
function

$$\rightarrow J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

Find  $\{r_{nk}\}$  and  $\{\mu_k\}$  such that  $J$  is minimized.

This can be solved by iterating:

- Minimize  $J$  wrt.  $r_{nk}$ , keeping  $\mu_k$  fixed ("E-step")
- Minimize  $J$  wrt.  $\mu_k$ , keeping  $r_{nk}$  fixed ("M-step")

### A.1) E-step for K-means

Constraint:  $\{r_{n1}, \dots, r_{nK}\}$  is a set of  $(K-1)$  zeros and a single 1.

Can optimize  $J$  wrt.  $r_{nk}$  for each  $n$  separately.

For each  $n$ , assign:

$$r_{nk} = \begin{cases} 1 & \text{if } k = \underset{j}{\operatorname{argmin}} \|x_n - \mu_j\|^2 \\ 0 & \text{else} \end{cases}$$

In words: Assign each data point to the closest cluster center.

### A.2) M-step for K-means

$$\frac{\partial J}{\partial \mu_k} = 2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0$$

$$\mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}}$$

In words: Set  $\mu_k$  equal to the mean of all data points assigned to cluster  $k$ .

### A.3) Convergence of the K-means algorithm.

- The E-step and M-step should be iterated until there is no further change in cluster assignments, or until some maximum number of iterations is exceeded.
- Each iteration reduces the objective function  $J$ , so the algorithm is guaranteed to converge.  
(We will prove this later in the context of the EM algorithm.)
- Convergence is guaranteed to a local (rather than a global) minimum of  $J$ .
- If there are multiple local optima, the particular local optimum reached depends on the initialization of the  $\{\mu_k\}$ .

[Let's look at an example of K-means in action.  
PRML Figures 9.1 and 9.2.]

# Mixture Models and EM

(Part 2)

18-698 / 42-632

Neural Signal Processing  
Prof. Byron Yu

## B) Mixtures of Gaussians

Let  $z \in \{1, \dots, K\}$  be a random variable

$$P(z=k) = \pi_k, \text{ where } k=1, \dots, K$$

(Note: For  $P(z)$  to be a valid probability distribution:  
 $0 \leq \pi_k \leq 1$   
 $\sum_{k=1}^K \pi_k = 1$ )

$$P(\underline{x} | z=k) = N(\underline{x} | \underline{\mu}_k, \Sigma_k)$$

Thus,

$$P(\underline{x}) = \sum_z P(\underline{x} | z) P(z) = \sum_{k=1}^K N(\underline{x} | \underline{\mu}_k, \Sigma_k) \cdot \pi_k$$

What is the graphical model for  $P(\underline{x}, z)$ ?



### B.1) Maximum likelihood parameter estimation

Goal: Fit  $\mu_k, \Sigma_k, \pi_k$  to training data

$x_1, \dots, x_N$ .

$$P(\{x\} | \theta) = \prod_{n=1}^N P(x_n)$$

shorthand for

$x_1, \dots, x_N$

shorthand for

$\mu_k, \Sigma_k, \pi_k (k=1, \dots, K)$

$$= \prod_{n=1}^N \sum_{k=1}^K N(x_n | \mu_k, \Sigma_k) \cdot \pi_k$$

$$\log P(\{x\} | \theta) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K N(x_n | \mu_k, \Sigma_k) \cdot \pi_k \right\}$$

call this  $\mathcal{L}$

i) Find  $\mu_k$

$$\frac{\partial \mathcal{L}}{\partial \mu_k} = \sum_{n=1}^N \frac{1}{\sum_{j=1}^K N(x_n | \mu_j, \Sigma_j) \cdot \pi_j} \cdot \pi_k \cdot \left( \frac{\partial}{\partial \mu_k} N(x_n | \mu_k, \Sigma_k) \right) \quad (1)$$

$$\left[ \begin{aligned} \frac{\partial}{\partial \mu_k} N(x_n | \mu_k, \Sigma_k) &= \frac{\partial}{\partial \mu_k} \left( \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_k|^{\frac{1}{2}}} e^{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)} \right) \\ &= N(x_n | \mu_k, \Sigma_k) \cdot \Sigma_k^{-1} (x_n - \mu_k) \end{aligned} \right]$$

Plugging into (1),

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mu_k} &= \sum_{n=1}^N \frac{N(x_n | \mu_k, \Sigma_k) \cdot \pi_k}{\underbrace{\sum_{j=1}^K N(x_n | \mu_j, \Sigma_j) \cdot \pi_j}_{\text{call this } \gamma_{nk}}} \cdot \Sigma_k^{-1} (x_n - \mu_k) \\ &= \Sigma_k^{-1} \sum_{n=1}^N \gamma_{nk} (x_n - \mu_k) \\ &= 0 \end{aligned}$$

Letting  $N_k = \sum_{n=1}^N \gamma_{nk}$ ,

$$\boxed{\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} x_n} \quad (2)$$

Let's try to understand this result intuitively.

$$\gamma_{nk} = P(z_n = k | x_n)$$

$$= \frac{P(x_n | z_n = k) P(z_n = k)}{P(x_n)}$$

$$= \frac{N(x_n | \mu_k, \Sigma_k) \cdot \pi_k}{\sum_{j=1}^K N(x_n | \mu_j, \Sigma_j) \cdot \pi_j} \quad \left( \begin{array}{l} \text{Note:} \\ 0 \leq \gamma_{nk} \leq 1 \\ \sum_{k=1}^K \gamma_{nk} = 1 \end{array} \right)$$

$\gamma_{nk}$  is the "responsibility" that cluster  $k$  takes for explaining the observation  $x_n$ .

Whereas  $\pi_k$  is the prior probability that  $z_n = k$ ,  
 $\gamma_{nk}$  is the posterior probability that  $z_n = k$   
 once we have observed  $x_n$ .

$N_k$  is the effective number of data points assigned to cluster  $k$ .

Thus, in (2),  $\mu_k$  is obtained by taking a weighted mean, where the weights are given by the responsibilities ( $0 \leq \gamma_{nk} \leq 1$ ).  
↑  
not responsible     ↑     very responsible

This is reminiscent of k-means, but weights can now range from 0 to 1 (soft assignments).



ii) Find  $\Sigma_k$

Setting  $\frac{\partial \mathcal{L}}{\partial \Sigma_k} = 0$ ,

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (x_n - \mu_k)(x_n - \mu_k)^T \quad (3)$$

(we won't show this here)

(3) has an intuitive interpretation like in (2).

iii) Find  $\pi_k$

Here, we need to perform a constrained optimization

since  $\sum_{k=1}^K \pi_k = 1$ .

Instead of maximizing  $\mathcal{L}$ , we will maximize

$$\mathcal{L}' = \mathcal{L} + \underset{\substack{\uparrow \\ \text{Lagrange multiplier}}}{\lambda} \left( \sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial \mathcal{L}'}{\partial \pi_k} = \frac{\partial \mathcal{L}}{\partial \pi_k} + \lambda$$

$$= \sum_{n=1}^N \frac{1}{\sum_{j=1}^K N(x_n | \mu_j, \Sigma_j) \pi_j} \cdot N(x_n | \mu_k, \Sigma_k) + \lambda$$

$$= 0$$

Multiplying both sides by  $\pi_k$ ,

$$\sum_{n=1}^N \delta_{nk} + \lambda \cdot \pi_k = 0$$

$$\pi_k = -\frac{N_k}{\lambda}$$

(4)

Enforcing the constraint  $\sum_{k=1}^K \pi_k = 1$ ,

$$-\frac{\sum_{k=1}^K N_k}{\lambda} = 1$$

$$\lambda = -N$$

Plugging into (4),

$$\boxed{\pi_k = \frac{N_k}{N}}$$

(5)

This is the effective number of data points assigned to class  $k$ , divided by the total number of data points.

### Very important note:

(2), (3), and (5) do not constitute a closed-form solution because the responsibilities  $\gamma_{nk}$  (which appear in (2), (3), and (5)) depend on  $\mu_k$ ,  $\Sigma_k$ , and  $\pi_k$ .

This suggests an iterative scheme, which turns out to be an instance of the Expectation-Maximization (EM) algorithm.

The EM algorithm is a powerful and general method for finding the maximum likelihood parameters for models with latent variables.

In the mixture of Gaussians,  $z$  is a latent variable because it is not observed (only  $x$  is observed).

## B.2) EM algorithm for Mixture of Gaussians

1) Initialize  $\mu_k, \Sigma_k, \pi_k$ .

2) E-step: Evaluate responsibilities given current parameter values

$$\gamma_{nk} = \frac{N(x_n | \mu_k, \Sigma_k) \pi_k}{\sum_{j=1}^K N(x_n | \mu_j, \Sigma_j) \pi_j}$$

3) M-step: Re-estimate parameters using current responsibilities

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \cdot x_n$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (x_n - \mu_k^{\text{new}})(x_n - \mu_k^{\text{new}})^T$$

$$\pi_k^{\text{new}} = \frac{N_k}{N}$$

where  $N_k = \sum_{n=1}^N \gamma_{nk}$

4) Evaluate log likelihood of data

$$\log P(\{x\} | \theta) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K N(x_n | \mu_k, \Sigma_k) \cdot \pi_k \right\}$$

and check for convergence of either the parameters or the log likelihood. If convergence criterion not satisfied, return to step 2).

- We will soon show that each iteration of the EM algorithm is guaranteed to increase the log likelihood  $\log P(\{x\} | \theta)$ .
- See PRML Figure 9.8 for an example of EM in action as applied to a mixture of Gaussians.



# Mixture Models and EM

18-698 / 42-632

Neural Signal Processing  
Prof. Byron Yu

(Part 3)

## C) Relating EM for Gaussian Mixtures to Classification ← call this ①

Call this ②

During the training phase, the goal in both cases is to estimate the model parameters from the training data.

Key difference: In ①, the class labels are known.

In ②, the class labels are not known.

In ①, we were able to estimate the model parameters in closed form:

$$\pi_k = \frac{N_k}{N}, \text{ where } N_k = \sum_{n \in C_k} 1 \quad (1)$$

$$\mu_k = \frac{1}{N_k} \sum_{n \in C_k} \underline{x}_n \quad (2)$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n \in C_k} (\underline{x}_n - \mu_k)(\underline{x}_n - \mu_k)^T \quad (3)$$

In ②, it is unknown whether  $n \in C_k$  (i.e. we don't have class labels), which leads to a "chicken and egg" problem.

- If the class labels were known, it would be easy to compute model parameters using (1)-(3).
- If the model parameters were known, it would be easy to compute a distribution on the class labels  $P(C_k | x)$ . (see p.9 of Classification notes)

Using the EM algorithm, we can "bootstrap" our way out of the chicken and egg problem.

1) Initialize model parameters

E-step

2) Estimate class labels given current model parameters.

M-step

3) Estimate model parameters given current class labels (represented as a distribution over class labels).

4) Go to 2).

In ②, the training phase uses the EM algorithm, in which:

- The E-step is reminiscent of test phase in ①
- The M-step is reminiscent of training phase in ①.

In ②, the test phase typically involves running a single E-step on the test data (no iterating needed).

## D) The EM Algorithm in General

### D.1) Motivation

You have a model that you would like to fit to your training data.

If the model has latent variables (very likely in neural signal processing), the EM algorithm provides a recipe for fitting the model.

Because of its generality, the EM algorithm is among the most important and widely used tools in machine learning.

We will be using it repeatedly in the rest of this course.

## D2) Decomposition of data likelihood

Let  $X$  denote all observed variables

$Z$  denote all latent variables

$\theta$  denote all model parameters.

Goal: Maximize  $p(X|\theta)$  w.r.t.  $\theta$ .

Let  $q(Z)$  be some distribution over latent variables.

For any choice of  $q(Z)$ , we can decompose

$\log p(X|\theta)$  into a sum of two terms:

$$\log p(X|\theta) = \mathcal{L}(q, \theta) + \text{KL}(q \| p), \quad (4)$$



where

$$\mathcal{L}(q, \theta) = \sum_z q(z) \log \frac{p(X, z | \theta)}{q(z)} \quad (5)$$

"Joint distribution"

$$KL(q \| p) = - \sum_z q(z) \log \frac{p(z | X, \theta)}{q(z)} \quad (6)$$

"Posterior distribution"

Notes:

- $\mathcal{L}(q, \theta)$  is a scalar that depends only on  $q$  and  $\theta$  because  $X$  is observed and  $z$  is summed out.
- $KL(q \| p)$  is the Kullback-Leibler divergence between  $q(z)$  and  $p(z | X, \theta)$ .

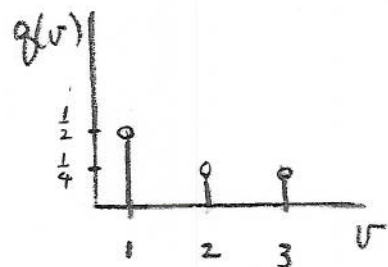
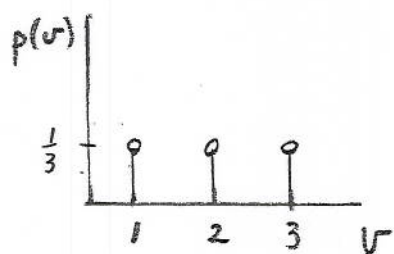
The KL divergence is a scalar that measures the distance between probability distributions.

In general, KL divergence is written:

$$KL(q \| p) = - \sum_v q(v) \log \frac{p(v)}{q(v)}$$

Let's do a simple example to illustrate KL divergence.





$$\begin{aligned}
 KL(q \parallel p) &= - \left( q(1) \log \frac{p(1)}{q(1)} + q(2) \log \frac{p(2)}{q(2)} + q(3) \log \frac{p(3)}{q(3)} \right) \\
 &= - \left( \frac{1}{2} \cdot \log \frac{\frac{1}{3}}{\frac{1}{2}} + \frac{1}{4} \log \frac{\frac{1}{3}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{3}}{\frac{1}{4}} \right) \\
 &= 0.0589
 \end{aligned}$$

$$KL(p \parallel p) = - \sum_v p(v) \log \frac{p(v)}{p(v)} = 0$$

Facts:  $KL(q \parallel p) \geq 0$  for any  $p$  and  $q$ .  
 $KL(q \parallel p) = 0$  iff  $p = q$ .

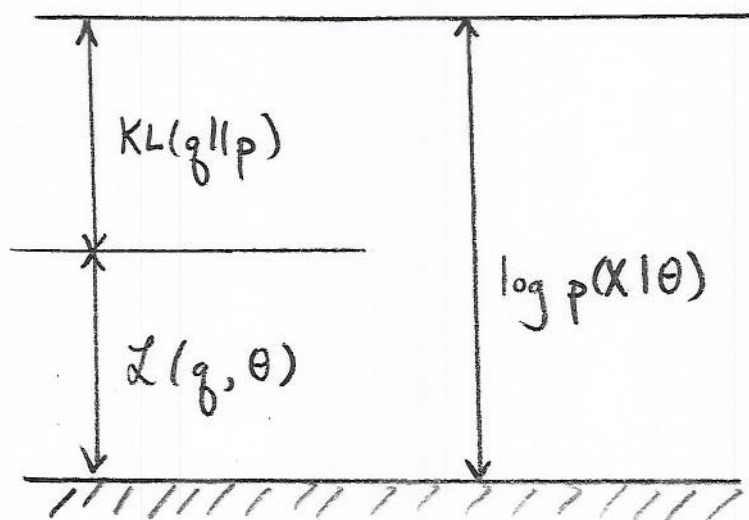
- In (6),  $KL(q \parallel p) \geq 0$ , with equality iff  $q(z) = p(z | X, \theta)$ .
- Thus,  $\log p(X | \theta) \geq \mathcal{L}(q, \theta)$   
 $\Rightarrow \mathcal{L}(q, \theta)$  is a lower bound for  $\log p(X | \theta)$
- The lower bound is tight (i.e.  $\log p(X | \theta) = \mathcal{L}(q, \theta)$ ) when  $q(z)$  is chosen to be  $p(z | X, \theta)$ .

Let's verify the decomposition (4):

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z q(z) \log \frac{p(z|x, \theta) p(x|\theta)}{q(z)} \\&= \sum_z q(z) \log \frac{p(z|x, \theta)}{q(z)} + \sum_z q(z) \log p(x|\theta) \\&= -KL(q||p) + \log p(x|\theta) \cdot \sum_z q(z) \end{aligned}$$

$$\Rightarrow \log p(x|\theta) = \mathcal{L}(q, \theta) + KL(q||p) //$$

Picture to have in mind



### D.3) The EM algorithm

Goal: Maximize  $\log p(X|\theta)$  w.r.t.  $\theta$

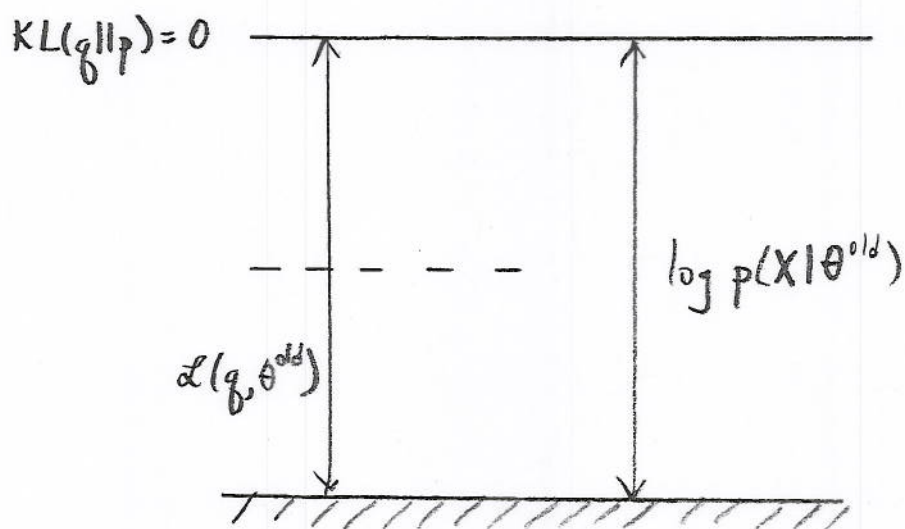
Approach: Iteratively maximize the lower bound  $\mathcal{L}(q, \theta)$  w.r.t.  $q$  and  $\theta$ .

E-step: maximize  $\mathcal{L}(q, \theta)$  w.r.t.  $q$  while keeping  $\theta$  fixed.

M-step: maximize  $\mathcal{L}(q, \theta)$  w.r.t.  $\theta$  while keeping  $q$  fixed.

E-step

With fixed  $\theta$ ,  $\mathcal{L}(q, \theta)$  is maximized when  $KL(q||p) = 0$ , corresponding to  $q(z) = p(z|X, \theta)$ .



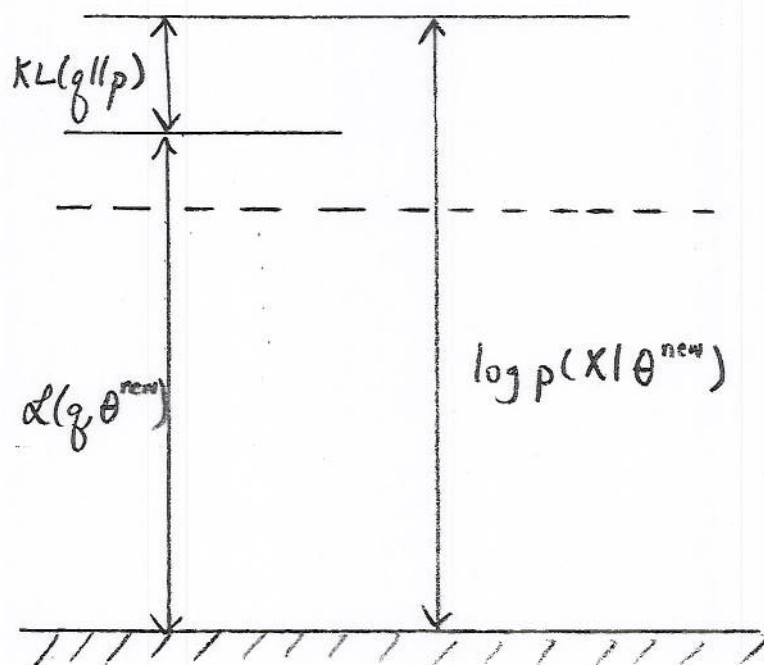
M-step

$$\mathcal{L}(q, \theta) = \sum_z q(z) \log p(x, z | \theta) - \sum_z q(z) \log q(z)$$

no  $\theta$  dependence

With fixed  $q$ , maximizing  $\mathcal{L}(q, \theta)$  is equivalent to maximizing  $\sum_z q(z) \log p(x, z | \theta)$ .

This is  $E_q[\log p(x, z | \theta)]$ , the "expected log joint distribution".

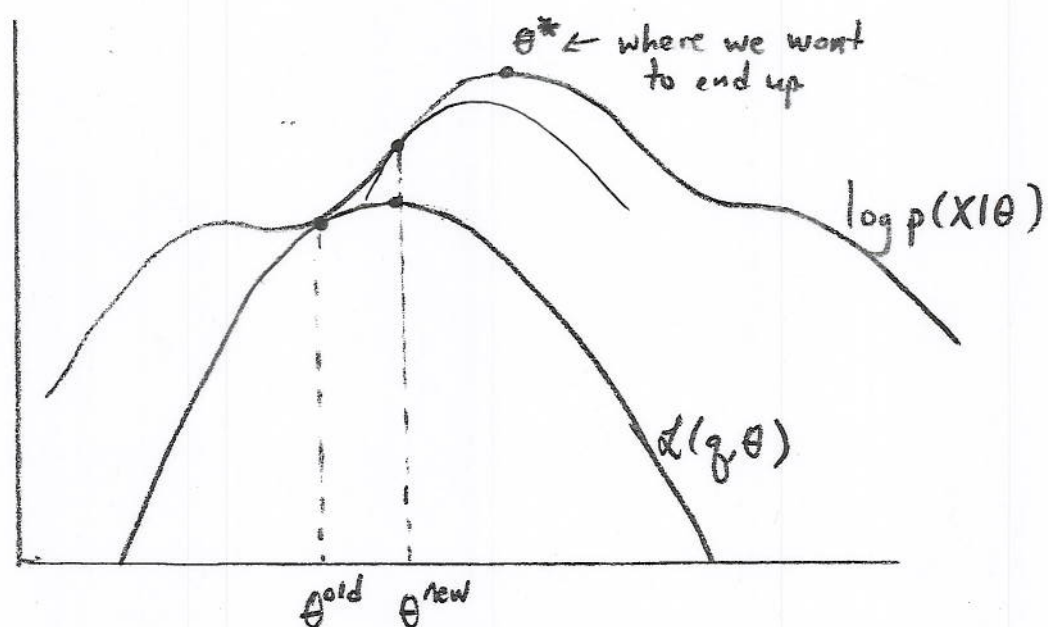


- Because  $\mathcal{L}(q, \theta)$  is lower bound on  $\log p(x|\theta)$ ,  $\log p(x|\theta)$  must increase at least as much as  $\mathcal{L}(q, \theta)$  does.
- Now, there's a nonzero  $KL(q||p)$ .

Based on the pictures on the previous 2 pages, we see that  $\log p(X|\theta)$  is guaranteed to be non-decreasing from one EM iteration to the next.

Thus, the EM algorithm is guaranteed to converge to a local optimum.

The EM algorithm in parameter space:

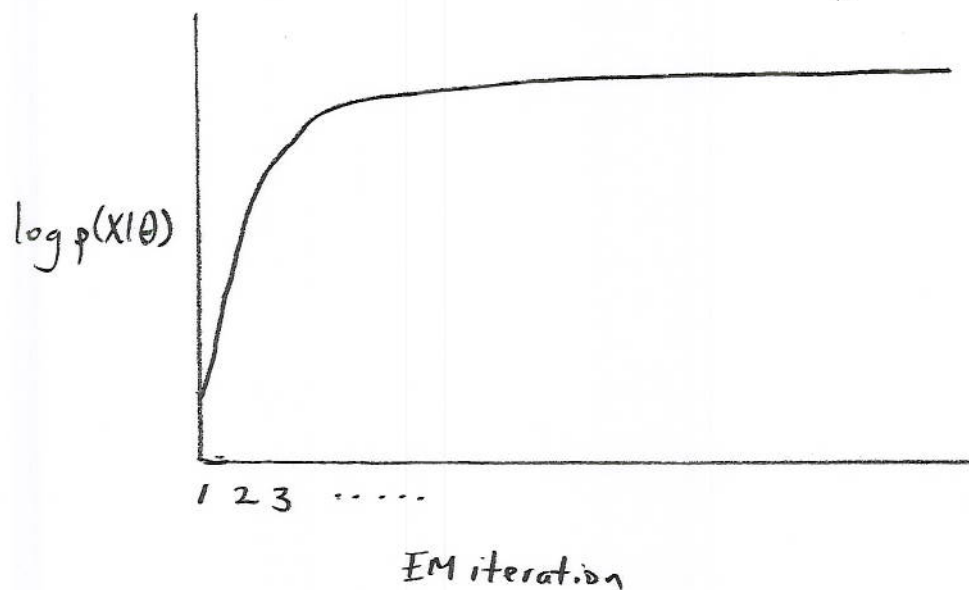


E-step: Make lower bound tight

M-step: Find peak of lower bound



Tracking the progress of the EM algorithm:



- This curve is guaranteed to be non-decreasing
- This is referred to as the "learning curve".

#### D.4) Summary of the general EM algorithm

1) Initialize parameters  $\theta^{old}$ , Compute  $\log p(X|\theta^{old})$ .

2) E-step:

Evaluate  $p(Z|X, \theta^{old})$

3) M-step:

$$\text{Find } \theta^{new} = \underset{\theta}{\operatorname{argmax}} \left( \sum_Z p(Z|X, \theta^{old}) \log p(X, Z|\theta) \right)$$

4) Compute  $\log p(X|\theta^{new})$ . Compare to  $\log p(X|\theta^{old})$ .

5) If not converged, assign

$$\theta^{old} \leftarrow \theta^{new}$$

and go to 2).