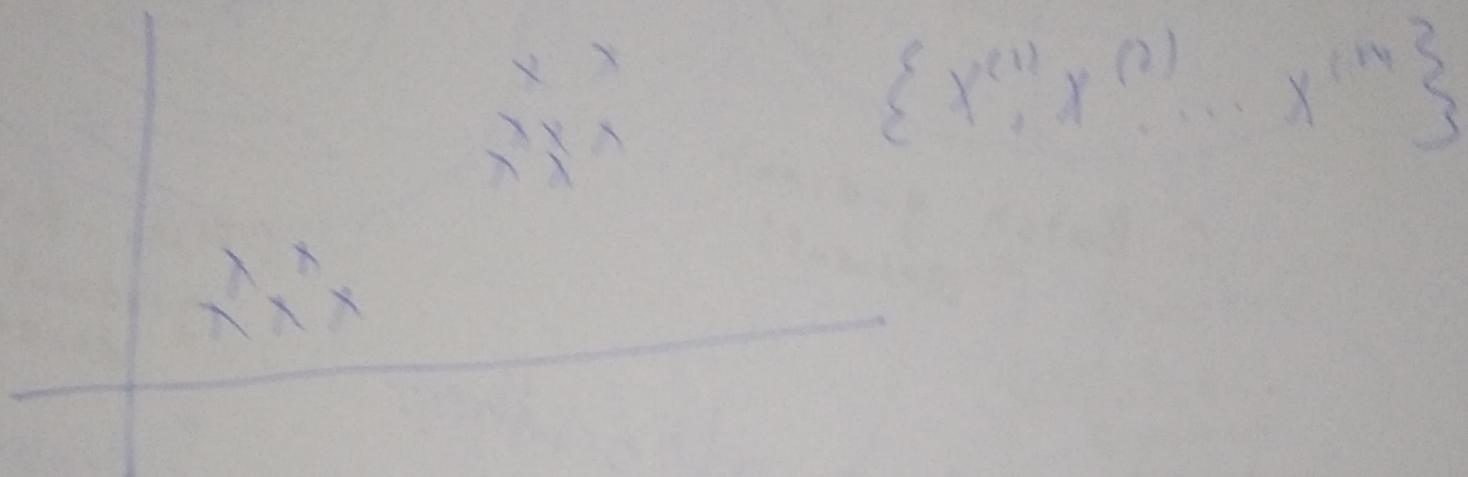


# ④ Unsupervised Learning:-



## K-Means Clustering:-

Data :  $\{x^{(1)}, \dots, x^{(m)}\}$

1. Initialize cluster centroids.

$\mu_1, \dots, \mu_k \in \mathbb{R}^n$   
randomly.

2. Repeat until convergence:

(a) Set  $C^{(t)} = \arg \min_j \|x^{(i)} - \mu_j\|$

(148)  
(149)

1(b): for  $j=1, \dots, k$

(Move  
the cluster  
centroids")

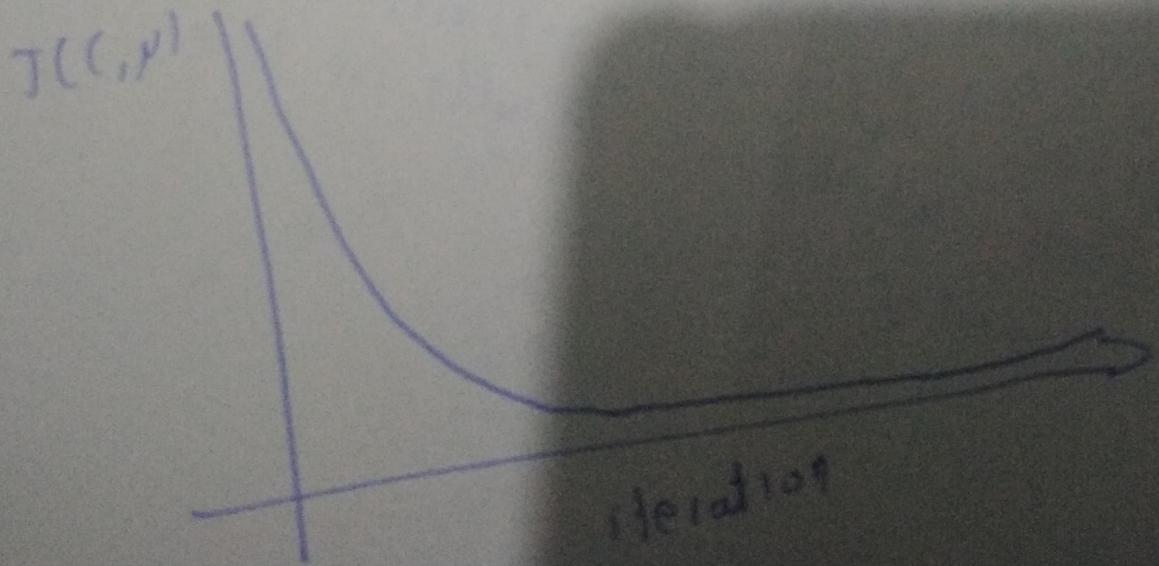
$$N_j = \frac{\sum_{i=1}^m \{c^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m \{c^{(i)} = j\}}$$

$$J(c, \mu) = \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

"Assignments"

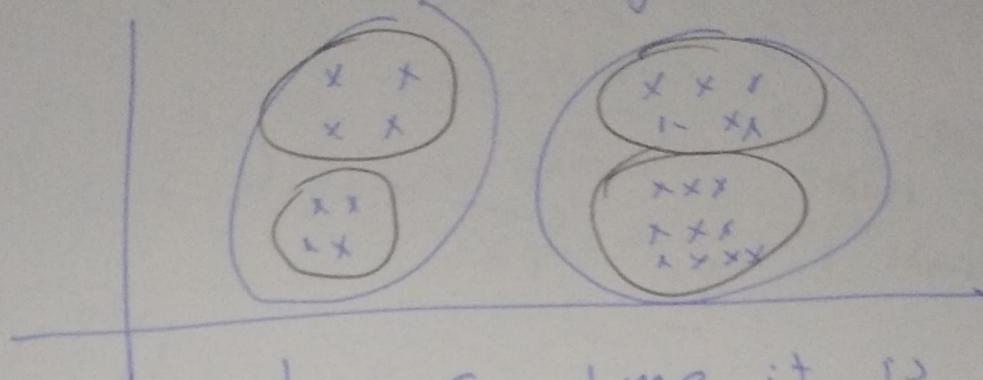
Centroids

On every iteration K-Means would drive down this cost. So it must converge because this function  $\Rightarrow$  can go below zero (it is a strictly non-negative function which goes down after every iteration).



\* How do you choose  $K$ ?

(1) Choose  $K$  by hand C

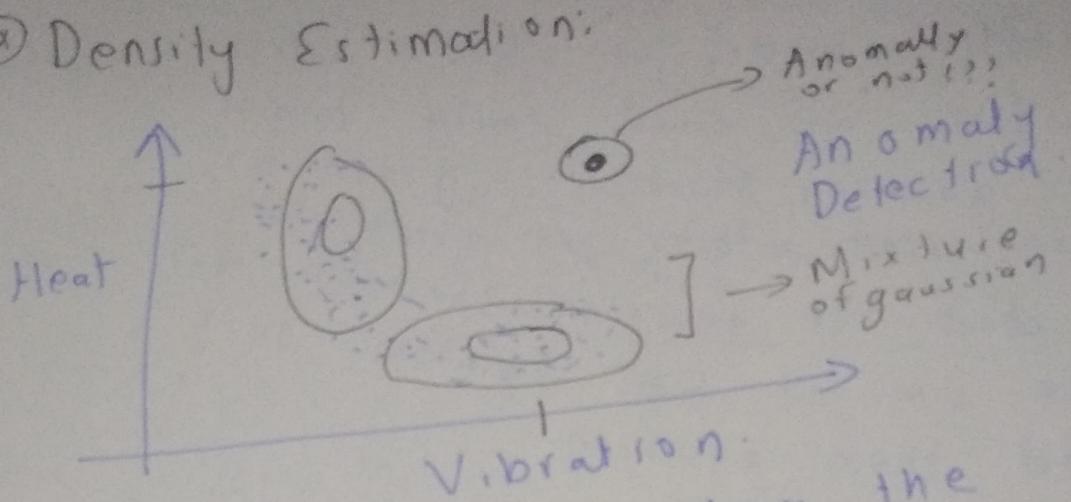


→ Sometime it is  
inherently ambiguous what  
is the right number of clusters.

(2) Criteria like AIC or  
BIC can be used but  
always try to choose  $K$  by  
hand depending on context,  
application & purpose.

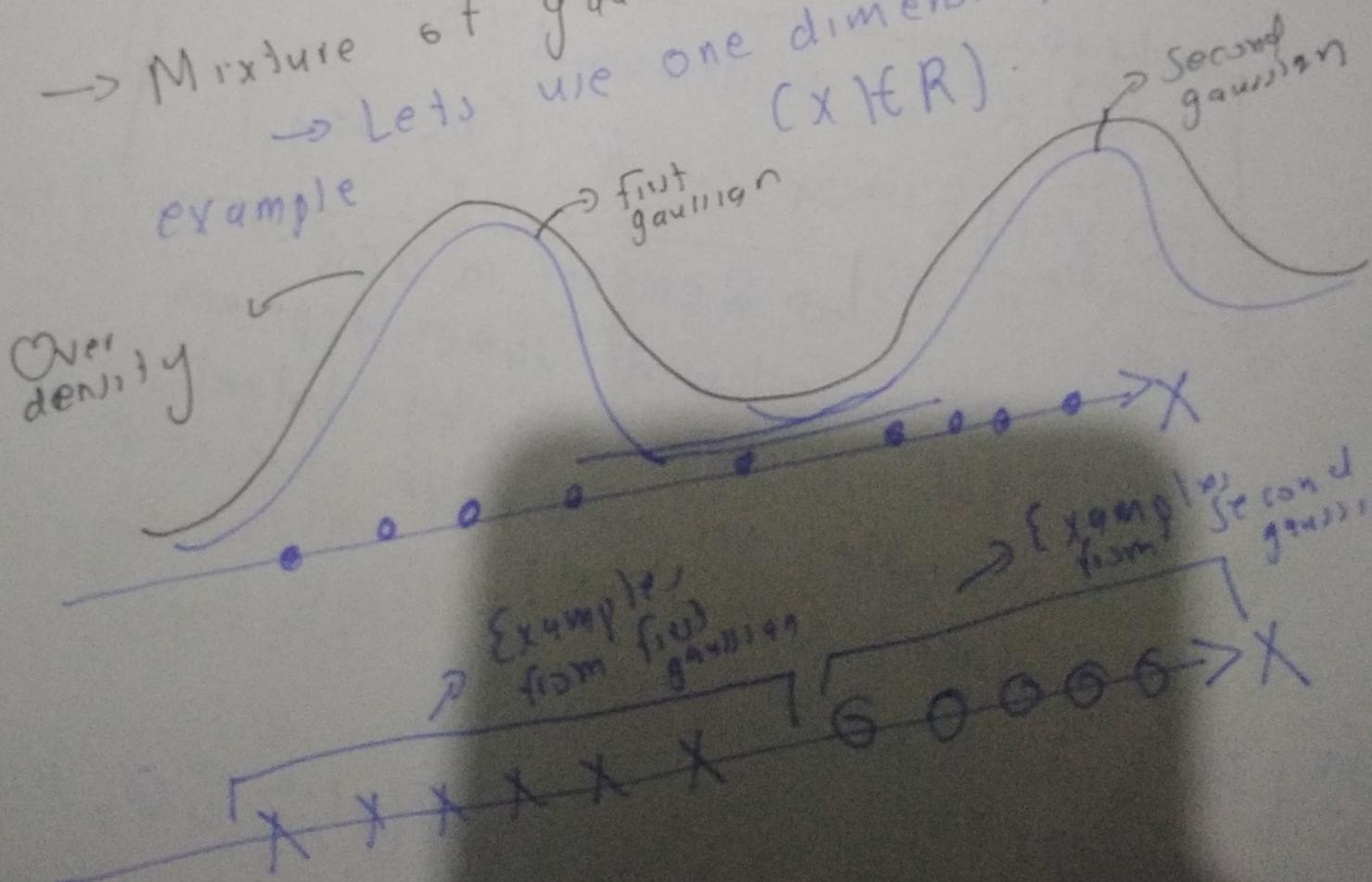
→ K-Means can stuck in local  
minima so run it many times  
from different initialization and  
pick whichever results in the  
lowest cost.

## ④ Density Estimation:



↳ Model  $p(x)$ : (Modelling the density from which  $x$  was drawn so if  $p(x)$  is really small than you draw it as an anomaly)

→ Mixture of gaussian model:  
→ Lets use one dimension ( $x \in \mathbb{R}$ )



(14)

In GDA model we knew labels therefore we could easily model the gaussian. However, given an unsupervised setting we can't be sure about the gaussian therefore will make use of EM algorithm to model these gaussians.

a) Mixture of Gaussian model:-

→ Suppose there is a latent (hidden) random variable  $z$  and  $x^{(i)} | z^{(i)}$  are distributed  $P(x^{(i)} | z^{(i)}) = P(x^{(i)} | z^{(i)}) \cdot p(z^{(i)})$

where

- ①  $z^{(i)} \sim \text{Multinomial}(\phi)$
- ②  $x^{(i)} | z^{(i)}=j \sim \mathcal{N}(\mu_j, \Sigma_j)$

↳ Remember in GDA we followed same approach but instead of using latent variable we used the observed labels. (Important)

So if we know  $z^{(i)}$ , we can use MLE:

$$\ell(\phi, \mu, \Sigma) = \sum_{i=1}^n \log P(x^{(i)} | z^{(i)}; \phi, \mu, \Sigma)$$

(48)

① We would take the  $\phi$

$\nabla_{\theta, \gamma} \ell(\theta, \gamma) = 0$  and would  
find:-

$$\phi_j = \frac{\sum_{i=1}^m \mathbb{1}\{z^{(i)} = j\}}{m}$$

$$N_j = \frac{\sum_{i=1}^m \mathbb{1}\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m \mathbb{1}\{z^{(i)} = j\}}$$

$$\Sigma_{ij} = \frac{\sum_{i=1}^m (x - N_j)^T (x - N_j)}{\sum_{i=1}^m \mathbb{1}\{z^{(i)} = j\}}$$

But we don't know  $z^{(i)}$ 's.

→ Therefore we will use  
EM (Expectation minimization algorithm)

EM (Expectation Minimization Algorithm):-

→ b) Expectation Minimization Algorithm:  
(called Bootstrap procedure):

→ EM algorithm has two  
main steps:-

for  
e

1) E-step (Guess value of  $z^{(1),j}$ )

(149)

$\rightarrow$  Set  $w_j^{(1)} = P(z^{(1)}=j | x^{(1)}; \theta, n, \Sigma)$

$$= \frac{P(x^{(1)} | z^{(1)}=j) \cdot P(z^{(1)}=j)}{\sum_{l=1}^k P(x^{(1)} | z^{(1)}=l) \cdot P(z^{(1)}=l)}$$

$$\Rightarrow \mathcal{N}(N_j, \Sigma_j) = \frac{1}{(2\pi)^{n/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2} (x^{(1)} - N_j)^T \Sigma_j^{-1} (x^{(1)} - N_j)\right)$$

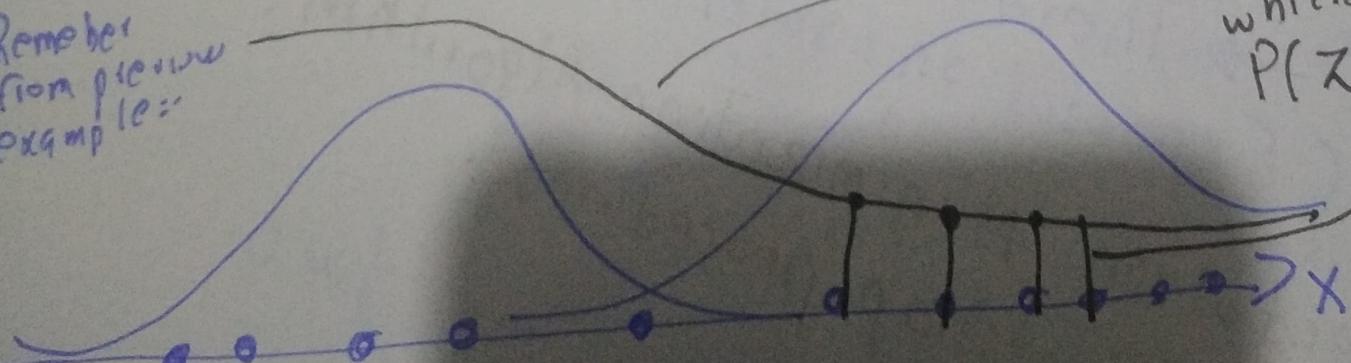
②  $\sim$  Multinomial distribution  $= \phi_j$

( $Z \sim$  Multinomially distributed).

$\rightarrow$  The terms in denominator would follow the same same distribution for different  $l$ 's.

Remember from previous example:

Sigmoid Curve which shows  $P(z^{(1)}=j | x^{(1)})$



→ You could compute  $w_{ij}$  (150)  
 for every training example with  
 every  $j$ .

2) M-Step:

$$\phi_j := \frac{1}{m} \sum_{i=1}^m w_{ij}^{(1)}$$

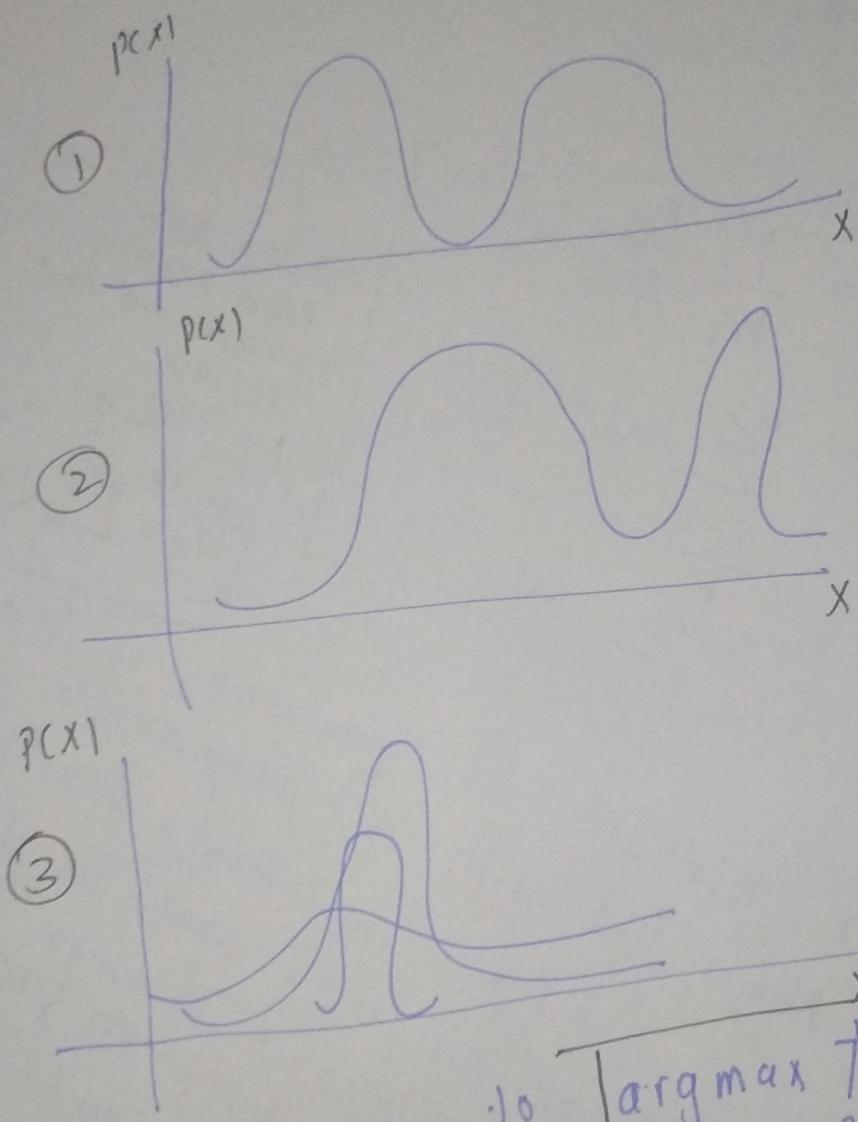
$$\mu_j = \frac{\sum_{i=1}^m w_{ij}^{(1)} x^{(1)}}{\sum_{i=1}^m w_{ij}^{(1)}}$$

→ See how  
 from HDA  
 equation we have  
 replace  
 $\sum z^{(1)} = j$   
 with  
 $w_j = E[\alpha z^{(1)} = j]$

→ this mixture of gaussian is like  
 K-Means but with soft-assignments.  
 ↳ EM implements a softer  
 way of assigning points to different  
 (centroids) because it assign a  
 specific probability for  $x^{(1)}$  for  
 each  $j$  (possible centers).

→ the algorithm Mixture of gaussian  
 will converge but with some caveats.

→ Mixture of gaussian can fit different and various kind of models:

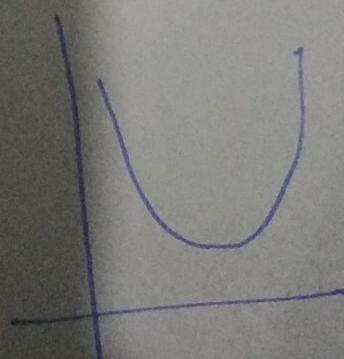


→ EM is trying to  $\underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m p(x^{(i)}; \theta)$

→ Jensen's Inequality:-

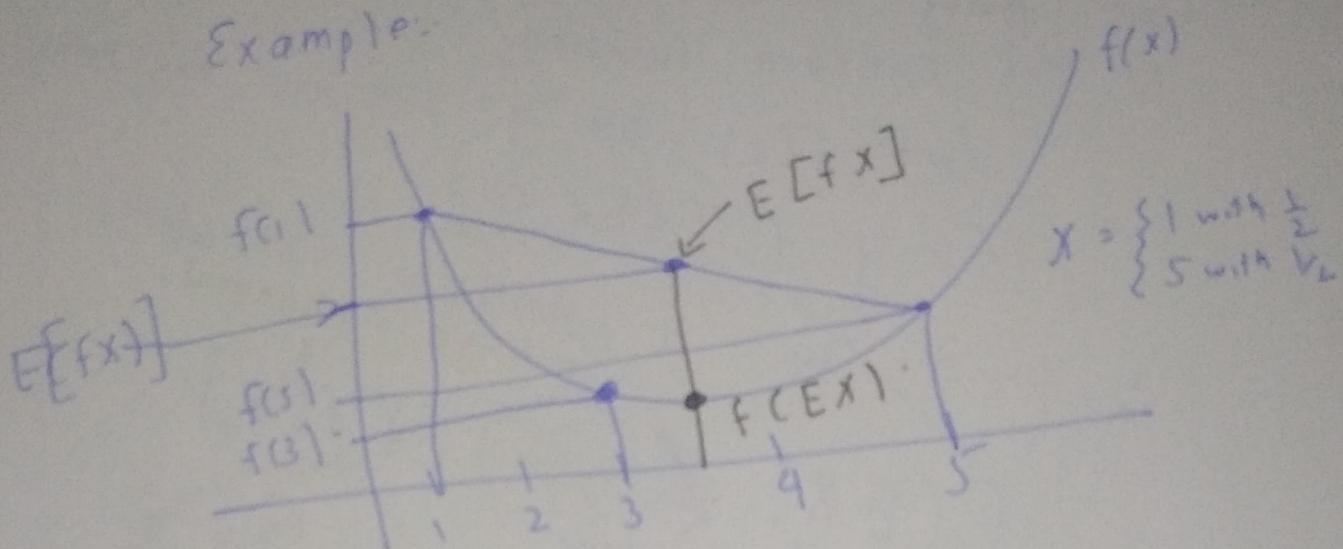
⑥ let  $f$  be a convex function  
 (e.g.  $f''(x) > 0$ )

→ Let  $X$  be a random variable



Then  $f(E\bar{x}) \leq E[f(\bar{x})]$

Example:



$$f(3) = f(E\bar{x})$$

$$E[f(\bar{x})] = \frac{1}{2}f(1) + \frac{1}{2}f(5)$$

then

$$f(E\bar{x}) \leq E[f(\bar{x})]$$

Further, if  $f''(x) > 0$  ( $f$  is strictly convex)

then:

$$E[f(\bar{x})] = f(E\bar{x}) \Leftrightarrow X \text{ is a constant}$$

$X = E[\bar{x}]$  with probability 1

→ Form of Jensen inequality we are going to use for EM algorithm is for concave function:-

Let  $f(x)$  be a concave function:-

(e.g.  $f''(x) < 0$ )

Let  $X$  be a random variable:-  
→ sign is changed.

$$f(\mathbb{E}X) \cancel{\geq} \mathbb{E}[f(X)]$$

(Round the m.)

→ Density Estimation problem:-

→ Have model for  $P(x, z; \theta)$ .

→ You only observe  $x$

$$\{x^{(1)}, \dots, x^{(m)}\}$$

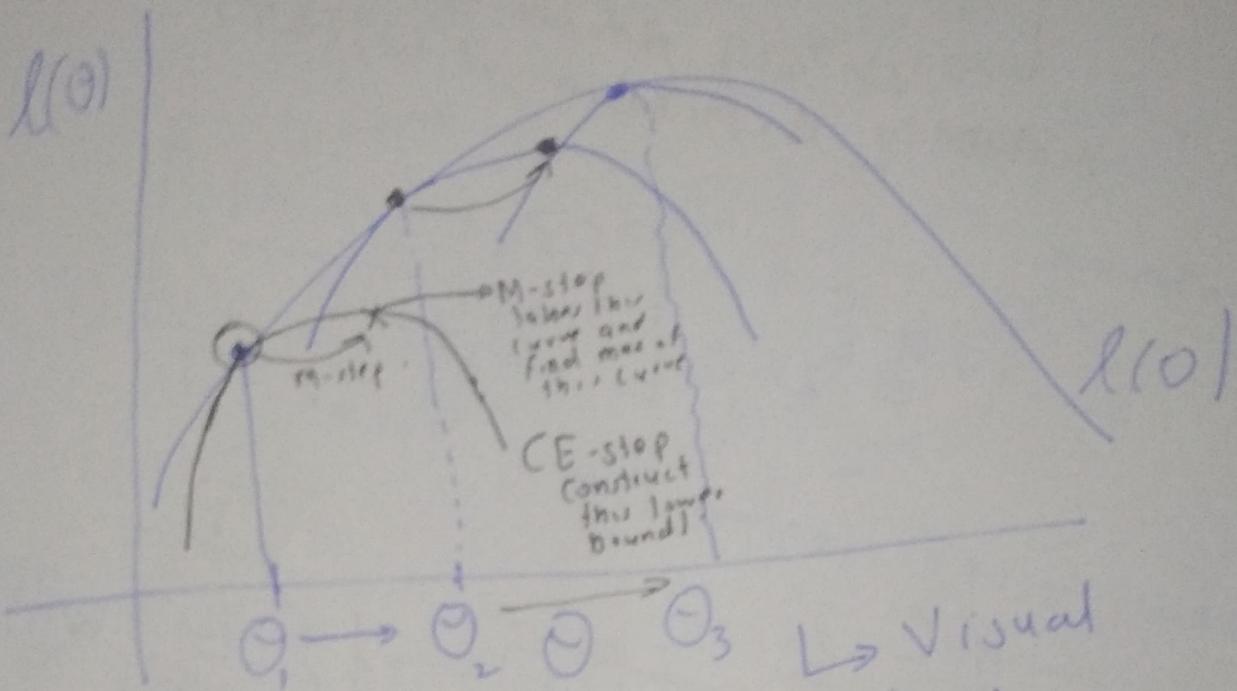
$$R(\theta) = \sum_{i=1}^m \log(P(x^{(i)}, z^{(i)}; \theta))$$

$$l(\theta) = \sum_{i=1}^m \log P(x^{(i)}; \theta)$$

$$l(\theta) = \sum_{i=1}^m \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)$$

④

Want: argmax <sub>$\theta$</sub>   $\ell(\theta)$



Representation of E-step which is  
construct the lower-bound curve and  
M-step taking the max of that  
curve.

$$\textcircled{4} \quad \max_{\theta} \sum_i \log p(x^{(i)}; \theta)$$

$$= \sum_i \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta)$$

$$= \sum_i \log \sum_{z^{(i)}} Q_i(z^{(i)}) \left[ \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

155

= where  $\boxed{Q_i(z^{(i)})}$  is a probability distribution i.e.  $\boxed{\sum_{z^{(i)}} Q_i(z^{(i)}) = 1}$ .

↳ We will decide later the probability distribution of  $Q_i(z^{(i)})$ .

Continuing:

$$= \sum_i \log E_{z^{(i)} \sim Q_i} \left[ \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

Now using concave form of Jensen's inequality:

(2)

$\downarrow f \circ x$

$$\sum_i \log E_{z^{(i)} \sim Q_i} \left[ \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] \stackrel{(2)}{\geq} \sum_i E_{z^{(i)} \sim Q_i} \left[ \log \left( \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right) \right]$$

$$f((E(x)))$$

where  $f(x) = \log(x)$

$f(x)$

$E(f x)$

$\log$

$\curvearrowright$

$$\textcircled{1} = \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \left( \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right)$$

(Unpacking the terms).

\* The above function  $\textcircled{1}$  is just a function of  $\theta$  where:

→ Function  $\textcircled{1}$  is lower bound (See the last graph) for log likelihood.

→ Function  $\textcircled{2}$  is log likelihood.

\* We have  $\textcircled{2} > \textcircled{1}$  but we want  $\textcircled{2} = \textcircled{1}$  as shown in graph for current value of theta.

On a given iteration of EM w.r.t parameters  $\theta$  we want want:

$$\log E_{z^{(i)} \sim Q_i} \left[ \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] = E_{z^{(i)} \sim Q_i} \left[ \frac{\log P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

For this to be true (remember Jensen inequality) we want:

$$\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} = \text{constant.}$$

(Ratio of denominator and numerator must be same).

Therefore Set  $Q_i(z^{(i)}) \propto P(x^{(i)}, z^{(i)}; \theta)$

$$\text{To ensure proportionality}$$

$$Q_i(z^{(i)}) = P(x^{(i)}, z^{(i)})$$

$$Q_i(z^{(i)}) = 1$$

$$\sum_{z^{(i)}} Q_i(z^{(i)}) = \frac{P(x^{(i)}, z^{(i)}; \theta)}{\sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)}$$

Couple More steps

$$Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

$$Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

↳ This will ensure the equality of expression ① & ②.

→ Giving a summary of EM algorithm:-

**E-step:**

Set

$$Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

**M-step**

$$\theta := \arg \max \sum_i P(x^{(i)})$$

$$= \arg \max_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$