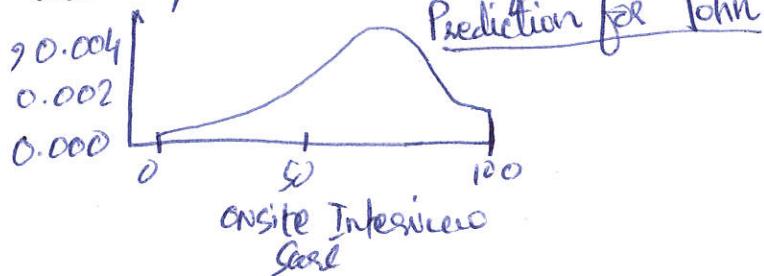


- \* Latent variable is neither observed during training nor testing. We will use EM algorithm that are used to train latent variable models.

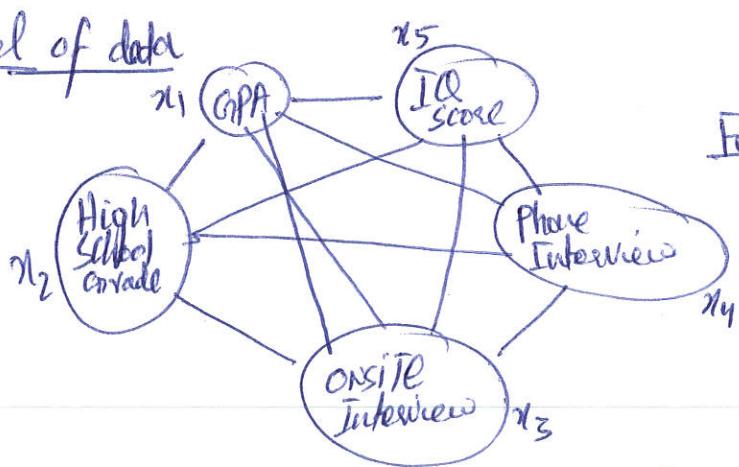
Why do we need probabilistic modelling of DATA?

	High School	GPA	IQ	Phone Interview	On Site Performance
John	4.0	4.0	120	30/4	???
Helen	3.7	3.6	N/A	4/4	???

We would like to predict onsite Performance. We can use historical data and do regression. But here we have missing variables as well. So, probabilistic modelling will help us deal with missing values & it can also provide us confidence over prediction. For instance,



Probabilistic Model of data



Fully Connected model

The joint distribution table for above model will be exponentially large. Therefore it is impractical to treat them as parameters.

High School	GPA	IQ	...	Probability
Exponential S				

## Option #2

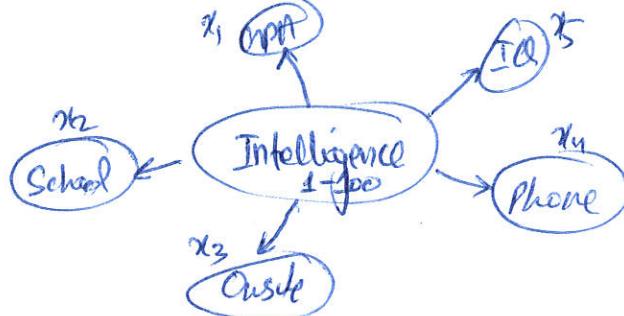
We can ~~but~~ assume parametric model.

$$P(x_1, x_2, x_3, x_4, x_5) = \frac{\exp(\omega^T x)}{Z}$$

The problem is that  $Z$  is a sum over exponential terms

## Option #3

We can introduce latent variable "Intelligence" in the model. It will simplify our model a lot.



$$\begin{aligned}
 P(x_1, x_2, x_3, x_4, x_5) &= \sum_{I=1}^{100} P(x_1, x_2, x_3, x_4, x_5, I) \quad (\text{Sum rule}) \\
 &= \sum_{I=1}^{100} \underbrace{P(x_1, x_2, \dots, x_5 | I)}_{\substack{\text{Conditional} \\ \text{Prob}}} \underbrace{P(I)}_{\substack{\text{Prior} \\ \text{Prob}}}
 \end{aligned}$$

This conditional prob will factorise into small prob because of the structure of model

$$= \sum_{I=1}^{100} P(x_1 | I) P(x_2 | I) \dots P(x_5 | I) P(I)$$

From huge joint prob table, we have much simpler table to work with.

## Summary of Latent Variable

### Pros:

- i) Introducing latent variables can simplify our model
- ii) They can be sometimes meaningful

### Cons:

- i) Hard to work with, require lots of maths

Probabilistic Clustering

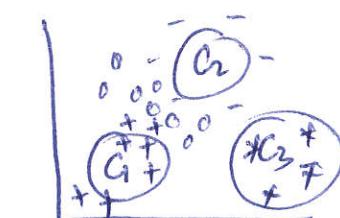
In this lecture, we will use latent variable model for clustering.

Soft Clustering

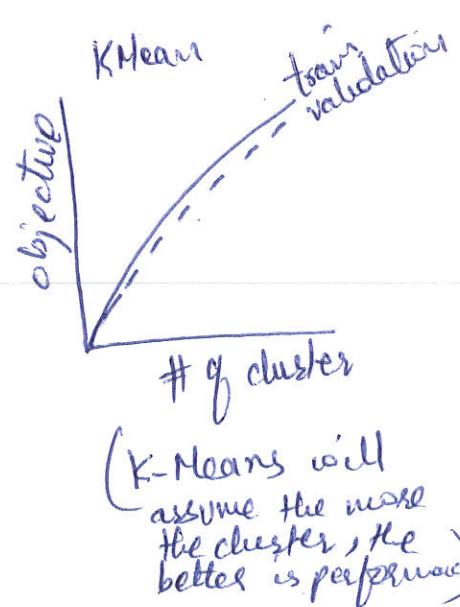
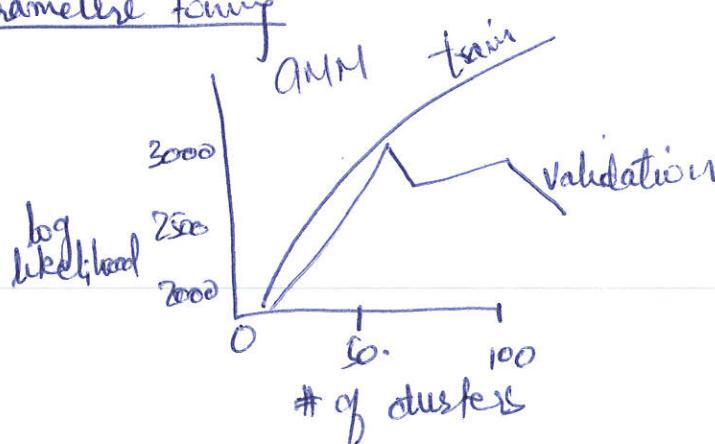
- \* A point can be assigned to multiple clusters using probabilities.

- \* Hard clustering  
A point is assigned single cluster  
 $\text{cluster-index} = f(x)$

- \* Why do we need to perform soft clustering probabilistically?
  - i) Because we want to deal with missing values/features
  - ii) We want to tune hyperparameter

Soft Clustering

Each point "x"  $\in \{p(c_1|x), p(c_2|x), p(c_3|x)\}$   
 i.e.  $\rightarrow p(\text{cluster index} | x)$

Hyperparameters tuning

Another reason we would like to train clustering probabilistically is the we may want to get generative model of our data. It can generate data (like images).

## Summary

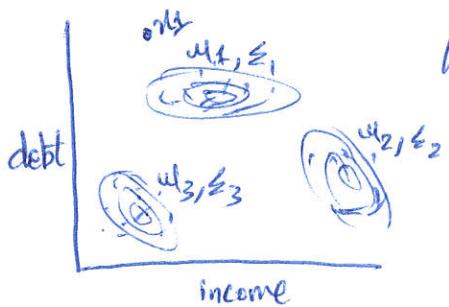
Soft clustering allows us

- i) to ~~sample~~ generate data (by getting generative model of data)
- ii) to tune hyperparameters

## W2-PI-L3

### GMM

We will build latent variable model for clustering probabilistically



In the above problem, we will use three Gaussians. Therefore we can assume that each datapoint comes from a weighted sum of each of these three gaussians:

$$p(x|\theta) = \pi_1 N(x|\mu_1, \Sigma_1) + \pi_2 N(x|\mu_2, \Sigma_2) + \pi_3 N(x|\mu_3, \Sigma_3)$$

i.e. the density of each data point equals to the sum of three gaussian distributions weighted

$$\pi_1 + \pi_2 + \pi_3 = 1$$

$$\Theta = \{\pi_1, \pi_2, \pi_3, \mu_1, \mu_2, \mu_3, \Sigma_1, \Sigma_2, \Sigma_3\}$$

This is called Gaussian Mixture model, where # of mixture = 3

(23)

How to find the parameters of GMM?

We can apply maximum likelihood estimation here

$$\max_{\theta} \underbrace{P(\mathbf{x}|\theta)}_{\text{Density of dataset given } \theta} = \prod_{i=1}^N P(x_i|\theta)$$

$$\text{dataset given } \theta = \prod_{i=1}^N \gamma_k N(x_i|\mu_k, \Sigma_k) + \gamma_1 N(x_i|\mu_1, \Sigma_1) + \gamma_2 N(x_i|\mu_2, \Sigma_2) + \gamma_3 N(x_i|\mu_3, \Sigma_3)$$

s.t. i)  $\gamma_1 + \gamma_2 + \gamma_3 = 1$  (this will ensure valid Prob. distribution)

ii)  $\gamma_k > 0; k=1,2,3$

and iii)  $\Sigma_k \geq 0 (k=1,2,3)$  PSD-matrix

Thus GMM is a flexible Probability distribution model.

from last figure

$$P(x_i|\theta) = \gamma_1 \underbrace{N(x_i|\mu_1, \Sigma_1)}_{\rightarrow \text{height}} + \gamma_2 \underbrace{N(x_i|\mu_2, \Sigma_2)}_{\rightarrow \text{low } \gamma_2} + \gamma_3 \underbrace{N(x_i|\mu_3, \Sigma_3)}_{\rightarrow \text{least } \gamma_3}$$

for  $x_1$  case:

$\rightarrow$  potentially  $\gamma_1$  is highest here

$\downarrow$   
 $\rightarrow$  low  $\gamma_2$

$\downarrow$   
 $\rightarrow$  least  
 $\rightarrow$  least  $\gamma_3$



How do we train GMM using a fact that we would like to avoid using stochastic gradient descent. We will introduce latent-variable that will help us using EM.

From last lecture, density of datapoint " $x$ " was given as follows:

$$\text{Eq A} \quad P(x|\theta) = \pi_1 N(x|\mu_1, \Sigma_1) + \pi_2 N(x|\mu_2, \Sigma_2) + \pi_3 N(x|\mu_3, \Sigma_3)$$

lets assume there is a latent-variable " $t$ " that causes " $x$ "

$$t \rightarrow x$$

In our case, this latent variable will be the cluster index. We will not observe latent variable during testing and training. Given this model, it is safe to assume that prior distribution on " $t$ " is  $\pi$  i.e:  $P(t=c|\theta) = \pi_c$  where  $c=1, 2, 3$

And  $P(x|t=c, \theta) = N(x|\mu_c, \Sigma_c)$  If we know from which cluster " $x$ " is coming from

Now, using sum rule (marginalize " $t$ "), we would like to get  $P(x|\theta)$  — the likelihood

$$P(x|\theta) = P(x|t=1, \theta) + P(x|t=2, \theta) + P(x|t=3, \theta)$$

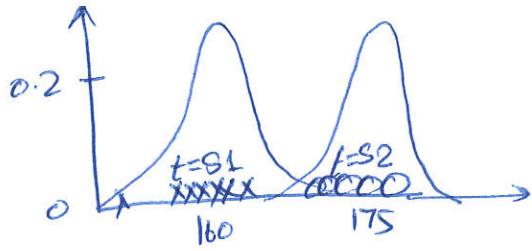
$$\text{Eq B} \quad P(x|\theta) = \sum_{c=1}^3 P(x|t=c, \theta) P(t=c|\theta) \quad \boxed{P(x|\theta) = P(x, t=1|\theta) + P(x, t=2|\theta) + P(x, t=3|\theta)}$$

$$\text{i.e. } P(x|\theta) = P(x, t=1|\theta) + P(x, t=2|\theta) + P(x, t=3|\theta)$$

$$\Rightarrow P(x|\theta) = \pi_1 N(x|\mu_1, \Sigma_1) + \pi_2 N(x|\mu_2, \Sigma_2) + \pi_3 N(x|\mu_3, \Sigma_3)$$

So, we have seen that EqA is exactly similar to EqB. Therefore we can proceed forward with this latent variable model in EqB. It will give us the same result as of EqA.

## Expectation Maximization



How to compute parameters " $\theta$ "?

Solution-1 (sources are known) → Hard assignment  
If we know the sources, i.e. we know from which Gaussian data " $x_i$ " is coming from  $s_i$  can easily evaluate parameters as follows:

$$p(x | t=S_1, \theta) = N(x | \mu_1, \sigma_1^2) \rightarrow \text{Data given source } S_1$$

$$\mu_1 = \frac{\sum_{S_1} x_{S_1}}{\# \text{ of } S_1 \text{ points}} \quad \sigma_1^2 = \frac{\sum_{S_1} (x_{S_1} - \mu_{S_1})^2}{\# \text{ of } S_1 \text{ points}}$$

This is the example of hard assignment. And if we have soft assignments i.e. some posterior distributions on " $t$ " (In this case any point can belong to all clusters but with some probability)

Known:  $p(t_i=1 | x_i, \theta) \& p(t_i=2 | x_i, \theta) \rightarrow \text{Posterior for } "t"$

then we can estimate the ~~parameters~~ as follows:

$$\text{Eq(C)} \quad \text{Then } \mu_1 = \frac{\sum_i p(t_i=1 | x_i, \theta) x_i}{\sum_i p(t_i=1 | x_i, \theta)} \quad \sigma_1^2 = \frac{\sum_i p(t_i=1 | x_i, \theta) (x_i - \mu_1)^2}{\sum_i p(t_i=1 | x_i, \theta)}$$

## Solution-2

If we know the parameters " $\theta$ ", then we can estimate the ~~Gaussian~~ sources using Bayes Rule.

i.e. Given:  $p(x | t=1, \theta) = N(-2, 1)$  i.e.  $\mu_1 = -2, \sigma_1^2 = 1$  for  $t=1$

then using Bayes Rule:

$p(t=1 | x, \theta)$  can be imputed

From Bayes Rule, the soft assignment to data " $x$ " can be imputed as follows:

$$p(t=1 | x, \theta) = \frac{p(x | t=1, \theta)}{Z} \leftarrow \begin{array}{l} \text{Joint} \\ x \rightarrow \text{normalizing constant} \end{array}$$

$$\text{posterior probability} = p(x | t=1, \theta) / Z = \frac{\text{likelihood}}{\pi}$$

$$X = P(t=1, \theta) + p(t=2, \theta) \rightarrow \text{just calculate}$$

(27)

### EM Algorithm :

- 1) Need Gaussian param  $(\mu, \sigma^2)$  to estimate source  $p(t|x, \theta)$
- 2) Need sources  $p(t|x, \theta)$  to estimate gaussian parameters using Eq(C).

### EM Algo

- 1) Initial params  $(\mu_1, \sigma_1^2)$  randomly &  $(\mu_2, \sigma_2^2)$

- 2) Until Convergence:

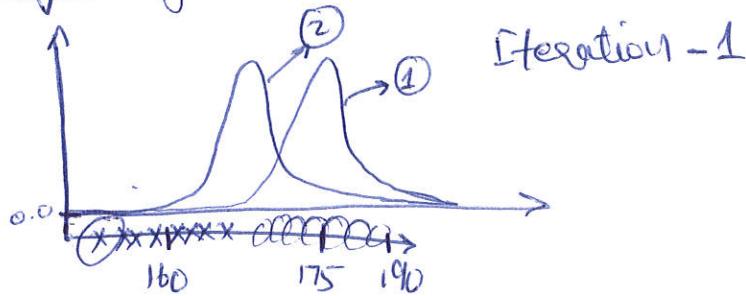
a)  $\cancel{P(t=c|x_i, \theta)}$   
 $P(t=1|x_i, \theta) = ??$  } for all data points "i"  
 $P(t=2|x_i, \theta) = ??$

b) Using  $p(t|x_i, \theta)$   
 Update " $\theta$ " using Eq(C)



## Example GMM-training using EM

(29)



Step-1

- i) Initialize " $\theta$ " for  $K=2$ . Now parameters are known so compute:

$$p(t_i=1 | x_i, \theta) \approx p(t_i=2 | x_i, \theta) \text{ for all points}$$

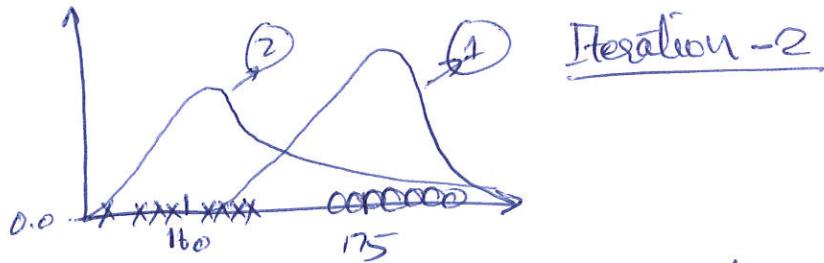
$$p(t_i=1 | x_i, \theta) \approx \frac{0.01 \leftarrow \text{density of } x_i \text{ using } \theta_1}{0.01 + 0.0002 \leftarrow \text{sum of densities}}$$

Step-2

- i) Once sources are known, re-estimate params using Bayes Rule using formula

$$\mu_1 = \text{---} \quad \sum \theta_1 = \text{---}$$

$$\mu_2 = \text{---} \quad \sum \theta_2 = \text{---}$$



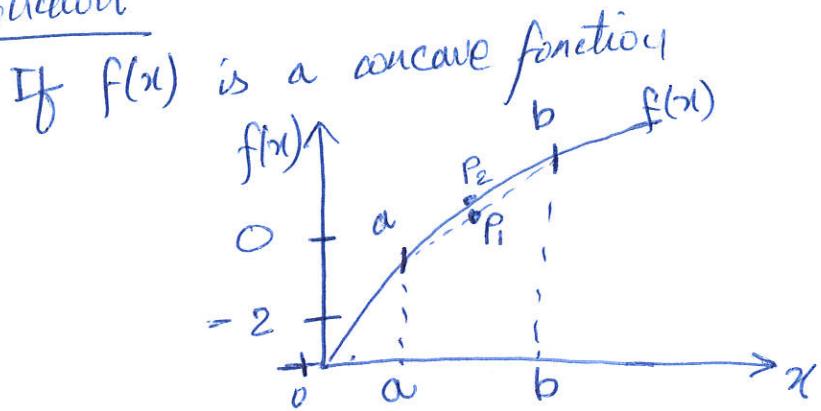
- \* GMM is probabilistic approach for soft clustering. You can also sample new data points.
- \* EM can replace Stochastic Gradient.
- \* EM gives local maxima (the exact solution is NP-hard).
- \* for this clustering problem.

## General form of EM

We will derive general form of EM-Algorithm that will help us train almost any latent variable model we can think of. But we will see some mathematical tools first:

- 1) Concave functions
- 2) Jensen Inequality

## Concave function



Then for any  $a, b \in \mathbb{R}$ :

$$f(\alpha a + (1-\alpha)b) \geq \alpha f(a) + (1-\alpha)f(b)$$

and  $0 \leq \alpha \leq 1$

i.e.  $\text{point}(P_\alpha) > \text{point}(P_2)$

whenever  $0 \leq \alpha \leq 1$  for two points "a" & "b" on concave function  $f(x)$

and for any number of points on concave functions (like logarithmic function), you can prove the same property:  
For e.g: For any three points  $a_1, a_2, a_3$  with weights  $\alpha_1, \alpha_2$  and  $\alpha_3$   
we can prove  $f(\alpha_1 a_1 + \alpha_2 a_2 + \alpha_3 a_3) \geq \alpha_1 f(a_1) + \alpha_2 f(a_2) + \alpha_3 f(a_3)$   
st  $\alpha_1 + \alpha_2 + \alpha_3 = 1$

We can think of above property in probabilistic term: first we can think of  $\alpha$ 's as probabilities ( $\alpha$ 's are weights for each point). Then we can assume latent variable "t" with the following probability distribution

$$P(t=a_1) = \alpha_1$$

$$P(t=a_2) = \alpha_2$$

$$P(t=a_3) = \alpha_3$$

Now we can rewrite Jensen's Inequality as follows:

$$\underbrace{f(\alpha_1 a_1 + \alpha_2 a_2 + \alpha_3 a_3)}_{E_p(t)t} \geq \underbrace{\alpha_1 f(a_1) + \alpha_2 f(a_2) + \alpha_3 f(a_3)}_{E_p(t)f(t)}$$

i.e.  $E_p(t=1)f(t=1) = \underbrace{\alpha_1 f(a_1)}_{\text{Expected value of } f(t=1)}$

This sums up Jensen's Inequality:

For any concave function  $f$ , and for any probability distribution " $p(t)$ " on " $\mathbb{R}$ :

$$f(E_p(t)t) \geq \underbrace{E_p(t)f(t)}_{\substack{\text{"f" of expected} \\ \text{value } t}} \quad \underbrace{\text{Expected value}}_{\text{of } f(t)}$$

This is true, for any number of points "t". For infinite number of points "t", this summation will turn into integral as follows:

$$f(E_p(t)t) \geq \int E_p(t)f(t) dt \quad (\text{check it})$$

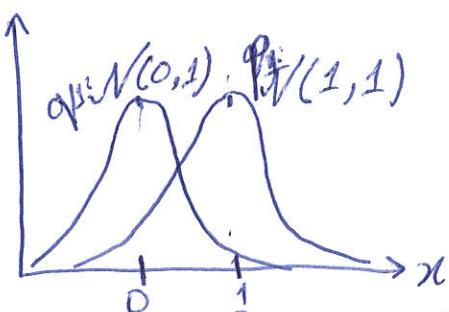
Final thing we need to know is KL-Divergence.

### Kullback-Leibler divergence

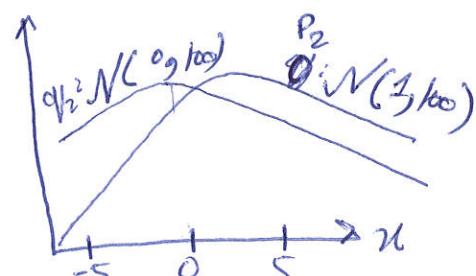
It calculates the difference b/w two probabilistic dist.

### Motivation

Consider following Gaussians:



- i) Difference in parameters = 1
- ii)  $KL(q_1 || p_1) = 0.5$



- i) Difference in parameters = 1
- ii)  $KL(q_2 || p_2) = 0.005$

- \* Intuitively, Gaussians on right picture are more similar which is captured by KL-Divergence but not by difference in parameters.

$$K_h(q||p) = \int q(x) \log \frac{q(x)}{p(x)} dx = \int \mathbb{E}_q [\log(\frac{q}{p})] dx$$

Expected value of logarithm of ratio.

$$\log \frac{q(x)}{p(x)} = \log q(x) - \log p(x) \rightarrow \text{measures how different these distributions are at current point } "x"\text{ in the log-scale}$$

Illustration

$$\int q(x) \log \frac{q(x)}{p(x)} dx \rightarrow \begin{array}{l} x_1 = -99.99 \dots \log q(x) - \log p(x) \rightarrow xq(x) \\ x_2 = -99.98 \dots \log q(x) - \log p(x) \rightarrow xq(x) \\ \vdots \quad \vdots \\ x_n = 99.99 \dots \log q(x) - \log p(x) \rightarrow xq(x) \end{array} \rightarrow \sum xq(x)$$

i.e. we are taking differences of  $q(x)$  &  $p(x)$  on entire range, and multiplying these differences by  $q(x)$

Summary

$$K_h(q||p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

$$1. K_h(q||p) \neq KL(p||q)$$

i.e. it's not a proper distance measure in strict mathematical sense

$$2. K_h(q||q) = 0$$

$$\begin{aligned} i.e. K_h(q||q) &= \int q(x) \log \frac{q(x)}{q(x)} dx \\ &= \int q(x) \log(1) dx \\ &= 0 \end{aligned}$$

$$3. K_h(q||p) \geq 0$$

To prove this: we should consider

$$-KL(q||p)$$

$$As \quad K_h(q||p) = \mathbb{E}_q (\log \frac{q}{p})$$

$$\Rightarrow -KL(q||p) = \mathbb{E}_q (-\log(\frac{q}{p})) \\ = \mathbb{E}_q (\log \frac{p}{q})$$

Using Jensen's Inequality:

$$\mathbb{E}_q (\log \frac{p}{q}) \leq \log (\mathbb{E}_q \frac{p}{q}) \text{ since } \log \text{ is concave}$$

So, we have  $\log(E_q(\frac{P}{Q})) = \boxed{0}$

(33)

$$E_q\left(\frac{P}{Q}\right) = \int q(x) \frac{P(x)}{q(x)} dx$$

$$\begin{aligned} \Rightarrow \log E_q\left(\frac{P}{Q}\right) &= \log \int q(x) \frac{P(x)}{q(x)} dx \\ &= \log \int p(x) dx^1 \\ &= \log(1) \\ &= 0 \end{aligned}$$

$$\Rightarrow \boxed{\log E_q\left(\frac{P}{Q}\right) = 0} \rightarrow \text{eq A}$$

This means: ~~ABD(Q||P) = 0~~

$$\begin{aligned} -KL(Q||P) &= E_Q(-\log \frac{Q}{P}) \\ &= E_Q(\log \frac{P}{Q}) \end{aligned}$$

Using Jensen's Inequality:

$$E_Q(\log \frac{P}{Q}) \leq \log(E_Q \frac{P}{Q})$$

And using (A)

$$E_Q[\log \frac{P}{Q}] \leq 0$$

i.e.  $-KL(Q||P)$  is always non-positive  
which implies  $KL(Q||P) \geq 0$ .  $\square$

## GENERAL form of Expectation Maximization

### Motivation

lets assume we have latent variable "t"  
for each "x":

$$(t_i) \rightarrow (x_i)$$

So, the marginal of "x<sub>i</sub>" is given as:

$$p(x_i|\theta) = \sum p(x_i, t_i|\theta) \rightarrow \text{Joint over } x_i, t_i$$

If  $t_i \in \{1, 2, 3\}$

$$\begin{aligned} &= \sum_{c=1}^3 p(x_i, t_i=c|\theta) \rightarrow \text{Marginalize "t<sub>i</sub>" using sum-rule} \\ &= \underbrace{\sum_{c=1}^3 p(x_i|t_i=c, \theta)}_{\text{likelihood}} \underbrace{p(t_i=c|\theta)}_{\text{prior}} \end{aligned}$$

Our problem here is:

$$\max_{\theta} p(x|\theta) \rightarrow \text{This is marginal likelihood as we have marginalized latent variable } "z".$$

$$\begin{aligned} \text{Now: } \max_{\theta} p(x|\theta) &= \max_{\theta} \log p(x|\theta) \max_{\theta} \prod_{i=1}^N p(x_i|\theta) \text{ (product)} \\ &= \max_{\theta} \log \left( \prod_{i=1}^N p(x_i|\theta) \right) \\ &= \max_{\theta} \log p(x|\theta) \end{aligned}$$

$$\begin{aligned} \text{And } \max_{\theta} \log p(x|\theta) &= \log \prod_{i=1}^N p(x_i|\theta) \quad \{ \cancel{\text{sum}} \} \\ &= \sum_{i=1}^N \log p(x_i|\theta) \quad \{ \text{to optimize} \} \end{aligned}$$

$$\text{Now: } \log p(x|\theta) = \sum_{i=1}^N \log p(x_i|\theta)$$

Adding latent variables here

$$\log p(x|\theta) = \sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, t_i=c|\theta)$$

Now,  $\max_{\theta} p(x|\theta)$  function with latent variables  $t_i$  is transformed into the following

$$\log p(x|\theta) = \sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, t_i=c|\theta)$$

This is what we would like to maximize. We can use stochastic gradient descent here.

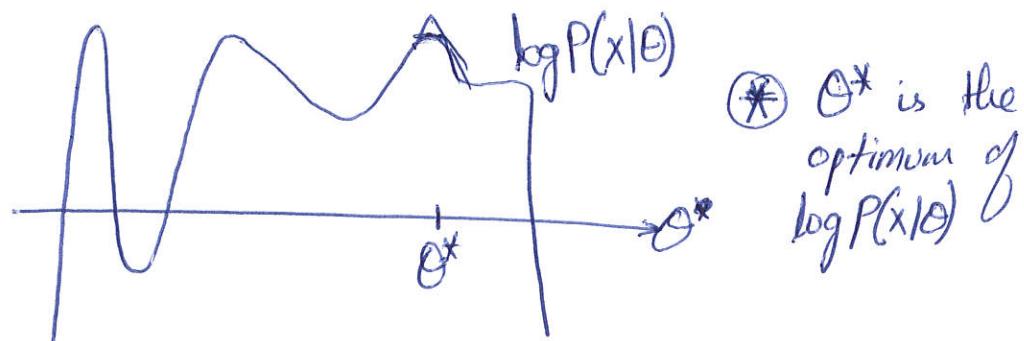
Now lets derive the lower bound of this function using Jensen's Inequality. Lets assume we came up with lower bound  $l(\theta)$ , then using Jensen's Inequality we know:

$$\sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, t_i=c|\theta) \geq l(\theta)$$

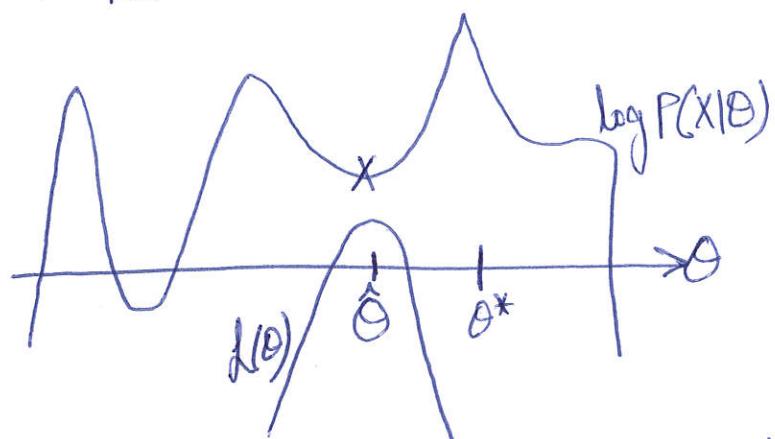
So, instead of maximizing  $\sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, t_i=c|\theta)$  we can maximize  $l(\theta)$ .

Let's assume that our marginal loglikelihood look as follows:

(35)



lets calculate  $L(\theta)$  for  $\log P(x|\theta)$ . Its value will be  $\leq \log P(x|\theta)$  at all points " $\theta$ ". Let assume we come up with the following  $L(\theta)$



The optimum of  $L(\theta)$  is at  $\hat{\theta}$  but we know that at  $\hat{\theta} = \hat{\theta}$ , we get local minimum of  $\log P(x|\theta)$  — our job was to find local maximum of  $\log P(x|\theta)$ . So, now let's generate family of lower bounds as follows. Let's introduce any distribution  $q(t_i=c)$  on our latent variable "t" as:  $q(t_i=c)$

Now multiply and divide  $\log P(x|\theta)$  by  $q(t_i=c)$

$$\begin{aligned}\Rightarrow \log P(x|\theta) &= \sum_{i=1}^N \log \frac{\sum_{c=1}^3 q(t_i=c)}{q(t_i=c)} p(x_i, t_i=c | \theta) \\ &= \sum_{i=1}^N \log \frac{\sum_{c=1}^3 q(t_i=c) \underbrace{p(x_i, t_i=c | \theta)}_{\alpha_1}}{q(t_i=c) \underbrace{\sum_{c=1}^3 q(t_i=c)}_{\alpha_2}} \underbrace{\sum_{c=1}^3 q(t_i=c)}_{\alpha_3} \\ &= \sum_{i=1}^N \log \left[ \underbrace{q(t_i=1)}_{\alpha_1} \underbrace{\frac{p(x_i, t_i=1 | \theta)}{q(t_i=1)}}_{\alpha_2} + \dots + \underbrace{\frac{q(t_i=3)p(x_i, t_i=3 | \theta)}{q(t_i=3)}}_{\alpha_3} \right]\end{aligned}$$

And we know that  
 $q(t_i=1) + q(t_i=2) + q(t_i=3) = 1$   
i.e.  $\alpha_1 + \alpha_2 + \alpha_3 = 1$

Thus :

$$\log P(X|\theta) = \sum_{i=1}^N \log \left[ \alpha_1 V_1^i + \alpha_2 V_2^i + \alpha_3 V_3^i \right]$$

where  $\alpha_i^i = q/(t_i=1)$

$$V_i^i = \frac{P(x_i, t_i=1 | \theta)}{q/(t_i=1)}$$

from Jensen's Inequality

$$\log \left( \sum_c \alpha_c V_c \right) \geq \sum_c \alpha_c \log(V_c)$$

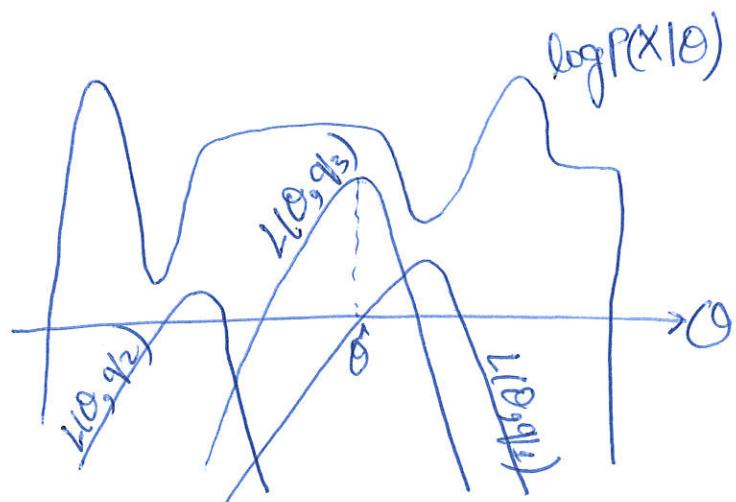
Therefore,  $\sum_c \alpha_c \log(V_c)$  is the lower bound of  $\log \left( \sum_c \alpha_c V_c \right)$ .  
This means, our lower bound of  $\log P(X|\theta)$  is given as follows:

$$\begin{aligned} \log P(X|\theta) &= \sum_{i=1}^N \log \left[ \sum_{c=1}^3 \underbrace{q/(t_i=c)}_{\alpha_c} \underbrace{\frac{P(x_i, t_i=c | \theta)}{q/(t_i=c)}}_{V_c} \right] \\ &\geq \boxed{\sum_{i=1}^N \sum_{c=1}^3 \underbrace{q/(t_i=c)}_{\alpha_c} \log \underbrace{\frac{P(x_i, t_i=c | \theta)}{q/(t_i=c)}}_{V_c}} \quad \rightarrow L(\theta, q) \end{aligned}$$

Now, using different values of "q", we can obtain family of lower bounds for  $\log P(X|\theta)$  given parameter " $\theta$ ".

Now our picture looks like this:

$$\log P(X|\theta) \geq L(\theta, q) \text{ for any } q$$



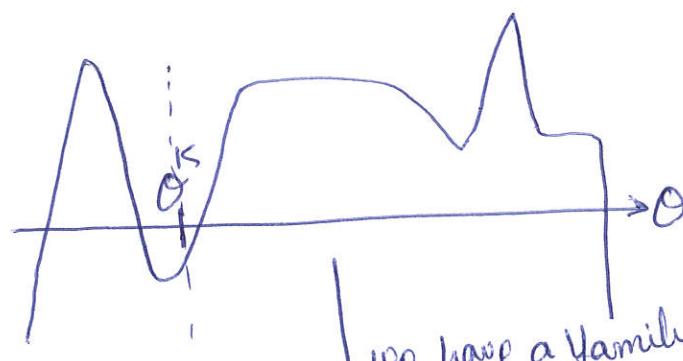
Since, we can choose  $\lambda(\theta, q_3)$ , get its maximum at  $\hat{\theta}$  and obtain  $\log P(X|\theta=\hat{\theta})$  which will be a good approximation of optimum for  $\log P(X|\theta)$ .

E-STEP:  $q^{k+1} = \underset{q}{\operatorname{argmax}} \lambda(\theta^k, q)$  at  $k^{\text{th}}$  iteration

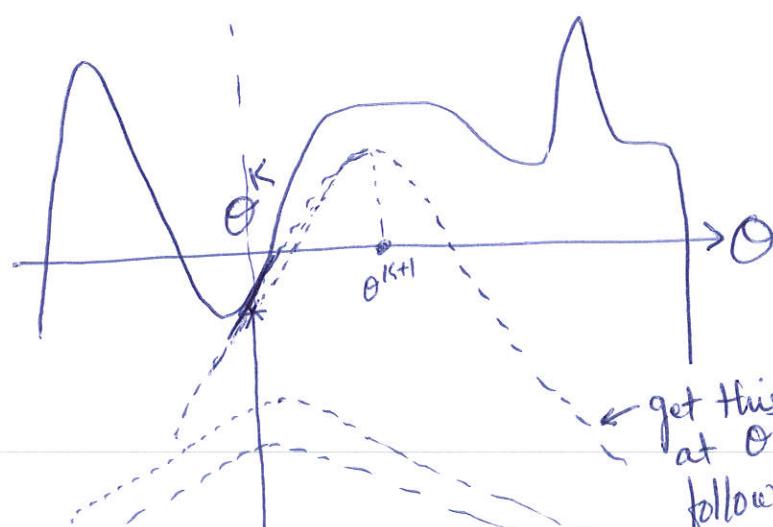
M-STEP:  $\theta^{k+1} = \underset{\theta}{\operatorname{argmax}} \lambda(\theta, q^{k+1})$

### Visualization

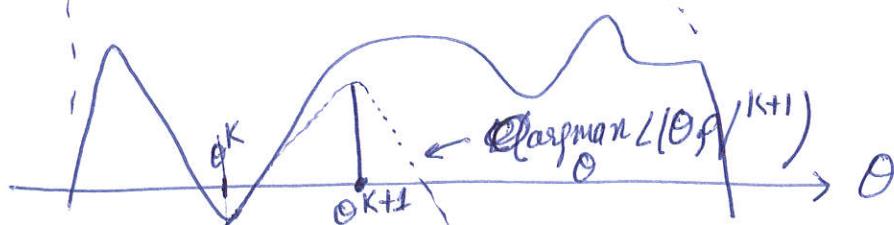
Iteration - 1<sup>st</sup>: Let initialize  $\theta^k$



we have a family of lower bound  $L(\theta^k, q)$ .  
lets find  $\underset{q}{\operatorname{argmax}} L(\theta^k, q)$

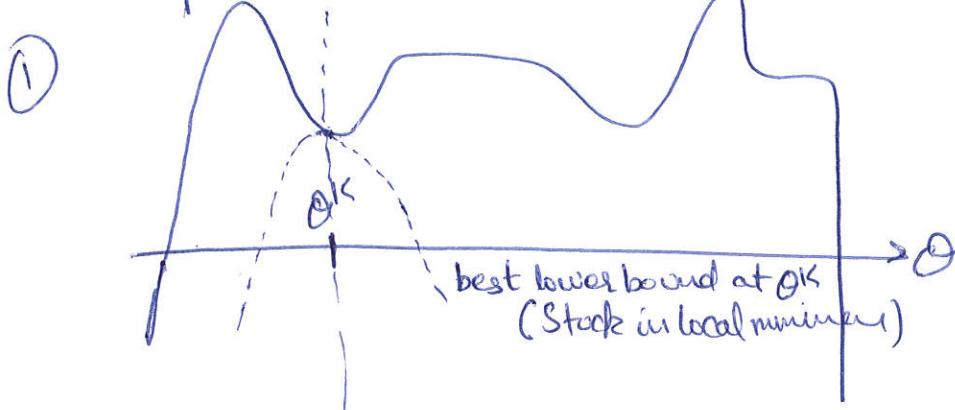


get this lower bound because  
at  $\theta = \theta^k$  it fulfills the  
following requirement  
 $q^{k+1} = \underset{q}{\operatorname{argmax}} L(\theta^k, q)$

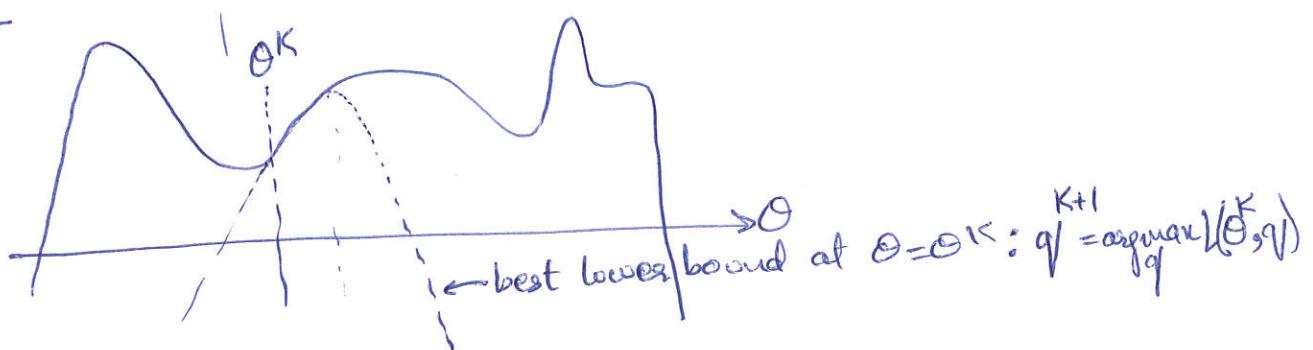


In  $\log p(x|\theta) \geq L(\theta, q)$  for any  $q$

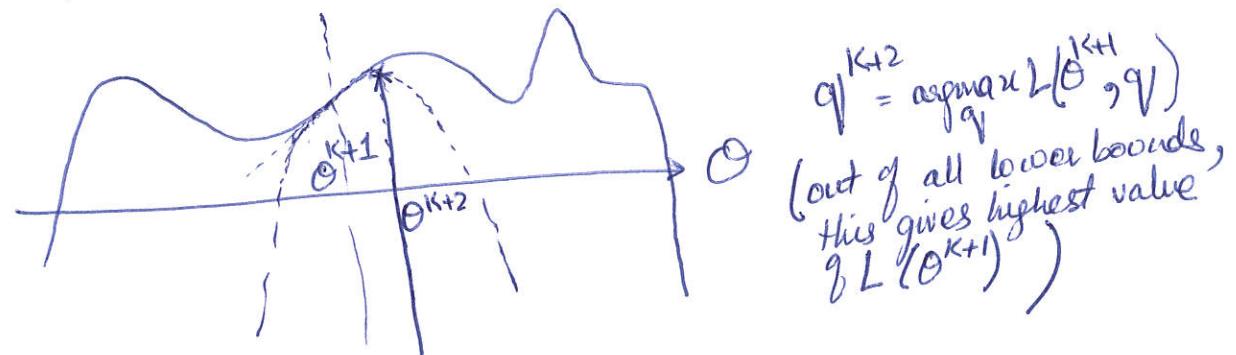
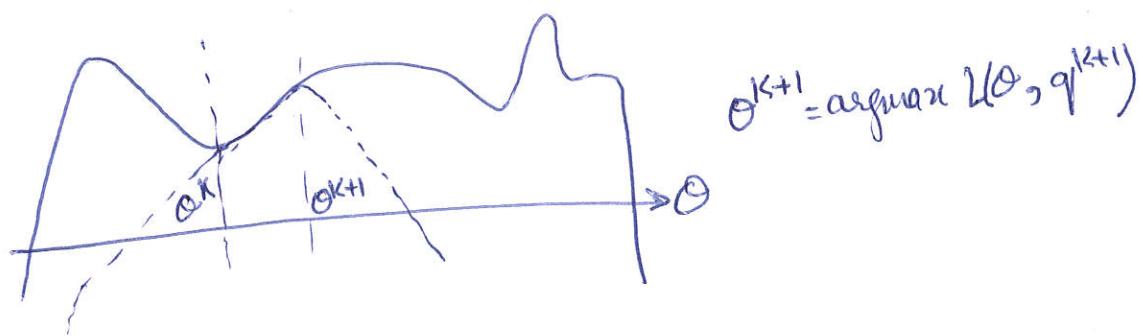
\*  $L(\theta, q)$  is called variational bound because you can vary it by varying " $q$ " & " $q$ " was not the part of original ~~model~~.



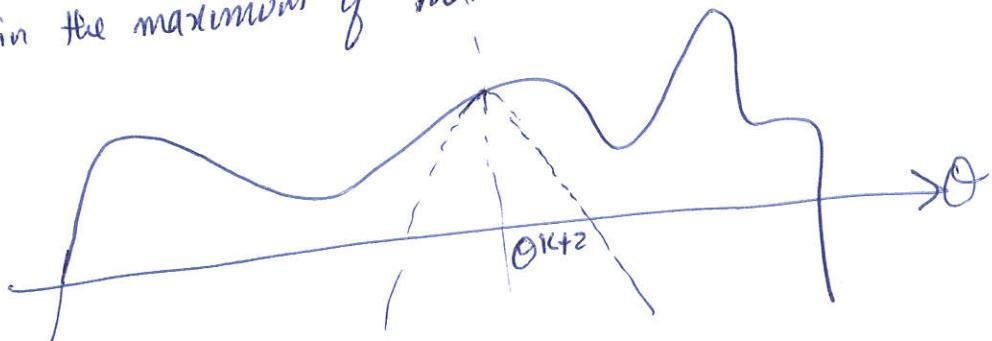
Restart



Get the max  $L(\theta, q^{K+1})$



Now, obtain the maximum of that lower bound in  $\theta$ .



$$\text{E-STEP} \circ \quad q^{(k+1)} = \underset{q}{\operatorname{argmax}} \ L(\theta^{(k)}, q)$$

Find the best lower bound

M-STEP  $\circ$

$$\theta^{(k+1)} = \underset{\theta}{\operatorname{argmax}} \ L(\theta, q^{(k+1)})$$

Get the  $\underset{\theta}{\operatorname{max}}$  of that lower bound

E-STEP DETAILS

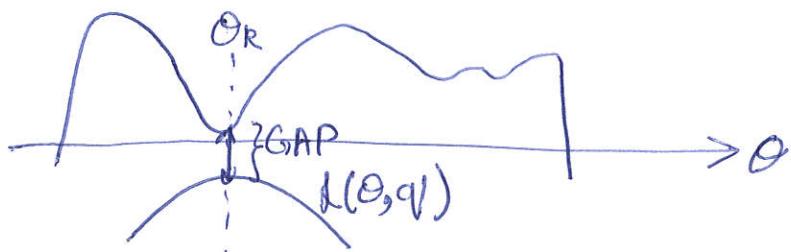
We now know that

$$\log P(X|\theta) \geq L(\theta, q)$$

variational lower bound which depends on "q". "q" is itself a distribution.

In E-STEP 8

$$q^{k+1} = \underset{q}{\operatorname{argmax}} L(\theta^k, q) \quad \left\{ \text{find the best lower bound} \right.$$



$$GAP = \underbrace{\log P(X|\theta)}_{\text{Marginal log likelihood}} - \underbrace{L(\theta, q)}_{\text{variational lower bound}} \quad \left\{ \begin{array}{l} \text{E-Step tries to minimize} \\ \text{this gap.} \end{array} \right.$$

$$\begin{aligned} &= \sum_{i=1}^N \log P(x_i|\theta) - \sum_{i=1}^N \sum_{c=1}^3 q(t_i=c) \frac{\log P(x_i, t_i=c|\theta)}{q(t_i=c)} \\ &= \sum_{i=1}^N \left[ \underbrace{\log P(x_i|\theta)}_{\text{Independent of } c} \times \sum_{c=1}^3 q(t_i=c) \right] - \sum_{i=1}^N \sum_{c=1}^3 q(t_i=c) \frac{\log P(x_i, t_i=c|\theta)}{q(t_i=c)} \\ &= \sum_{i=1}^N \left[ \sum_{c=1}^3 \log P(x_i|\theta) q(t_i=c) - \sum_{i=1}^N \sum_{c=1}^3 q(t_i=c) \frac{\log P(x_i, t_i=c|\theta)}{q(t_i=c)} \right] \\ &= \sum_{i=1}^N \sum_{c=1}^3 q(t_i=c) \left[ \log P(x_i|\theta) - \log \frac{P(x_i, t_i=c|\theta)}{q(t_i=c)} \right] \\ &= \sum_{i=1}^N \sum_{c=1}^3 q(t_i=c) \log \left[ \frac{P(x_i|\theta)}{P(x_i, t_i=c|\theta)} \times q(t_i=c) \right] \end{aligned}$$

IMP Simplification (Q.viz)

$$\frac{P(x_i|\theta)}{P(x_i, t_i|\theta)} = \frac{P(x_i|\theta)}{P(t_i|x_i, \theta) P(x_i|\theta)} = \frac{1}{P(t_i|x_i, \theta)}$$

Using definition of conditional probability.

$$p(x_i, t_i=c | \theta) = \underbrace{p(t_i=c | x_i, \theta)}_{\text{conditional}} \underbrace{p(x_i | \theta)}_{\text{prior distribution}}$$

$$\begin{aligned} &= \sum_{i=1}^N \sum_{c=1}^3 q(t_i=c) \log \frac{p(x_i | \theta)q(t_i=c)}{p(t_i=c | x_i, \theta) p(x_i | \theta)} \\ &= \sum_{i=1}^N \sum_{c=1}^3 q(t_i=c) \log \frac{p(t_i=c)}{\underbrace{p(t_i=c | x_i, \theta)}} \\ &\quad \underbrace{\text{KL}\left(q(t_i) \parallel p(t_i=c | x_i, \theta)\right)} \end{aligned}$$

$$\Rightarrow \text{GAP} = \log P(X|\theta) - \lambda(\theta, q)$$

$$= \sum_{i=1}^N \text{KL}\left(q(t_i) \parallel \underbrace{p(t_i | x_i, \theta)}_{\text{posterior distribution}}\right)$$

→ We want to minimize GAP

$$\begin{aligned} &\rightarrow \text{we want to maximize } \lambda(\theta, q) \text{ w.r.t } q \\ &\rightarrow \text{this means we have to } \min_q [-\lambda(\theta, q)] \end{aligned}$$

$$\text{GAP} = \underbrace{\log P(X|\theta)}_{\text{Independent of } q} - \underbrace{\lambda(\theta, q)}_{\max_q \min_q}$$

Thus  $\max_q \lambda(\theta, q)$  is equivalent to minimizing  $\log P(X|\theta) - \lambda(\theta, q)$ .

Thus minimizing  $\log P(X|\theta) - \lambda(\theta, q)$  is equivalent to minimizing  $\sum_{i=1}^N \text{KL}\left(q(t_i) \parallel p(t_i | x_i, \theta)\right)$  or sum of KL divergences.

This implies  $\max_q \lambda(\theta, q)$  is equivalent to  $\min_q \sum_{i=1}^N \text{KL}\left(q(t_i) \parallel p(t_i | x_i, \theta)\right)$

Recall that KL-divergence had two properties:

$$\text{i) } \text{KL}(q || p) \geq 0$$

$$\text{ii) } \text{KL}(q || q) = 0$$

Thus  $\sum_{i=1}^N \text{KL}\left(q(t_i) \parallel p(t_i | x_i, \theta)\right) \geq 0$  (as it can never be negative)

In order to  $\min_{q} \sum_{i=1}^K KL[q_i(t_i) || p(t_i | x_i, \theta)]$ ,  
 we will set  $q_i(t_i) = \underbrace{p(t_i | x_i, \theta)}_{\text{posterior distribution}}$  as  $KL(q || q) = 0$ . It  
 will return us global minimum in "q"

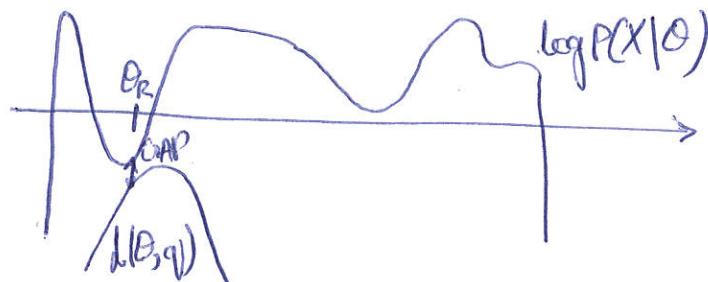
### Outcome

Minimizing sum of KL-divergence for E-step simply asks us to set  $q_i(t_i) = p(t_i | x_i, \theta)$ . It will return us the best lower-bound  $L(\theta^*, q)$  at  $\theta = \theta^*$

### E-Step

We know that:  $\log P(X|\theta) \geq L(\theta, q)$

$$\text{E-step: } q^{k+1} = \max_q L(\theta^k, q)$$

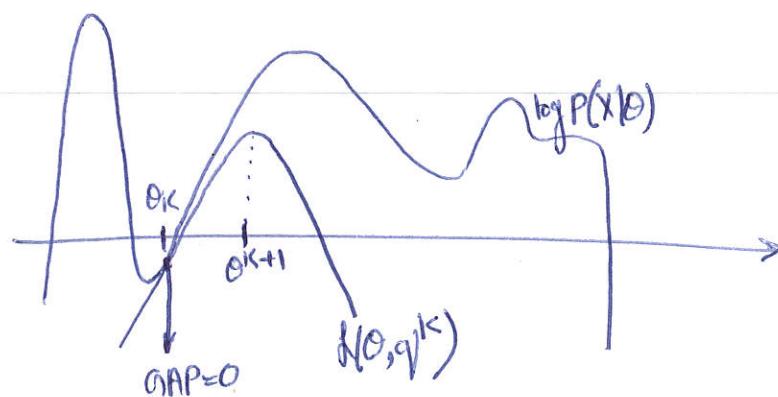


Conclusion To get the best lower bound, we define gap.

$$\text{gap: } \log P(X|\theta) - L(\theta, q) = \leq \sum_i KL[q_i(t_i) || p(t_i | x_i, \theta)]$$

This will minimize gap but maximize lower bound  $L(\theta^k, q)$ . To set gap=0, use  $KL(q || q) = 0 \Rightarrow \text{set } q_i(t_i) = p(t_i | x_i, \theta)$

$$\text{E-Step: } \arg \max_{q_i(t_i)} L(\theta^k, q) = p(t_i | x_i, \theta)$$



M-STEP

In M-step, we want to maximize variational lower bound  $\mathcal{L}(\theta, q)$  w.r.t " $\theta$ "

$$\begin{aligned}\mathcal{L}(\theta, q) &= \sum_i \sum_c q(t_i=c) \log \frac{p(x_i, t_i=c | \theta)}{q(t_i=c)} \\ &= \underbrace{\sum_i \sum_c q(t_i=c) \log p(x_i, t_i=c | \theta)}_{\text{Eq log } p(X, T | \theta)} - \underbrace{\sum_i \sum_c q(t_i=c) \log q(t_i=c)}_{\text{constant w.r.t "q"}}$$

$\text{Eq log } p(X, T | \theta)$  = expected value of  $\log$  of joint distribution of  $X \in T$  w.r.t variational distribution " $q$ ".

$$= \underbrace{\text{Eq log } p(X, T | \theta)}_{\text{function}} + \text{const}$$

We choose  $p(X, T | \theta)$  such that its  $\log$  becomes concave and  $\text{Eq}[\text{concave function}]$  is also concave. So you can easily optimize it. It ensures that this function has global maxima

SKETCH of EM-Algorithm

for  $K=1 \dots$  until convergence

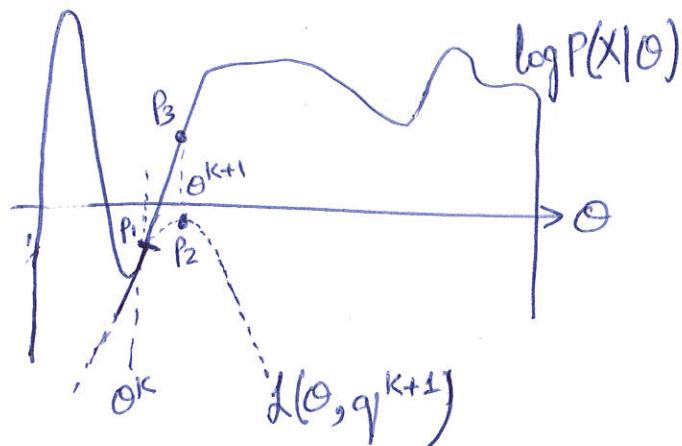
$$\begin{aligned}\text{E-STEP: } q^{(K+1)} &= \underset{q}{\operatorname{argmax}} \mathcal{L}(\theta^k, q) \\ &= \underset{q}{\operatorname{argmin}} \text{KL} \left[ \underbrace{q(T)}_{\text{Variational distribution}} \middle\| \underbrace{p(T | X, \theta^k)}_{\text{posterior distribution}} \right] \\ &\stackrel{<=>}{=} q^{(K+1)}(t_i) = p(t_i | x_i, \theta^k)\end{aligned}$$

M-STEP:

$$\begin{aligned}\theta^{K+1} &= \underset{\theta}{\operatorname{argmax}} \mathcal{L}(\theta, q^{(K)}) \\ &= \underset{\theta}{\operatorname{argmax}} \left[ \underset{q}{\operatorname{E}} \log p(x_i, t_i | \theta) \right] \\ &= \underset{\theta}{\operatorname{argmax}} \underbrace{\underset{q}{\operatorname{E}} q^{(K+1)} \log p(X, T | \theta)}_{\text{Expected value of joint distribution}}\end{aligned}$$

# Convergence of EM

(45)



$$\underbrace{\log p(X|\theta^{K+1})}_{P_3} \geq \underbrace{L(\theta^{K+1}, q^{K+1})}_{P_2} \geq \underbrace{L(\theta^K, q^{K+1})}_{P_1}$$

Can we compare value of  $L(\theta^K, q^{K+1})$  with log-likelihood  $\log p(X|\theta^K)$  in general case?

Ans: Yes, since for any "q" the following inequality holds true  
 $L(\theta^K, q) \leq \log p(X|\theta^K)$ , but for  $q^{K+1}$  is the maximizer of the lower bound  $L(\theta^K, q^{K+1}) = \max_q L(\theta^K, q)$  and in the best case it became as large as  $\log p(X|\theta^K)$ .

Hence  $\underbrace{\log p(X|\theta^{K+1})}_{P_3} \geq \cancel{L(\theta^K, q^{K+1})} \geq \underbrace{\log p(X|\theta^K)}_{P_1}$

i.e. the marginal log likelihood never decreases during an iteration of EM.

- On each iteration, EM doesn't decrease objective (good for debugging)
- Guaranteed to converge to a local maximum (or saddle point).



w2-p1-l1: Latent Variable Models

(19)



W2-P2-25

(4)

Example: EM-for discrete mixture, E-Step

