

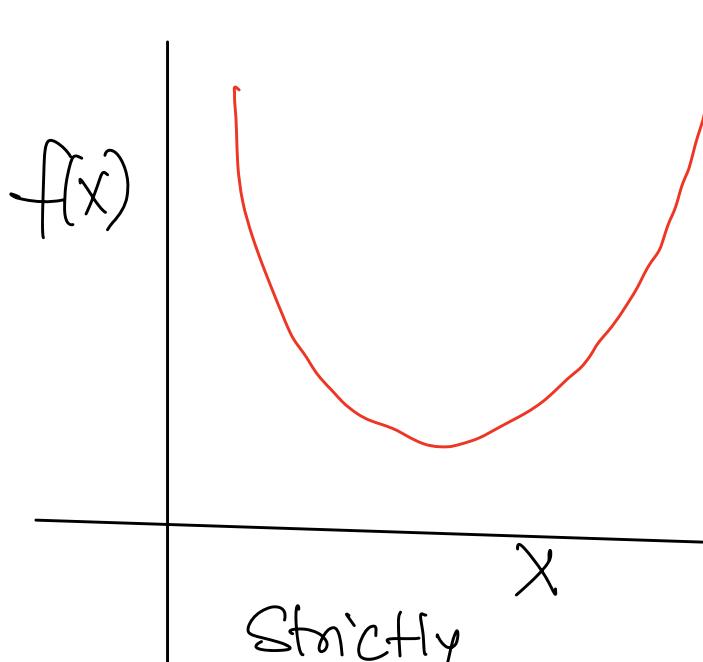
# E.M. Algorithm

MLE in presence of latent variables

True Model:  $p(x, z; \theta)$

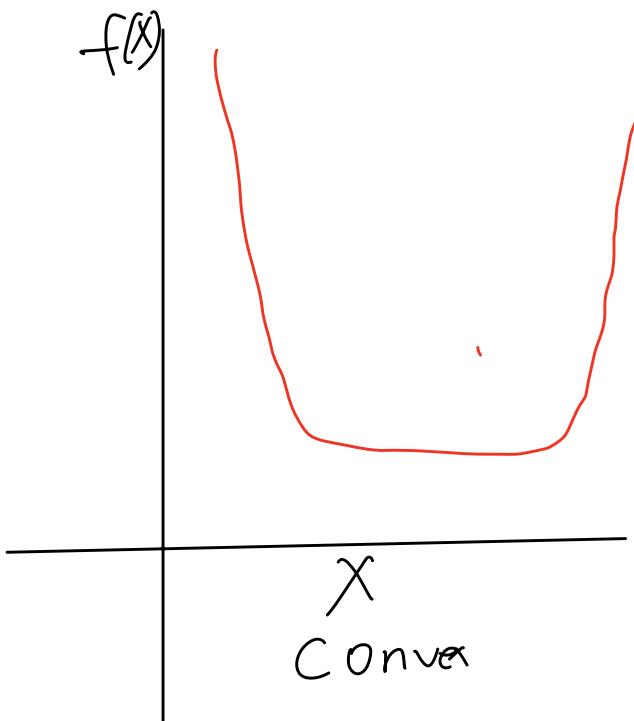
when  $z$  is unobserved.

Maximization  $\ell(\theta) = \log p(x; \theta)$



Strictly  
Convex

$$f''(x) > 0$$



Convex

$$f''(x) \geq 0$$

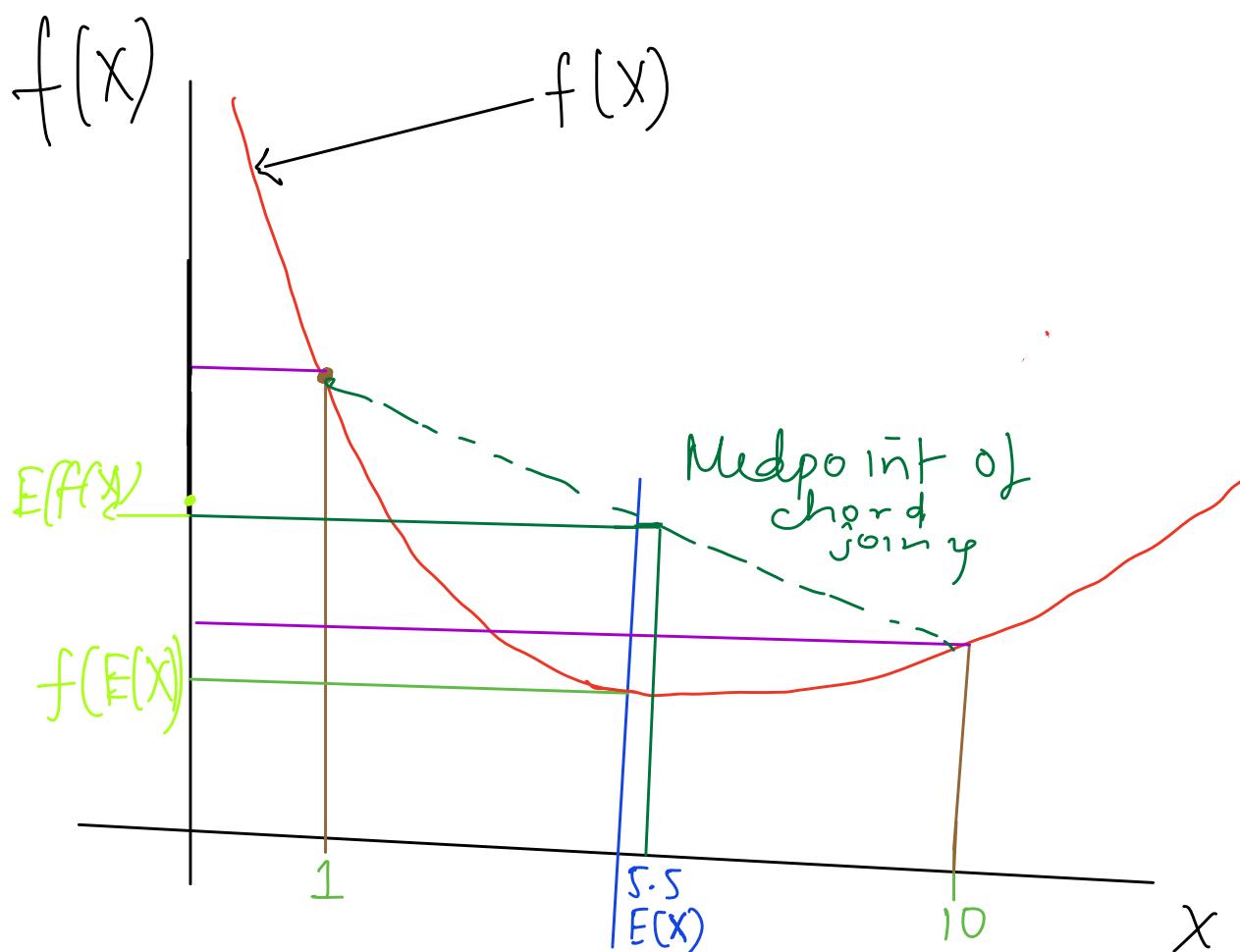
Jensen Inequality

$E[f(X)] \geq f(E[X])$  where  $f$  is convex

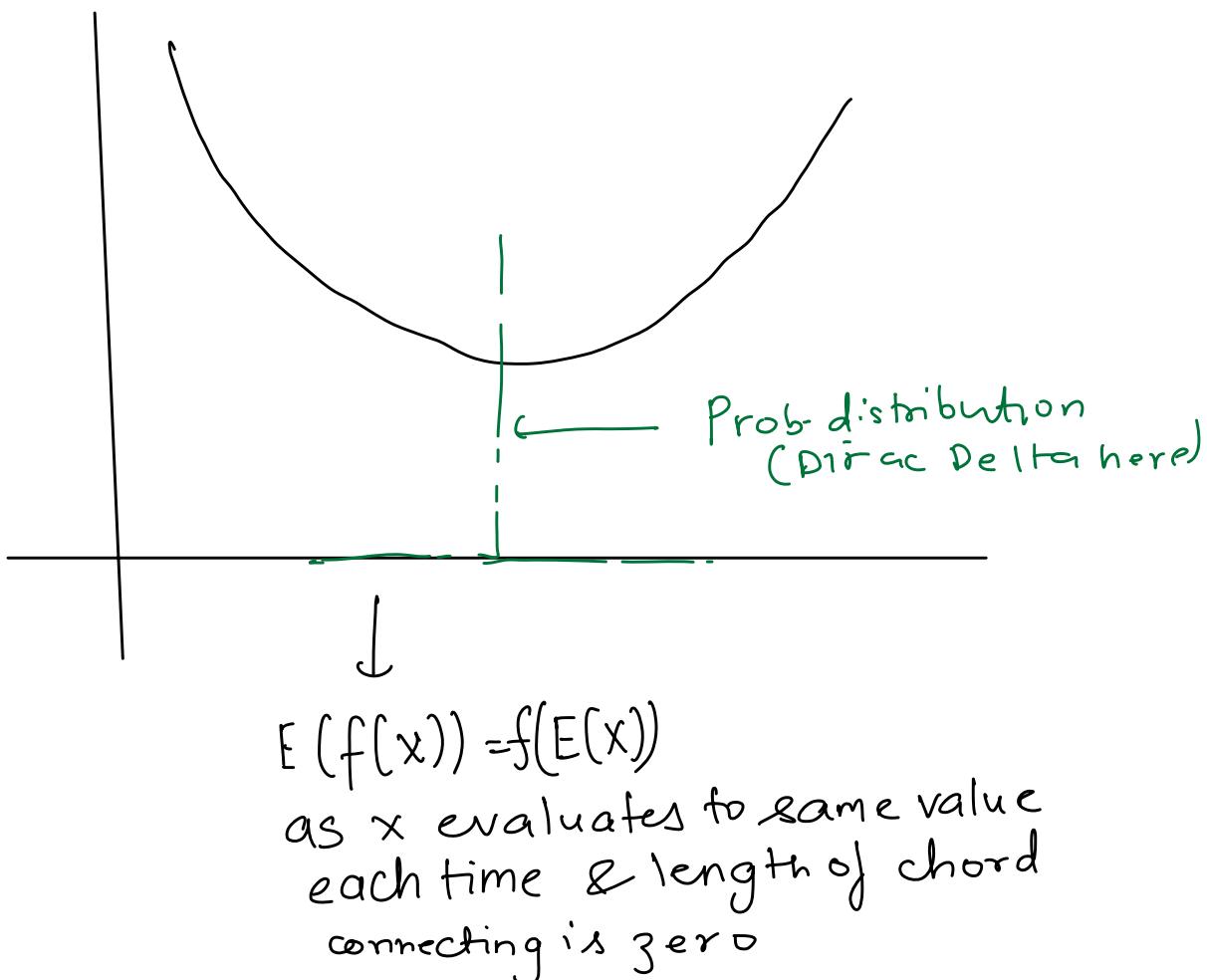
If  $f$  is strictly convex,  
then,

$$\text{if } E[f(X)] = f(E[X])$$

then  $X = E(X)$  w.p. 1  
i.e.  $X \rightarrow \text{constant}$



$p(x)$  which is discrete i.e. can be 1 or 0 with  $p = \frac{1}{2}$



$f''(x) = 0$  so both convex & concave

Convex	Concave	strict
$x^2$	$-x^2$	✓
$mx+c$	$mx+c$	✗
$e^x$	$-e^x$	✓
$-\log x$	$\log x$	✓

If  $f$  is concave

$$\text{e.g. } f(x) = \log(x)$$

$$E[\log(X)] \leq \log E[X]$$

Goal is max

$$\sum_{i=1}^n \log p(x^{(i)}; \theta) \quad z \rightarrow \text{unobserved}$$

$$\max \log p(x; \theta)$$

$$\log \sum_z p(x, z; \theta)$$

$$\log \sum_z \frac{Q(z) p(x, z; \theta)}{Q(z)}$$

[where  $Q(z)$  is  
arbitrary prob-  
dist over  $z$ )  
 $Q(z) > 0$  for all  $z$ ]

$$\log \mathbb{E}_{z \sim Q} \left[ \frac{p(x|z; \theta)}{Q(z)} \right]$$

$$\geq \underset{z \sim q}{\mathbb{E}} \left[ \log \frac{P(x, z; \theta)}{q(z)} \right]$$

$$:= \text{ELBO}(x; \vartheta, \theta)$$

Evidence lower bound

i.e. setting  
good  $\vartheta$

so if we find  $\max \text{ELBO}$ , then  
 $\log p(x)$  will be greater

Are there cases  $\log p(x; \theta) = \text{ELBO}(x; \vartheta, \theta)$ ?

→ Yes when  $\frac{P(x, z; \theta)}{q(z)} \rightarrow \text{constant}$

constant w.r.t  $z$

$$Q(z) = \sum_{\mathcal{Z}} p(x, z; \theta)$$

$$Q(z) \propto p(x, z; \theta)$$

Since  $Q(z)$  is prob. distribution  
it should be summed to 1

$$\begin{aligned} Q(z) &= \frac{p(x, z; \theta)}{\sum_{\mathcal{Z}} p(x, z; \theta)} \\ &= \frac{p(x, z; \theta)}{p(x; \theta)} \\ &= p(z|x; \theta) \end{aligned}$$

$$Q(z) = p(z|x; \theta) \Rightarrow \log p(x) = E(\text{BO})$$

# EM Algo

Random  $\theta$

## E Step

For each  $i$ , set

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

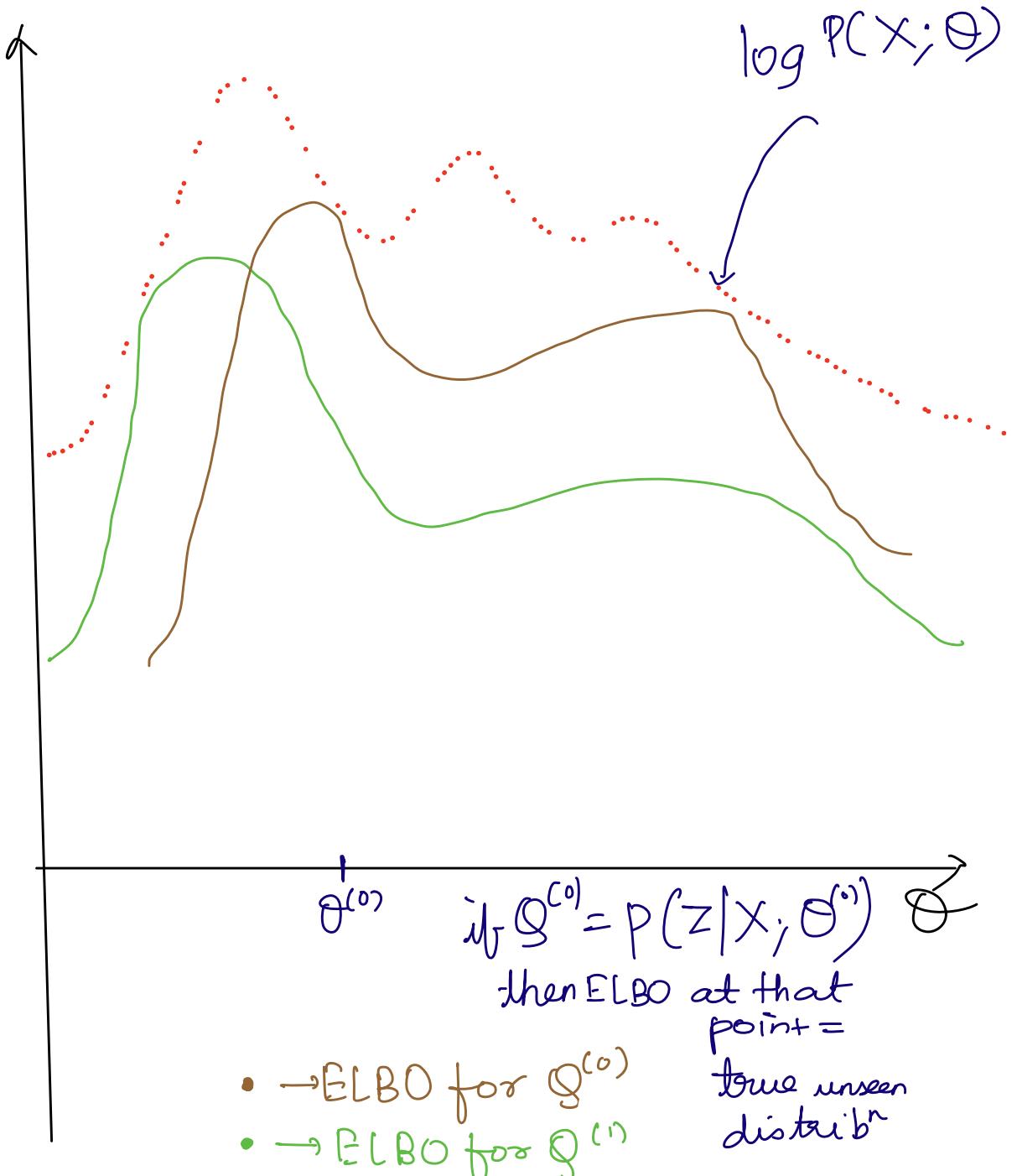
## M Step

Set

$$\theta := \arg \max_{\theta} \sum_{i=1}^n \text{ELBO}(x^{(i)}; Q_i, \theta)$$
$$= \sum_{i=1}^n \sum_z Q_i(z^{(i)}) \log \frac{P(x, z; \theta)}{Q(z^{(i)})}$$

$Q$  is constant

Although it has  $\theta$ , we assume it's free from previous.



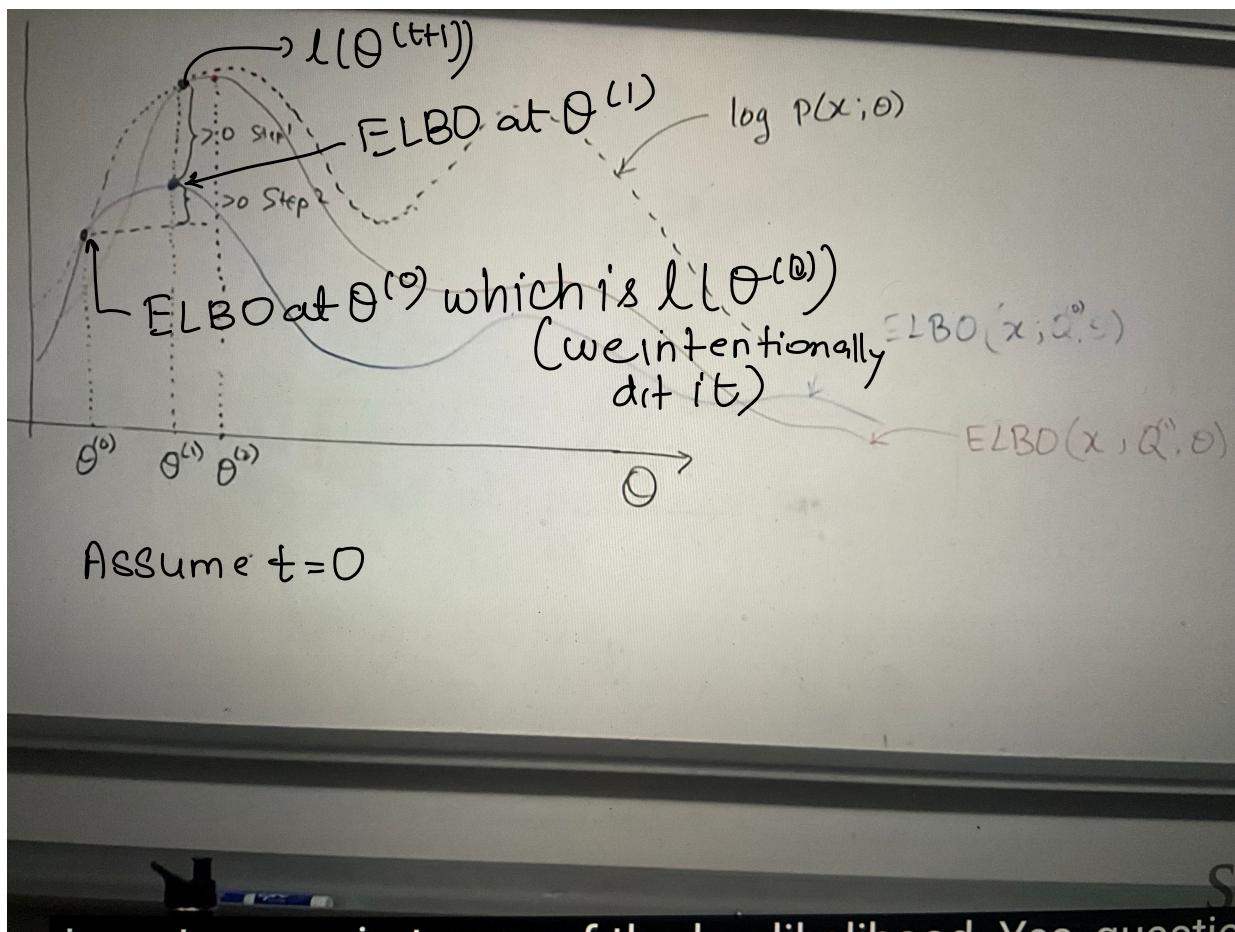
using ELBO, we indirectly max  $\log p(x; \theta)$   
 & find  $\theta$  which is a local optima

$$\text{GPI}_Z^{(3)} = \begin{bmatrix} 0.1 \\ 0.9 \\ 0.01 \end{bmatrix} - \text{Cluster 1}$$

↑  
3rd ex

$$- \text{Cluster 12}$$

Proof that EM converges



For every t

$$l(\theta) = \log P(x; \theta)$$

$$l(\theta^{(t+1)}) \geq l(\theta^{(t)})$$

$$l(\theta^{(t+1)}) \geq \text{ELBO}(x; Q^{(t)}, \theta^{(t+1)}) \rightarrow \text{Jensen}$$

$$\geq \text{ELBO}(x; Q^{(t)}, \theta^{(t)}) \rightarrow \text{M-step}$$

$$= l(\theta^{(t)}) \rightarrow \text{Jensen}$$

GMM via EM

① First write out model  $P(X, Z; \theta)$

$Z \sim \text{Multinomial}$

$X|Z \sim N(\mu_Z, \Sigma_Z)$

② Then write what is latent & what are observed

$x$  - Evidence

$z$  - Latent

[Unknown  $z \rightarrow$  example specific

Unknown parameters  $\rightarrow$  global]

③ In E-step, we estimate  $z$ , for each eg.  
holding parameter fixed

④ In M-step, we update parameters holding  
 $z$  fixed

## E - Step

For each  $i$ ,

$$w_j^{(i)} = Q_i(z^{(i)}=j) = p(z^{(i)}|x^{(i)}; \phi, \mu, \Sigma)$$

$$= \frac{p(x|z)p(z)}{p(x)}$$

$$= \frac{p(x|z)p(z)}{\sum_j p(x|z=j)p(z=j)}$$

$$= \frac{\frac{1}{(2\pi)^{d/2}} |\Sigma_j|^{1/2} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right\} \times \phi_j}{\sum_j \left( \quad \right)}$$

M-Step

$$\mu, \Sigma, \phi = \arg \max_{\mu, \Sigma, \phi} \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log \frac{p(x^{(i)}, z^{(i)}; \mu, \Sigma, \phi)}{w_j^{(i)}}$$

$$= \arg \max_{\mu, \Sigma, \phi} \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log \frac{p(x^{(i)} | z^{(i)}=j) p(z^{(i)}=j)}{w_j^{(i)}}$$

[Here for M-step hold  $w$  constant]

Solve it then take gradient solve for 0  
then you get G.M.M. update rules

[ $k$  diff  $\mu, \Sigma$ ]

[If we take mode instead of expectation,  
we get K-Means]

FACTOR ANALYSIS

$$x^{(i)} \in \mathbb{R}^d$$

$$S = \{x_i^{(i)}\}_{i=1}^n \quad d > n > k$$

$$\Sigma \in \mathbb{R}^{d \times d}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i^{(i)} - \mu)(x_i^{(i)} - \mu)^T$$

$$\Sigma \rightarrow \text{rank } n$$

but actually it is rank d

$$k < n$$

k-low dimensional subspace

$$z^{(i)} \sim N(0, I) \quad I - k \text{ dim}$$

$$z^{(i)} \in \mathbb{R}^k$$

$$x^{(i)} | z^{(i)} \sim N(\mu + Lz, \Psi)$$

$$L \in \mathbb{R}^{d \times k}$$

$$\Psi \in \mathbb{R}^{d \times d} \rightarrow \text{diagonal matrix}$$

[In GMM, each  $x$  is discrete and has separate  $\mu, \Sigma$

Here  $\mu, \Sigma$  constant across egs]

Here, we are not doing classification

We want to find a lower dimensional subspace in which  $z$  resides and find  $x$  for each  $z$

Suppose there are  $d$  temp. sensor & we get temp values at time  $t$

$$x^{(t)} \in \mathbb{R}^d$$

$d$  very high & all values are not independent. In one room, value by sensor is almost same.

So we try to estimate a few  $k$  factors which decide temp. & then make int.

$$z^{(t)} \in \mathbb{R}^k$$

## Model

$$z \sim N(0, I)$$

$$x|z \sim N(\mu + Lz, \Psi)$$

$$z \sim N(0, I)$$

$$\varepsilon \sim N(0, \Psi)$$

$$x = \mu + Lz + \varepsilon \rightarrow \text{Location Scale Transform}$$

$$E(x) = \mu$$

Joint distribution

$$\begin{bmatrix} z \\ x \end{bmatrix} \sim N\left(\mu_{zx}, \Sigma\right)$$

$$\mu_{zx} = \begin{bmatrix} \vec{0} \\ \mu \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \Sigma_z & \Sigma_{zx} \\ \Sigma_{zx}^T & \Sigma_x \end{bmatrix} = \begin{bmatrix} I & L^T \\ L & LL^T + \Psi \end{bmatrix}$$

Observe -  $x$

Latent -  $z$

Parameters :  $\mu, L, \Psi$

NOTE:

$$l(\mu, L, \Psi) = \log p(x; \rightarrow)$$

⋮

$\curvearrowright$  diagonal matrix

$$x \sim N(\mu, LL^T + \Psi)$$

$L$ , low rank

$$\log p(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |LL^T + \Psi|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu)^T (LL^T + \Psi)^{-1} (x - \mu)\right)$$

This is correct but you won't get closed form sol'n here.

E Step

$$p(z|x) \sim N \left[ L^T (L L^T + \Psi)^{-1} (x - \mu), \underbrace{I - L^T (L L^T + \Psi)^{-1} L}_{\text{schur complement}} \right]$$

$$Q_i(z^{(i)}) = N \left[ \quad, \quad \right]$$

M-Step

$$\mu, L, \Psi = \arg \max_{\mu, L, \Psi} \sum_{i=1}^n \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \mu, L, \Psi)}{Q_i(z^{(i)})}$$

$$= \sum_{i=1}^n E \left[ \log \frac{p(x^{(i)}, z^{(i)}; \mu, L, \Psi)}{Q_i(z^{(i)})} \right] \quad \begin{matrix} \text{constant} \\ \text{always} \end{matrix}$$

$$= \sum_{i=1}^n E \left[ \log p(x^{(i)}|z^{(i)}) + \boxed{\log p(z^{(i)})} - \cancel{\log(Q_i(z^{(i)})}) \right]$$

In this case  
z holds no info about  
parameters

$$L = \sum_{i=1}^n \left( x^{(i)} - \mu \right) \mu_{Z|x}^T \left( \sum_{i=1}^n \mu_{Z|x} \mu_{Z|x}^T + \Sigma_{Z|x} \right)$$

$$\hat{M} = \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

$$\text{Top} = \frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)T} - x^{(i)} M_{2|X}^T L^T - L M_{2|X} x^{(i)T} + L \left( M_{2|X} M_{2|X}^T + \sum_{z|x} L^T \right)$$

$$\Psi_{ii} = \underline{\Phi}_{ii}$$

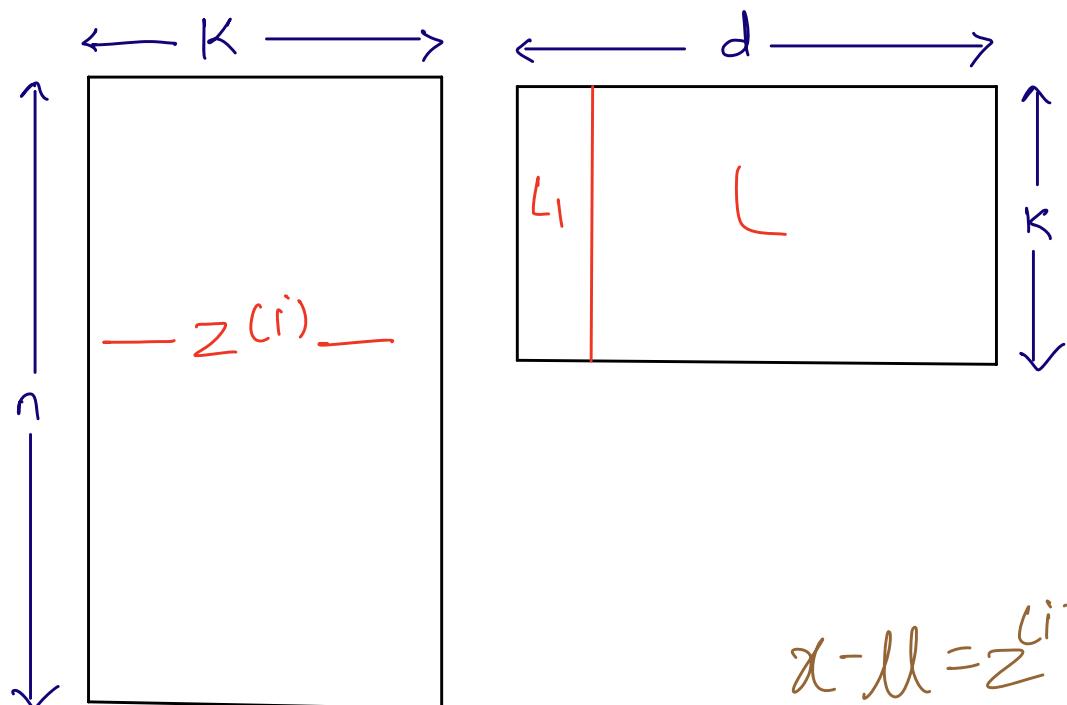
FACTOR ANALYSIS

$$\text{Model: } z \sim N(0, I) \in \mathbb{R}^k$$

$$x|z \sim N(\mu + Lz, \Psi) \in \mathbb{R}^d$$

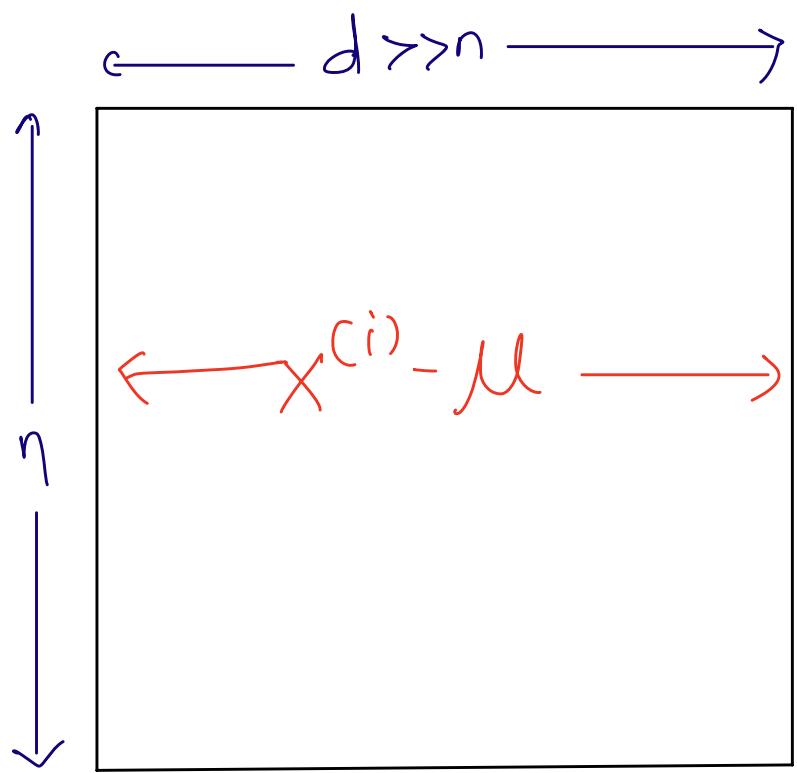
Parameters:

- $\mu \in \mathbb{R}^d$
- $L \in \mathbb{R}^{d \times k}$
- $\Psi \in \mathbb{R}^{d \times d}$  → diagonal matrix  
(each dimension in d has diff noise added to it)



$$x - \mu = z^{(i)\top} L + \Psi_{ii}$$

?



Think of it as  $d$  linear regression  
happening simultaneously

$\mu_{z^{(i)}|x^{(i)}} \approx \hat{z}^{(i)}$  best estimate of  $z^{(i)}$

$$\begin{aligned} L \mu_{z^{(i)}|x^{(i)}} &= L \hat{z}^{(i)} \approx x^{(i)} - \mu \\ &:= \hat{x}^{(i)} - \mu \end{aligned}$$

$$L = \left[ \sum_{i=1}^n y^{(i)} z^{(i)T} \right] \left[ Z^T Z + \Sigma \right]^{-1}$$

$$L^T = (Z^T Z + \sim)^{-1} Z^T y$$

L) Regularised L.R.

$$\Psi_{ii} \approx \frac{\sum_{i=1}^n (x^{(i)} - \bar{x}^{(i)}) (x^{(i)} - \bar{x}^{(i)})^T}{[\sum_{i=1}^n L^T]_{ii}}$$

$$y = \theta^T x + \varepsilon$$

$$\varepsilon = \frac{1}{n} \sum_{i=1}^n (y - \theta^T x)^2$$

We take  $\Psi_{ii}$  because we don't care about covariance of noise bet' 1 LR & other. We are evaluating each seperately

Given  $x$ , we have to estimate  $z$  &  $L$ .

$z$  &  $x$  have affine relationship

P.C.A.