

Lecture-15

* EM algorithm and mixture of
Gaussian model:

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

$z^{(i)} \sim \text{Multinomial } (\phi_j)$

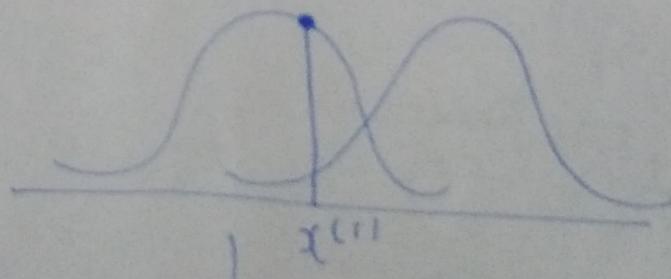
$$[P(z^{(i)}=j) = \phi_j].$$

$$x^{(i)}|z^{(i)} \sim N(\mu_j, \Sigma_j)$$

→ Deriving:-

① E-step:-

$$\boxed{w_j^{(i)}} = Q_j(z^{(i)}=j) = \frac{P(z^{(i)}=j|x^{(i)}; \phi, \Sigma)}{\sum P(z^{(i)}=j)}$$



$w_j^{(i)}$ from where
 $x^{(i)}$ is likely to come
(which Gaussian)

② M-Step:-

$$\max_{\phi, \mu, \Sigma} \sum_i \sum_{j \neq i} Q_j(z^{(i)}) \log \left(\frac{P(x^{(i)}, z^{(i)}; \phi, \Sigma)}{Q_j(z^{(i)})} \right)$$

$$= \sum_i \sum_j Q_i(z^{(i)}) \log \frac{1}{(\sigma\pi)^{1/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^\top \Sigma_j^{-1} (x^{(i)} - \mu_j)\right)$$

$w_j^{(i)}$

ϕ_j

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

$w_j^{(i)}$

\rightarrow This formula

$$\nabla_{w_j} (\dots) = 0 \Rightarrow \nu_j = \frac{\sum_i w_j^{(i)} (x^{(i)})}{\sum_i w_j^{(i)}}$$

\rightarrow Remember $w_j^{(i)}$ is the strength with which $x_j^{(i)}$ is a signal to Gaussian j ($P(z_j^{(i)} | x^{(i)}, \dots)$).

\hookrightarrow Formula for mixture of Gaussian model derived from EM algorithm.

$$\nabla_\phi (\dots) = 0 \Rightarrow \phi_j = \frac{1}{m} \sum_i w_j^{(i)}$$

$\nabla_\Sigma (\dots) = 0 \Rightarrow \Sigma_j = \frac{\sum_{i=1}^m (x - \mu_j)^\top (x - \mu_j)}{\sum_{i=1}^m w_j^{(i)}}$

→ Factor analysis model:-

→ In this case $Z^{(i)}$ is continuous random variable rather than discrete random variable.

→ One equivalent view of EM-algorithm:-

$$CJ(\theta, Q) = \sum_i \sum_{Z^{(i)}} Q_i(z^{(i)}) \log \left(\frac{P(x^{(i)}, Z^{(i)}, \theta)}{Q_i(z^{(i)})} \right).$$

• We know $\ell(\theta) \geq J(\theta, Q)$
for any θ, Q .

E-step \Rightarrow Maximize J w.r.t Q

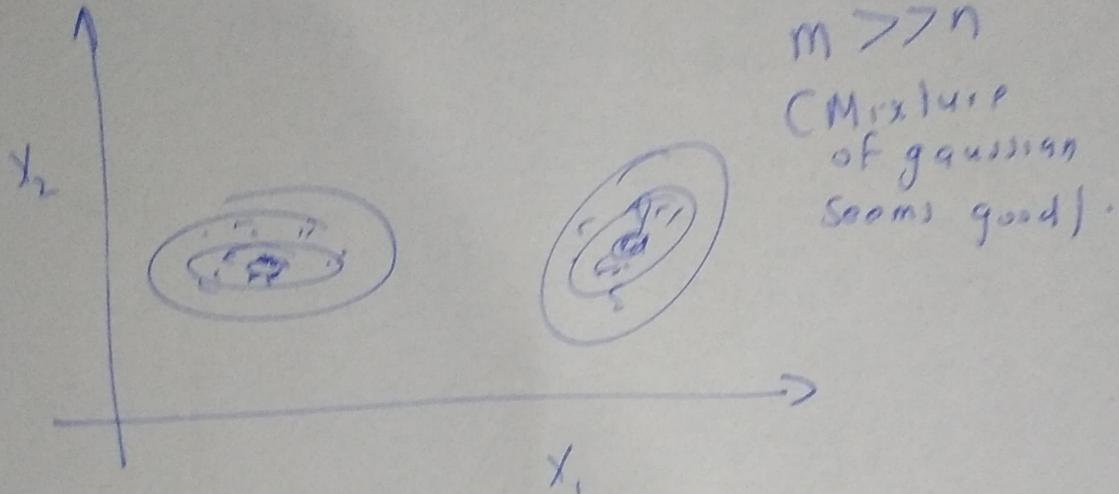
M-step \Rightarrow Maximize J w.r.t. θ .

→ This procedure is called coordinate ascent because you are maximizing $J(\theta, Q)$ w.r.t to one parameter at a time. (EM is a coordinate ascent algorithm).

(16)

→ What is Factor Analysis and
comparison with mixture of gaussian?

$$n=2 \quad m=100$$



$$m >> n$$

(Mixture
of gaussian
Seems good).

$$m \approx n, \text{ or } m \ll n$$

(Factor Analysis is better in this
case as compare to mixture of gaussian).

Say $m=30, n=100$ (100 dimension with
30 examples).

↳ Problem with applying Gaussian is
as follow in this case:

You can model it as **single**
gaussian ($m=30, n=100$).

$$\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

MLE:

$$\boldsymbol{\mu} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)}$$

$$\boldsymbol{\Sigma} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \boldsymbol{\mu})(\mathbf{x}^{(i)} - \boldsymbol{\mu})^T$$

If $m \leq n$ then $\boldsymbol{\Sigma}$ matrix will be singular | Non-invertible

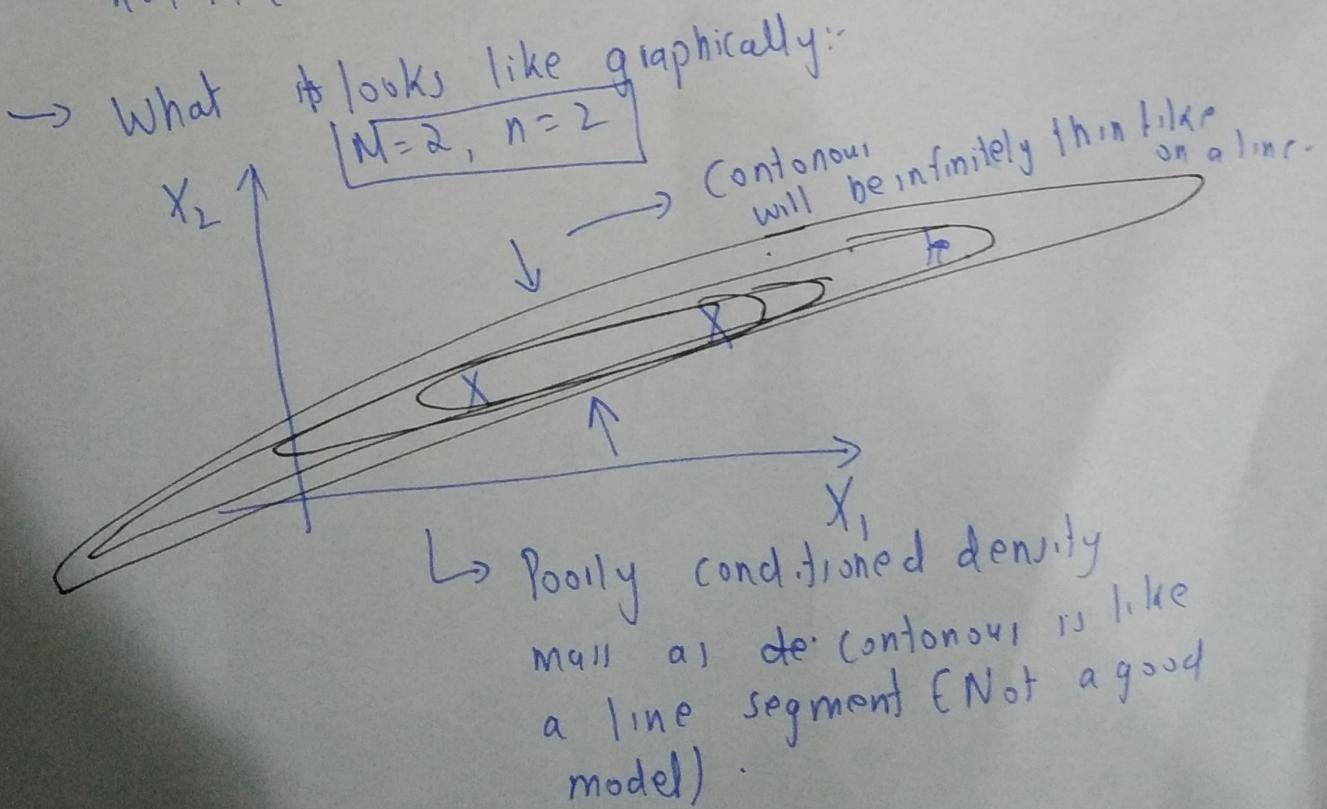
$$\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} \text{tr}(\Sigma^{-1} \mathbf{x}\mathbf{x}^T)\right)$$

Being singular \rightarrow
this term will be
 $2^{e^{10}}$.

Being singular \rightarrow
this will be
undefined.

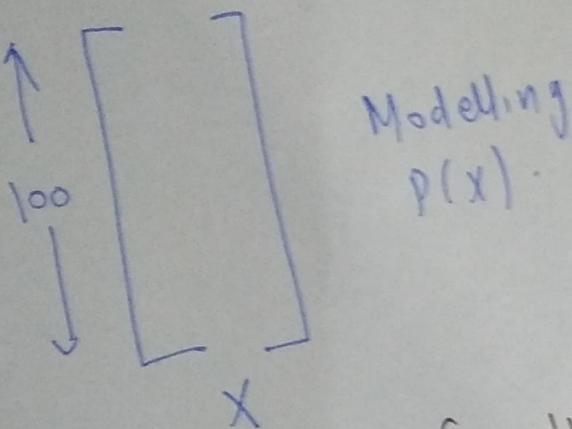
(when $m \ll n$)

- ④ So the problem in applying mixture of gaussian model is that you will end up with Σ that is singular hence non-invertible.



~ **Factor Analysis** to the rescue in
this case: ($m < n$)

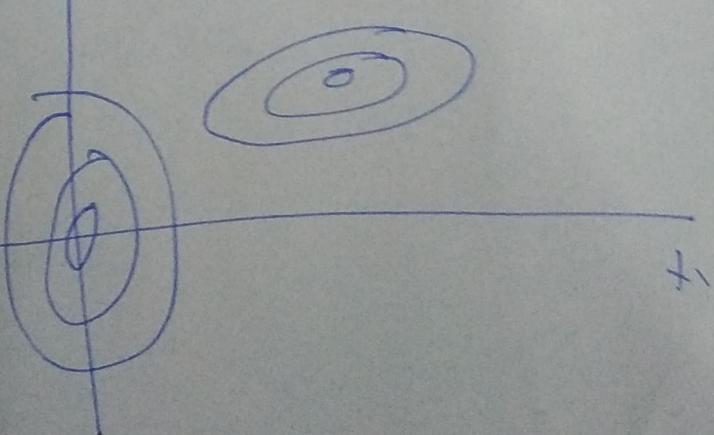
one case of $m < n$:
 ← [→ 100 psychological attributes
 → 30 persons]



~ Some other alternatives for the case $m < n$:
 Ⓛ **Option 1:** Constraint Σ to be diagonal

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & & & \\ & \sigma_2^2 & & & \\ & & \ddots & & \\ & & & & \sigma_n^2 \end{bmatrix}$$

Example of contours with Σ diagonal matrix with off diagonal 0's.



* With Σ as diagonal matrix

MLE:

$$\hat{\sigma}_j^2 = \frac{1}{m} \sum_i (x_j^{(i)} - \mu_j^{(i)})^2.$$

$\rightarrow \Sigma$ has n diagonal parameters

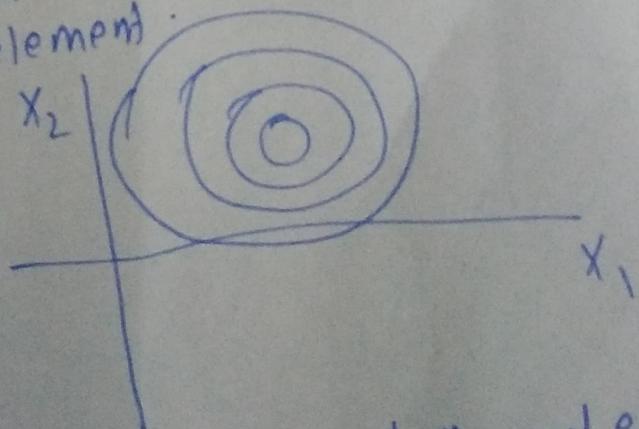
\rightarrow The problem with this ^{optimal} is that it assumes that features are uncorrelated.

② Option 2:- Constraint Σ to be

$$\Sigma = \sigma^2 I$$

$$\text{i.e. } \Sigma = \begin{bmatrix} \sigma^2 & & & \\ & \ddots & & 0 \\ & & \ddots & \\ 0 & & & \sigma^2 \end{bmatrix}$$

\rightarrow Constraining Σ matrix not only to be diagonal but also have the same diagonal element.



\rightarrow Very Strong Assumption and poor

$$\sigma^2 = \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^{n_j} (x_{ij} - \mu_{ij})^2$$

→ We won't use any of this options because of their strong assumptions but instead go with Factor Analysis which captures correlation and does not go in problem of invertibility.

→ Factor Analysis explanation:-

$$P(Y, Z) = P(Y|Z) \cdot P(Z)$$

Z is a hidden variable

$$\begin{aligned} d &= 3 \\ n &= 100 \\ m &= 30 \end{aligned}$$

$$Z \sim N(0, I), \quad z \in \mathbb{R}^d \quad (d < n)$$

$$X = \mu + Az + \varepsilon \quad (1)$$

$$\text{where } \varepsilon \sim N(0, E)$$

Parameters:

$$\mu \in \mathbb{R}^n, A \in \mathbb{R}^{n \times d}, \Psi \in \mathbb{R}^{d \times d \text{ diagonal}}$$

(166) → The expression ① can be equivalently written as:

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

→ z is basically the power full which we assume is deriving the features of

x .

→ So, x is basically the linear combination of mean and the power ful which we assume is deriving it.

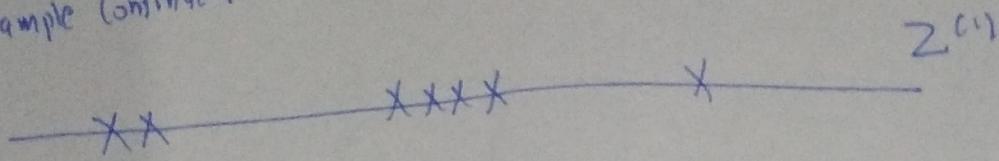
→ Ψ being a diagonal covariance matrix assumes that noise added at each feature of x is independent of the noise of other feature. (Each feature has independent noise).

→ Examples of data which factor analysis can model:

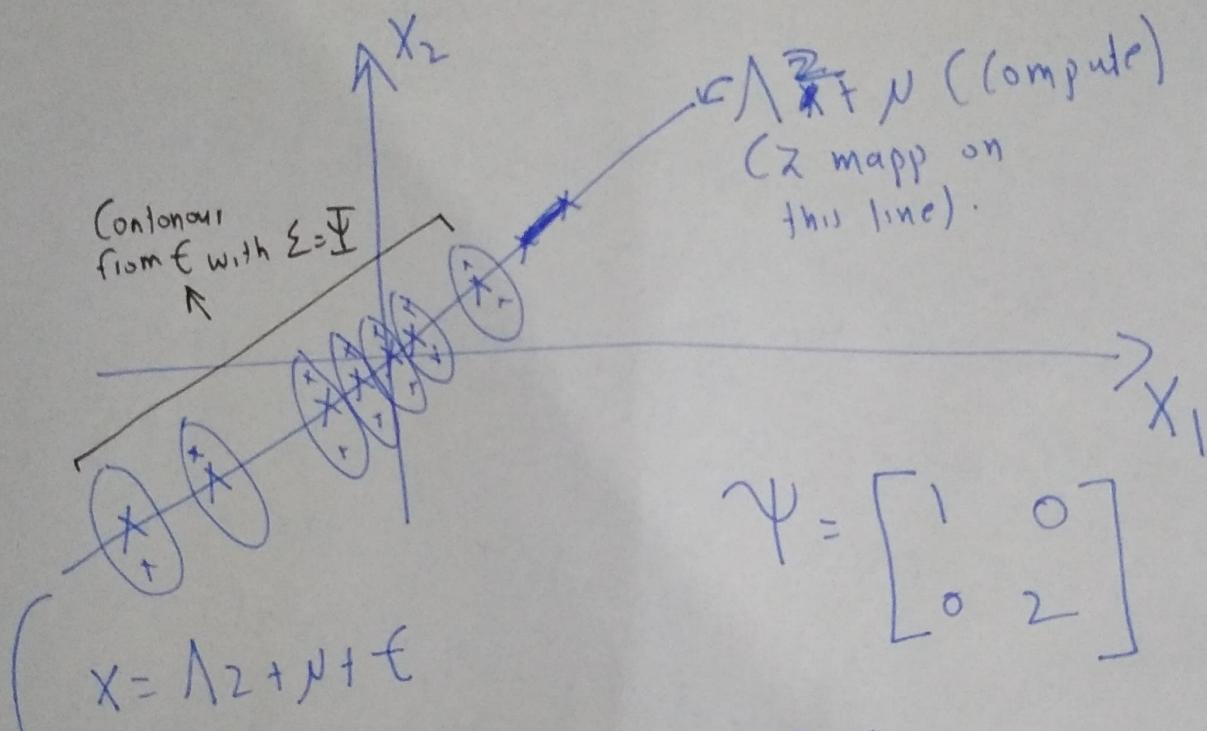
$$\textcircled{1} \quad z \in \mathbb{R}^d \times \mathbb{R}^m \quad d=1 \\ n=2 \\ m=7$$

Example (Continued)

(161)

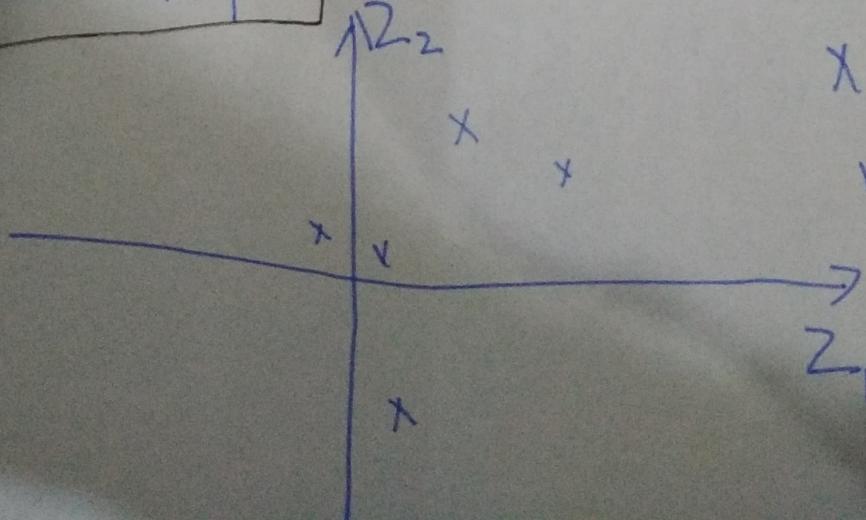


$$\text{Say } \Lambda = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \quad N = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



→ Samples from each of the
contours (Normal distribution $X|z$) .

→ Another Example:



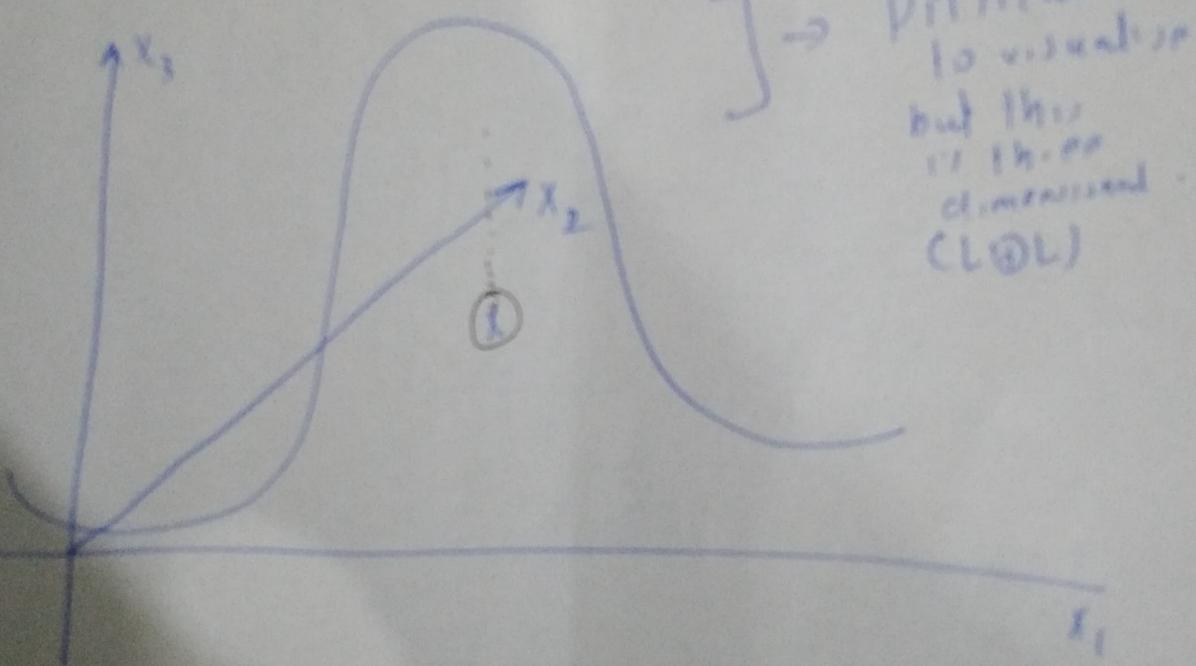
$$Z \in \mathbb{R}^2 \quad d=2$$
$$X \in \mathbb{R}^3 \quad n=3$$
$$m=5 \quad M=5$$

→ Now we will compute $A_2 + \mu$ Mapping to higher dimensional space

$$\begin{pmatrix} A_2 & \mu \\ 1 & 3x \\ 3 & 2x \\ 2 & 1 \end{pmatrix}$$

→ Now we are going to take 2 and map to higher dimensional space.

→ Now each point in three dimensional i.e. all five points would have a gaussian bump.



- ④ Factor analysis can take very high dimensional data and model the data roughly lying on lower dimension subspace with a little bit of fuzz on that lower dimension

* How to derive ε -step and M-step of Factor Analysis and model grouping
to be tricky):

→ Multivariate Gaussian (Proposition)

Partitioned
vector!!!

$$\begin{bmatrix} X \\ X_2 \end{bmatrix}$$

$$\uparrow r$$

$$\downarrow s$$

$$X, F R^r, X_2 \perp R^s$$

$$X | F R^{r+s}$$

So:

$$X \sim N(\mu, \Sigma)$$

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \uparrow r \quad \uparrow s \rightarrow \text{Partitioned}$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \uparrow r \quad \uparrow s$$

Marginal $P(X_1) = ?$

$$P(X) = P(X_1, X_2)$$

$$\int_{X_2} P(X_1, X_2) dX_2 = P(X_1)$$

$$\int_{X_2} \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} \begin{pmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \end{pmatrix}^\top \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \begin{pmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \end{pmatrix}\right)$$

→ Of course not going to integrate
but results will be as follows:

$$P(X_1) = X \sim N(\mu_1, \Sigma_{11})$$

(Not shocking)

Conditional: $P(X_1 | X_2) = ?$

$$\frac{P(X_1, X_2)}{P(X_2)} = P(X_1 | X_2)$$

$$\leftarrow X_1 | X_2 \sim N(\mu_{1|2}, \Sigma_{1|2})$$

Simplifying
the above expression:

where

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (X_2 - \mu_2)$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

★

* Using these above properties of multivariate gaussian, we will derive EM go through how you derive EM algorithm for factor analysis :-

① Derive $P(X, Z)$

Stacking them up $\begin{bmatrix} Z \\ X \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}_{x,z}, \Sigma)$

$$\left. \begin{array}{l} Z \sim \mathcal{N}(0, I) \\ X = \boldsymbol{\mu} + \Lambda Z + \boldsymbol{\varepsilon} \end{array} \right\} \rightarrow \begin{array}{l} \text{Original} \\ \text{definition} \end{array}$$

$$E_Z = 0, E_X = \boldsymbol{\mu} + \Lambda Z + \boldsymbol{\varepsilon}^T \\ = \boldsymbol{\mu}$$

therefore

$$\boldsymbol{\mu}_{x,z} = \begin{bmatrix} 0 \\ \boldsymbol{\mu} \end{bmatrix} \begin{matrix} \uparrow d \\ \uparrow n \end{matrix}$$

Similarly: You can compute:-

partitioned matrix Σ \downarrow $n \times n$ $\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$

$$\Sigma = \begin{bmatrix} E(z-E_2)(z-E_2)^T & E(z-E_2)(x-E_x)^T \\ E(x-E_x)(z-E_2) & E(x-E_x)(x-E_x)^T \end{bmatrix}$$

→ From above you can derive $\Sigma_{11}, \Sigma_{12}, \Sigma_{21}, \Sigma_{22}$ one by one: lets do Σ_{22} (rest are similar)

$$\Sigma_{22} = E(x-E_x)(x-E_x)^T$$
$$= E[(\Lambda z + \epsilon - \bar{\lambda})(\Lambda z + \epsilon - \bar{\lambda})^T] \quad \text{[From previous]}$$

= Doing the Quadratic expansion

$$= E[\Lambda z \Lambda^T + \Lambda \epsilon \epsilon^T + \epsilon \Lambda^T + \epsilon \epsilon^T]$$

$$= E[\Lambda z \Lambda^T + \Sigma \Sigma^T]$$

$$= E[\Lambda z z^T] + E[\varepsilon \varepsilon^T]$$

$$= \Lambda E[z z^T] \Lambda^T + \Psi$$

$\Sigma_{22} = \Lambda \Lambda^T + \Psi$ (Lower Right Block)
 \rightarrow Similarly other blocks are as follow:

$$\Sigma = \begin{bmatrix} I & \Lambda^T \\ X & \Lambda \Lambda^T + \Psi \end{bmatrix}$$

\rightsquigarrow So we have figure out

$$\begin{bmatrix} z \\ X \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ X & \Lambda \Lambda^T + \Psi \end{bmatrix}\right)$$

$P(X^{(i)}) \rightarrow$ You could write down

$P(X^{(i)})$ and ~~then~~ find derivative of log-likelihood of above expression then set it to zero and then find the parameter but turns out there is no closed form solution.

T74

→ In order to fit we are going
to resort to EM:

→ **E-step**

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

(But $z^{(i)}$ is
continuous density)

So

$$\leftarrow z^{(i)} | X^{(i)} \sim N(\mu_{z^{(i)} | X^{(i)}}, \Sigma_{z^{(i)} | X^{(i)}})$$

~~Through~~
where

$$\mu_{z^{(i)} | X^{(i)}} = \vec{\Omega} + \Lambda^T (\Lambda \Lambda^T + \bar{E})^{-1} (X^{(i)} \cdot \nu)$$

$$\Sigma_{z^{(i)} | X^{(i)}} = I - \Lambda^T (\Lambda \Lambda^T + \bar{E})^{-1} \Lambda$$

From previous equations

↳ In e-step you
compute these.

→ M-step:

→ Derivation of M-step is
quite long & complicated but
I would mention just a
key algebraic trick you need
to use when deriving M-step.

$$Q_i(z^{(i)}) = \frac{1}{2(\pi)^{d/2}} \exp\left(-\frac{1}{2}(\dots)\right)$$

↓ from
 above
 expression

↑ from
 above
 expression

$$= \int_{z^{(i)}} Q_i(z^{(i)}) z^{(i)} dz^{(i)} \quad (\text{There will be steps where you need to compute this.})$$

$$= \int_{z^{(i)}} \frac{1}{2\pi^{d/2}} |\Sigma| \exp\left(-\frac{1}{2}(\dots)\right) z^{(i)} dz^{(i)} \quad (\text{There is much simpler})$$

→ There much simple way to compute this integral:-

$$\int_{z^{(i)}} Q_i(z^{(i)}) z^{(i)} dz^{(i)} \quad (\text{This is expected value})$$

$$E_{z^{(i)} \sim Q_i} [z^{(i)}] = \boxed{\mu_{z^{(i)} | x^{(i)}}} \rightarrow \text{We have already computed this}$$

→ M-step after derivations would be
a) follow

$$\Theta_i: \arg \max \sum_i \int_{Z^{(i)}} Q(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)})}{Q(z^{(i)})} dz^{(i)}$$

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

Plug in
gaussian
density.

After plugging in
gaussian density and then
you take the derivative of above
expression with respect to
the parameter set to zero and
find the parameters

$$\nabla_{\theta} (\dots) = 0 \quad (\text{and solve})$$