

- Some properties of MVNs
- Review of EM algo. + GMMs
- Factor analysis model
- Solving the FAM w/ EM
- Comparing FA to PCA

## The Factor Analysis Model

### Multivariate Normals (Facts about MVNs w/o proofs)

- For a random vector,  $X \in \mathbb{R}^n$  that is Gaussian distributed with an invertible covariance matrix,  $\Sigma \in \mathbb{R}^{n \times n}$  (w/  $\Sigma$  symmetric PSD), its PDF is given by:

$$P(X; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (X-\mu)^T \Sigma^{-1} (X-\mu)\right\},$$

w/  $\mu \in \mathbb{R}^n$ . Note that for a Gaussian random vector,  $\Sigma$  may not be invertible, in which case we say the Gaussian is degenerate (as we will see later). We write this as:

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



It is sometimes convenient to write this in block (or partitioned form):

$$\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right),$$

$\begin{matrix} \text{r} & \text{r}-1 & \text{r}-s-1 \\ \text{r} & \text{r}-1 & \text{r}-s-1 \end{matrix}$

$X_1 \in \mathbb{R}^r, X_2 \in \mathbb{R}^s$

where  $r+s=n$ . This is just convenient shorthand, eg:

$$X = \begin{bmatrix} X_1^T \\ \vdots \\ X_r^T \\ X_{r+1}^T \\ \vdots \\ X_n^T \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}, \text{ and } \Sigma_{11} = \begin{bmatrix} \sigma_{11}^2 & \dots & \sigma_{1r}^2 \\ \vdots & \ddots & \vdots \\ \sigma_{r1}^2 & \dots & \sigma_{rr}^2 \end{bmatrix}. \quad \text{also}$$

$\Sigma_{21} = \Sigma_{12}^T$  since we know  $\Sigma$  must be symmetric.

It is not hard to show that multiplication w/ block vectors/matrices behaves the same way as if you treated them as scalars in normal matrix multiplication.

marginal and conditional dist.

The expression above represents the joint density of  $X_1, X_2$ :  $P(X_1, X_2) = P(X_1^1, \dots, X_1^r, \dots, X_2^1, \dots, X_2^s)$ . It can be shown that the marginal density of  $X_1$  ( $P(X_1) = \int P(X_1, X_2) dX_2$ ) and  $X_2$  are also normal:

$$\{X_1 \sim \mathcal{N}(\mu_1, \Sigma_{11})\}; \{X_2 \sim \mathcal{N}(\mu_2, \Sigma_{22})\}$$



and the conditional dists  $\left( p(x_1|x_2) = \frac{p(x_1, x_2)}{\int_{x_1 \in \mathbb{R}} p(x_1, x_2) dx_1} \right)$  (3)

are normal as well.

$$\boxed{\begin{aligned} X_1 | X_2 &\sim \mathcal{N}(\mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}) \\ X_2 | X_1 &\sim \text{as above w/ } 1 \leftrightarrow 2. \end{aligned}}$$

This is not difficult to show and relies on the matrix form of the completion of squares trick to turn quadratics in the exponential to the familiar gaussian form:

$$e^{ax^2 + bx + c} = e^{a(x-h)^2 + h} \propto e^{a(x-h)^2}.$$

the sum of indep. normally dist. rvs is also normal.

$$\text{if } X \sim \mathcal{N}(\mu_x, \Sigma_x)$$

$$Y \sim \mathcal{N}(\mu_y, \Sigma_y)$$

and  $Z \equiv X + Y$ , then

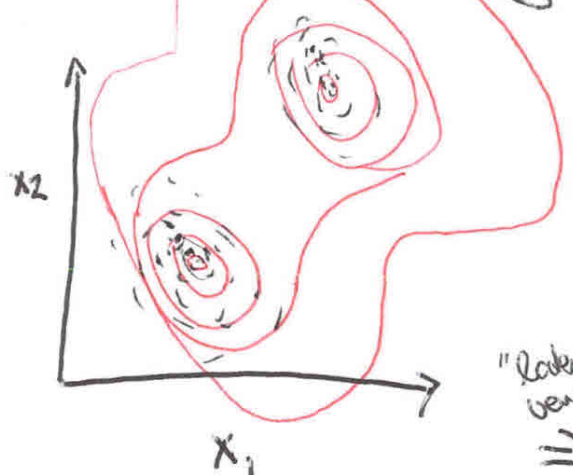
$$\boxed{Z \sim \mathcal{N}(\mu_x + \mu_y, \Sigma_x + \Sigma_y)}$$



- Finally, MNV ~~and~~ rvs under an affine transformation<sup>(4)</sup> are also normal (affine property):

- let  $X \in \mathbb{R}^n$  be  $X \sim \mathcal{N}(\mu, \Sigma)$  and  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ ,  
then  $\{AX + b \sim \mathcal{N}(A\mu + b, A\Sigma A^T)\}$

## GMM and EM algo



### model

① Flip an unfair coin to determine class

② draw a rv.  $\mathcal{I}$  from the appropriate Gaussian

"Latent" variable

$$\Rightarrow Z^{(i)} \sim \text{Bern}(\phi)$$

$$X^{(i)} | Z^{(i)} = j \sim \mathcal{N}(\mu_j, \Sigma_j)$$

(or multinomial: depending on # of classes)

Can use for:

1) estimation of density field

2) soft clustering of data, where we obtained soft guesses on cluster assignment through inference:

$$P(Z^{(i)} = j | X^{(i)})$$

(Not Bayesian inference though!)

- For a known (modeled) conditional  $p(Z^{(i)} | X^{(i)}, \Theta)$  and joint density dist.  $p(X^{(i)}, Z^{(i)}, \Theta)$ , we may use the EM algo. to find the MLE of  $\Theta$ ,  $\Theta^*$ :

EM Algo.

① initialize  $\Theta^*$

② repeat til convergence {

(e-step) 2a)  $Q_i(Z^{(i)}) := p(Z^{(i)} | X^{(i)}, \Theta^*)$

(m-step) 2b)  $\Theta^* := \arg \max_{\Theta}$

$$\sum_i \sum_{Z^{(i)}} Q_i(Z^{(i)}) \log \frac{p(X^{(i)}, Z^{(i)}, \Theta)}{Q_i(Z^{(i)})}$$

}

- note that if  $Z^{(i)}$  is continuous, then  $\sum_{Z^{(i)}} \rightarrow \int_{Z^{(i)}}$

- note that we usually will solve for 2b with the standard calculus maximization technique.

# Factor Analysis $\mathcal{D} = \{x^{(1)}, \dots, x^{(m)}\}; x^{(i)} \in \mathbb{R}^n$

⑥

- Provides:

- ① density estimation (also provides a valid/invertible cov. matrix if needed)
- ② dimensionality reduction (mostly used for this)

I will motivate the problem through density estimation, and then move to Feature red.  
Suppose, like GMM we want to fit a MVN, but  $m < n$ . We know this will be problematic b/c it will be hard to fit a pdf when we have very few data points/unit volume.

- In fact, we can't even construct a MVN <sup>using MLE, in which case the MLE of  $\Sigma$  = sample covariance</sup> if we use the sample covariance as its covariance b/c it is not invertible, which we need for the MVN pdf.

eg: 2 points in  $\mathbb{R}^3$

$$\hat{\Sigma} = \frac{1}{m} X^T X = \frac{1}{m} \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = \frac{1}{m} \begin{bmatrix} 17 & 22 & 27 \\ 22 & 29 & 36 \\ 27 & 36 & 45 \end{bmatrix}$$

but,  $-\begin{bmatrix} 17 \\ 22 \\ 27 \end{bmatrix} + 2\begin{bmatrix} 22 \\ 29 \\ 36 \end{bmatrix} = \begin{bmatrix} 27 \\ 36 \\ 45 \end{bmatrix} \Rightarrow$  cols of  $\hat{\Sigma}$  not lin. indep.  
and thus  $\hat{\Sigma}$  not invertible.

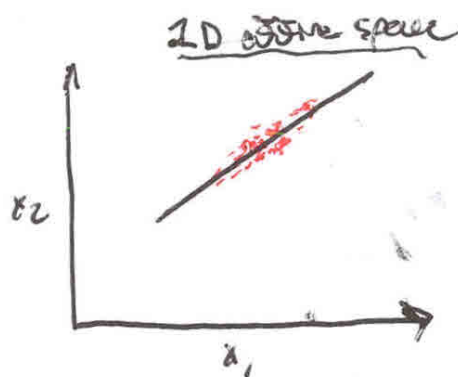
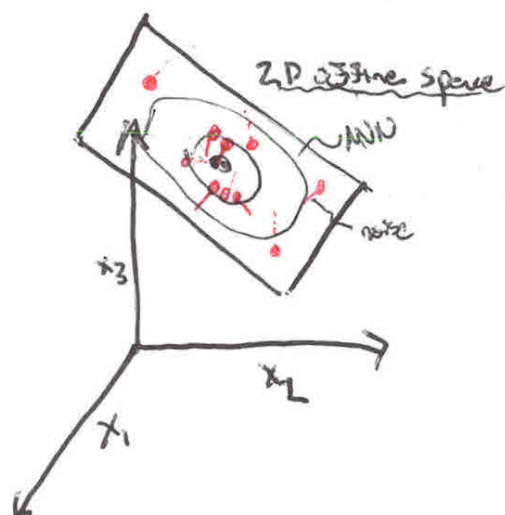




# Solution

(7)

- Fit a MVN to much smaller portion of the space, i.e. to an affine subspace (w/ a degenerate Gaussian), and model deviation from the affine space w/ Gaussian noise (off that plane). (We can easily fit on affine plane to the data since for an m-dim. plane, the data will go right through the plane).
- In fact, even when  $m > n$ , if we suspect our data lie approx. on an affine subspace w/ a Gaussian dist we can use this model.

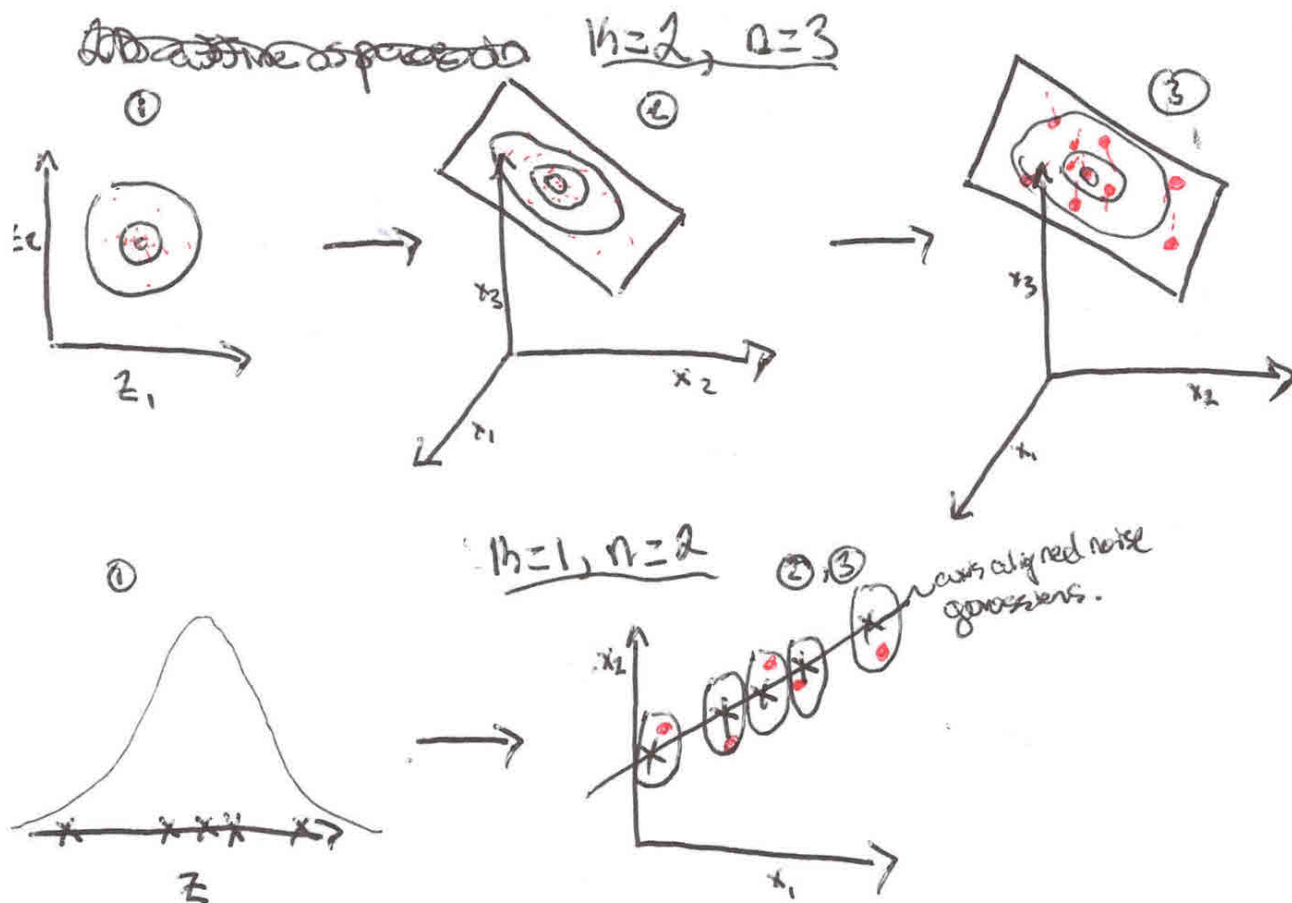


- In Math:
- ①  $\tilde{Z}^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$  ;  $\tilde{Z} \in \mathbb{R}^k$  ;  $k \leq n$    
 (latent variable)
  - ②  $\mu + \Lambda \tilde{Z}^{(i)} \sim \mathcal{N}(\mu, \Lambda \Lambda^T)$  ,   
  $\mu \in \mathbb{R}^n$  ;  $\Lambda \in \mathbb{R}^{n \times k}$  (by affine property)
  - ③  $X^{(i)} \equiv \mu + \Lambda \tilde{Z}^{(i)} + \epsilon$  , w/  $\epsilon$  axis aligned Gaussian noise   
 ( $\epsilon \sim \mathcal{N}(\mathbf{0}, \Psi)$ )   
  $\Psi \in \mathbb{R}^{n \times n}$  (diagonal)
- noise is necessary because we typically have  $k \ll m$ , so that only some will not exactly fit through plane*

In words:

- ① draw a MVN  $\mu$  in  $\mathbb{R}^h$  from a  $0, I$  gaussian
- ② map that vector to the affine space in  $\mathbb{R}^n$   
by multiplying by  $\Delta$  and adding  $\mu$  (dist. ~~is~~  
will still be Gaussian on the affine subspace)  
(by affine property)
- ③ add axis aligned noise (which can be thought of  
as measurement noise)

In pictures



→



- we solve for  $\mu, \Lambda, \Psi$  using their max. likelihood estimates, and we therefore need to derive a few more dists. to write down  $\ell$ .

- Another way to write this process is that:

$$\begin{cases} Z^{(i)} \sim \mathcal{N}(0, I) & (\text{noisy. of } Z) \\ X^{(i)} | Z^{(i)} \sim \mathcal{N}(\mu + \Lambda Z^{(i)}, \Psi) & (\text{cond. of } X) \end{cases}$$

(much like GMM)

Since sum of 2 indep Gaussian RV is Gaussian:

$$X^{(i)} \sim \mathcal{N}(\mu, \Lambda \Lambda^T + \Psi) \quad (\text{noisy. of } X)$$

(this thus shows us that we are modeling the observed density w/ a MVN)

cov. from  $\Lambda$

covariance from  $\mu + \Lambda Z$

- joint is also Gaussian

$$\begin{bmatrix} Z \\ X \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu_Z \\ \mu_X \end{bmatrix}, \begin{bmatrix} \Sigma_{ZZ} & \Sigma_{ZX} \\ \Sigma_{XZ} & \Sigma_{XX} \end{bmatrix} \right)$$

known from marginals

Since:

$$\begin{aligned} X^{(i)} | Z^{(i)} &\sim \mathcal{N}(\mu_X + \Sigma_{XZ} \Sigma_{ZZ}^{-1} (Z - \mu_Z), \Sigma_{XX} - \Sigma_{XZ} \Sigma_{ZZ}^{-1} \Sigma_{ZX}) \\ &\sim \mathcal{N}(\mu_{\text{obs}} + \Sigma_{XZ} \mathbf{I}^{-1}(Z), \Lambda \Lambda^T + \Psi - \Sigma_{XZ} \mathbf{I}^{-1} \Sigma_{XZ}^T) \\ &\sim \mathcal{N}(\mu + \Sigma_{XZ} Z, \Lambda \Lambda^T + \Psi - \Sigma_{XZ} \Sigma_{XZ}^T) \end{aligned}$$

$\Rightarrow$  Comparing this dist. to the cond. dist. For  $X$ , we see that:

$$\Sigma_{XZ} = \Lambda, \text{ so:}$$

$$\begin{bmatrix} Z^{(i)} \\ X^{(i)} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ \mu \end{bmatrix}, \begin{bmatrix} I_h & \Lambda^T \\ \Lambda & \Lambda \Lambda^T + \Psi \end{bmatrix} \right) \quad (\text{joint}).$$

We may also easily write down the cond. of  $Z^{(i)}$ :

$$Z^{(i)} | X^{(i)} \sim \mathcal{N}(\mu_{Z^{(i)} | X^{(i)}}, \Sigma_{Z^{(i)} | X^{(i)}}) \quad (\text{cond. of } Z)$$

$$\omega / \mu_{Z^{(i)} | X^{(i)}} = \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} (X^{(i)} - \mu)$$

$$\Sigma_{Z^{(i)} | X^{(i)}} = I_h - \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} \Lambda$$

### dim. reduction

- recall from the GMM, the latent variables represented the unobserved class labels for each data point.
- to obtain class probabilities for each data point, we used inference by calculating  $P(Z^{(i)} = j | X^{(i)})$  for each data point, for each class,  $j$ .

(we got this for free in the E-step)

- in FA, we posit that there is some unobserved vector (of unobserved features),  $z^{(i)}$  (w/  $\dim = k < n$ ) that fundamentally describes each data point (like how the unobserved class in the GMM fundamentally describes each data point).

- we again use inference (note this is not Bayesian inference). Our prior ~~on~~ on  $z^{(i)}$  is  $N(0, I_k)$  and ~~we obtain the posterior~~ the posterior,  $z^{(i)} | x^{(i)}$  is the <sup>where now the posteriors ~~are~~ not int</sup> MVN w/ mean  $\mu_{z^{(i)} | x^{(i)}}$ ,  $\Sigma_{z^{(i)} | x^{(i)}}$ . We use the values of  $\Lambda, \mu, \Phi$  obtained from MLE.

- We have a dist. over the posterior of  $z^{(i)}$ . To ascribe a single value for this data point, we <sup>(as we did in the GMM)</sup> therefore use the mode of the posterior.

$$x_{FA}^{(i)} = \mu_{z^{(i)} | x^{(i)}} = \hat{\Lambda}^T (\hat{\Lambda} \hat{\Lambda}^T + \hat{\Phi})^{-1} (x^{(i)} - \hat{\mu})$$

where the hats denote the MLEs.

# Solving the MLE w/ EM

The log-likelihood is:

$$l(\mu, \Lambda, \Psi) = \log p(X^{(1)}, \dots, X^{(n)}; \mu, \Lambda, \Psi)$$

$$= \log \prod_{i=1}^n p(X^{(i)}; \mu, \Lambda, \Psi)$$

$$= \log \prod_{i=1}^n \frac{1}{(2\pi)^{n/2} |\Lambda\Lambda^T + \Psi|^{1/2}} \exp\left\{-\frac{1}{2} (X^{(i)} - \mu)^T (\Lambda\Lambda^T + \Psi)^{-1} (X^{(i)} - \mu)\right\}$$

$$= -\sum_{i=1}^n \log((2\pi)^{n/2}) + \log |\Lambda\Lambda^T + \Psi|^{1/2}$$

$$+ \frac{1}{2} (X^{(i)} - \mu)^T (\Lambda\Lambda^T + \Psi)^{-1} (X^{(i)} - \mu)\}$$

which is intractable to maximize in closed form w/ calculus techniques. We therefore appeal to the EM algo:

E-step:  $Q_i(Z^{(i)}) = p(Z^{(i)} | X^{(i)}; \mu^*, \Lambda^*, \Psi^*) =$

we derived this cond. dist. before

$$\frac{1}{(2\pi)^{n/2} |\Sigma_{Z^{(i)}|X^{(i)}}|^{1/2}} \exp\left\{-\frac{1}{2} (Z^{(i)} - \mu_{Z^{(i)}|X^{(i)}})^T \Sigma_{Z^{(i)}|X^{(i)}}^{-1} (Z^{(i)} - \mu_{Z^{(i)}|X^{(i)}})\right\}$$

w/  $\mu_{Z^{(i)}|X^{(i)}} \equiv \Lambda^{*T} (\Lambda^* \Lambda^{*T} + \Psi^*)^{-1} (X^{(i)} - \mu)$   
 $\Sigma_{Z^{(i)}|X^{(i)}} \equiv I_n - \Lambda^{*T} (\Lambda^* \Lambda^{*T} + \Psi^*)^{-1} \Lambda^*$

- The M step update is more difficult, but doable. (13)

M-step:

$$\Lambda^*, \mu^*, \Psi^* = \underset{\substack{\Lambda \in \mathbb{R}^{n \times k} \\ \mu \in \mathbb{R}^n \\ \Psi \in \mathbb{R}^{n \times n} (\text{diag})}}{\operatorname{argmax}} \sum_{i=1}^m \int_{Z^{(i)} \in \mathbb{R}^k} Q_i(Z^{(i)}) \log \frac{P(X^{(i)}, Z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(Z^{(i)})} dZ^{(i)}$$

- Expressions for  $\Lambda^*, \mu^*, \Psi^*$  can be found in closed form by taking gradients,  $\nabla_{\Lambda}, \nabla_{\mu}, \nabla_{\Psi}$ , setting the expression to the zero vector/matrix, and solving (using several matrix and trace identities).

### - Connection between FA and PCA

PCA: geometric approach

- directly project data down to subspace of  $k$  eigenvectors of sample cov.

FA: probabilistic approach

- posit the existence of a  $k$  dim. latent variable and perform inference to estimate its value.

- Recall 2 assumptions made by PCA

- ① Principle axes are  $\perp$
- ② principle axes are linear

$\Rightarrow$  data distributed as an ellipsoidal blob. I.e., data distr. in  $\mathbb{R}^n$  as a MVN (just like the marginal of  $x$  in FA).

- In fact if  $\Psi = \sigma^2 I$  (isotropic noise) we may solve for  $\hat{\Lambda}$  analytically. In limit  $\sigma \rightarrow 0$  (i.e., data become deterministic), the soln. of for  $X_{FA}$  is identical to that of PCA. (called probabilistic PCA)

- I.e., PCA is just the deterministic limit of a special case of FA.