

# The EM Algorithm

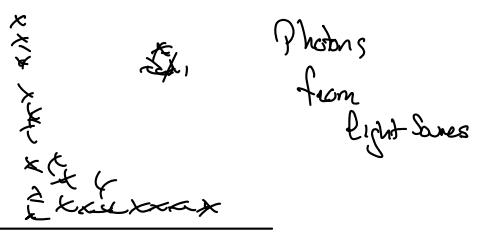
Recall Unsupervised Algorithms

K-means  $\neq$  Gmm

Big idea latent variable  $\rightarrow$  "fraction of points from a source"

Algorithm "guess" latent variable

Estimate other parameters (center and shape)

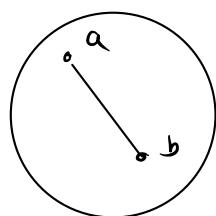


Today: EM Algorithm for latent variable models

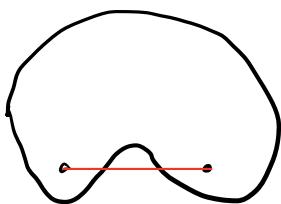
- + Technical Detour: Convexity  $\neq$  Jensen's Inequality
- + EM Algorithm as MLE
- + Gmm as EM Algorithm
- + Factor Analysis (EM for a new setting)

Detour Convexity  $\nsubseteq$  JENSEN (This is a key result, we'll go slowly)

A SET  $\Omega$  IS CONVEX if for any  $a, b \in \Omega$  the line joining  $a, b$  is in  $\Omega$  as well.



Convex



NOT convex!

IN symbols,

$$\forall \lambda \in [0,1], a, b \in \Omega$$

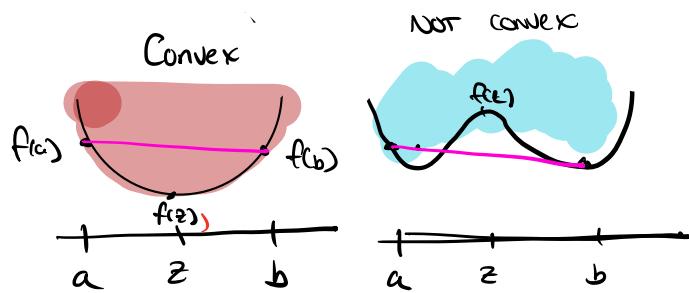
$$\lambda a + (1-\lambda)b \in \Omega$$

(NEED TO CHECK  $\lambda a + (1-\lambda)b \in \Omega$ )

GIVEN a function  $f$ , the graph of  $f$   $G_f$  is defined as

$$G_f = \{ (x, y) : y \geq f(x) \}$$

A function is convex if its graph is convex (as a set)



IN symbols,  $\forall \lambda \in [0,1]$

$$\lambda f(a) + (1-\lambda)f(b) \in \Omega$$

or let  $z = \lambda a + (1-\lambda)b$

$$\lambda f(a) + (1-\lambda)f(b) \geq f(z)$$

"Every curve is above function"

If  $f$  twice differentiable,  $\forall x \quad f''(x) \geq 0 \Rightarrow f$  is convex

$$\text{pf} \quad f(z) = f(a) + f'(a)(a-z) + f''(\zeta_a)(a-z)^2 \quad \zeta_a \in [a, z]$$

$$f(b) = f(z) + f'(z)(b-z) + f''(\zeta_b)(b-z)^2 \quad \zeta_b \in [z, b]$$

$$\lambda f(a) + (1-\lambda)f(b) = f(z) + f'(z)(\lambda a + (1-\lambda)b - z) + c \quad c \geq 0$$

i.e.  $\lambda f(a) + (1-\lambda)f(b) \geq f(z) \quad \square$  def d z.

We say  $f$  is strongly convex if  $\forall x \in \text{Dom}(f) \quad f''(x) > 0$ .

Ex:  $f(x) = x^2 \Rightarrow f''(x) = 2 \Rightarrow$  Strongly convex

$f(x) = x^2(x-1)^2$  : graph above (not convex)

JENSEN'S INEQUALITY  $\mathbb{E}[f(x)] \geq f(\mathbb{E}[x])$  for convex  $f$ .

Ex:  $x$  takes value  $a$  with prob  $\lambda$

. . . takes value  $b$  with prob  $1-\lambda$

$$\mathbb{E}[f(x)] = \lambda f(a) + (1-\lambda)f(b)$$

$$f(\mathbb{E}[x]) = f(z) \quad z = \lambda a + (1-\lambda)b$$

N.B. CAN PROVE finitely  
Supported Distribution  
By induction

for convex  $f$ , definition implies this in this case!

Stronger if  $f$  is strongly convex, and  $\mathbb{E}[f(x)] = f(\mathbb{E}[x])$

$\Rightarrow x$  is a constant (except: almost surely)

WE NEED CONCAVE FUNCTIONS  $g$  concave iff  $-g$  is convex

$$\mathbb{E}[g(x)] \leq g(\mathbb{E}[x])$$

Ex:  $g(x) = \log(x) \Rightarrow g''(x) = -x^{-2}$  on  $(0, \infty)$  NEGATIVE



WHAT ABOUT  $f(x) = ax + b$  CONVEX & CONCAVE since  $f''(x) = 0$ .

END DETOUR

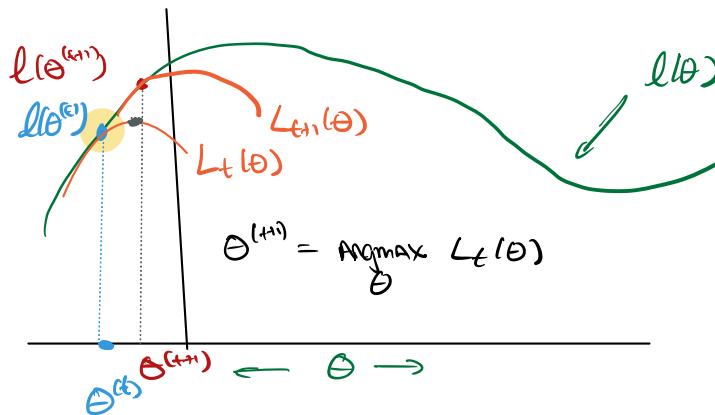
## EM Algorithm as max likelihood

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

PARAMETERS  
DATA

WE ASSUME  $P(x; \theta) = \sum_z P(x, z; \theta)$  (e.g. GMM latent variable)

## Picture of Algorithm



$$L_t(\theta) \leq l(\theta) \quad (\text{lowerbound})$$

$$L_t(\theta^{(t)}) = l(\theta^{(t)}) \quad (\text{tight})$$

Hope  $L_t$  easier to optimize than  $l(\theta)$

## Rough Algo (ITERATIVE)

(E-step) 1. Given  $\theta^{(t)}$  find  $L_t$

(M-step) 2. Given  $L_t$ , set  $\theta^{(t+1)} = \operatorname{argmax}_{\theta} L_t(\theta)$

## How do we construct $L_t$ ?

(let's look at single point term,  $\log P(x; \theta)$ )

$$\log \sum_z P(x, z; \theta) = \log \sum_z \frac{Q(z) P(x, z; \theta)}{Q(z)} \quad \text{for any } Q(z)$$

WE Pick  $Q(z)$  s.t.  $\sum_z Q(z) = 1$  AND  $Q(z) = 0 \iff$

$$= \log \mathbb{E}_{z \sim Q(z)} \left[ \frac{P(x, z; \theta)}{Q(z)} \right] \quad (\text{Symbol Pushing})$$

$$\geq \mathbb{E}_z \left[ \log \frac{P(x, z; \theta)}{Q(z)} \right] \quad \text{JENSEN?} \quad (\log \text{ is concave})$$

$$= \sum_z Q(z) \log \frac{P(x, z; \theta)}{Q(z)} \quad (\text{DEF of } \mathbb{E})$$

Key step holds for any such  $Q$ : (a)

This gives a family of lower bounds, one for each choice of  $Q$  ( $L \leq l$ )

How do we make it tight? Select  $Q$  to make Inequality tight

What if...  $\log \frac{P(x, z; \theta)}{Q(z)} = c$  for some constant, then JENSEN's is Equality!

$$P(x, z; \theta) = P(z|x; \theta) P(x; \theta)$$

so,  $Q(z) = P(z|x; \theta)$  then

$c = \log P(x; \theta)$  does not depend on  $z$ , so constant!

NB:  $Q(z)$  does depend on  $\theta + x$  - we will select a  $Q^{(i)}(z)$  for every point independently.

WE DEFINE Evidence-based Lower Bound (ELBO), sum over  $z$

$$\text{ELBO}(x, Q, z) = \sum_z Q(z) \log \frac{P(x, z; \theta)}{Q(z)}$$

WE'VE SHOWN  $\ell(\theta) \geq \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}; \theta)$  for any  $Q^{(i)}$  satisfying (a)  
lower bound

$$\ell(\theta^{(t)}) = \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}, \theta^{(t)}) \quad \text{for choice of } Q^{(i)} \text{ above.}$$

## WRAAP up

1. (E-STEP)  $Q^{(t)}(z) = P(z^{(t)})x^{(t)}; \theta)$  for  $i=1\dots n$
2. (M-STEP)  $\theta^{(t+1)} = \underset{\theta}{\operatorname{Argmax}} L(\theta)$   
in which  $L(\theta) = \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(i)}; \theta)$

WHY DOES THIS TERMINATE?  $L(\theta^{(t+1)}) \geq L(\theta^{(t)})$

IS IT GLOBALLY OPTIMAL? (NOPE, SEE PICTURE)

WE DERIVED EM ALGORITHM IN TERMS OF MLE.



## EM for Mixture of Gaussians

### RESTATE EM

(E-STEP) for  $i=1..n$  SET  $Q_i(z) = P(z^{(i)} | x^{(i)}, \theta^{(t)})$

$$\begin{aligned} (\text{M-STEP}) \quad \theta^{(t+1)} &= \underset{\theta}{\operatorname{Argmax}} \mathcal{L}_t(\theta) \\ &= \underset{\theta}{\operatorname{Argmax}} \sum_{i=1}^n \text{ELBO}(x^{(i)}, Q^{(t)}, \theta) \end{aligned}$$

DATA  $\notin$  INPUT PARAMS

WARM UP: Mixture of Gaussians. EM RECOVERS OUR AD-HOC ALGORITHM

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

$$z^{(i)} \sim \text{Multinomial}(\mathbb{I}) \quad \phi_i \geq 0 \quad \sum \phi_i = 1 \quad \text{"in cluster j"}$$

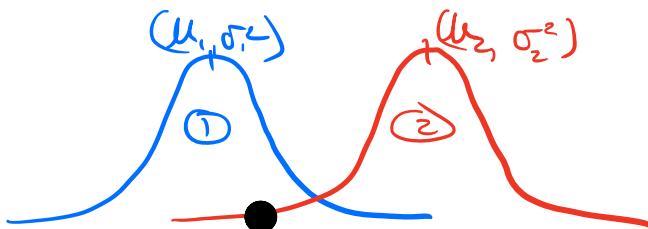
$$x^{(i)} | z^{(i)} = j \sim N(\mu_j, \sigma_j^2) \quad \text{"cluster means"}$$

$z^{(i)}$  is our LATENT VARIABLE.

### WHAT IS EM HERE?

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

WE SAW THAT COULD COMPUTE VIA Bayes Rule  $P(x^{(i)} | z^{(i)} = j; \theta)$



1. much more likely for ① than ②
  2. But if we knew,  $\phi_2 \gg \phi_1$ , maybe we'd think likely from ②
- Bayes Rule AUTOMATES this REASONING

M-STEP: Compute DERIVATIVES ..

$$\max_{\phi, \mu, \Sigma} \sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{P_i(\theta)}$$

→ WRITE  $\Theta$  FOR NOTATION ABOVE

$$\text{Rewrite } \omega_j^{(i)} \triangleq Q_i(z=j)$$

$$P(x^{(i)}, z^{(i)}; \Theta) = P(x^{(i)} | z^{(i)}) P(z^{(i)})$$

$$f_i(\Theta) = \sum_j \omega_j^{(i)} \log \left( \frac{\frac{1}{2\pi |\Sigma_j|^{1/2}} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right\}}{\omega_j^{(i)}} \cdot \phi_j \right)$$

GAUSSIAN ( $\mu, \Sigma$ )

$$\nabla_{\mu_j} f_i(\Theta) = \sum_i \nabla_{\mu_j} \left( \omega_j^{(i)} - \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right)$$

$$= -\frac{1}{2} \sum_i \omega_j^{(i)} \Sigma_j^{-1} (x^{(i)} - \mu_j) = -\frac{1}{2} \sum_j^{-1} \left( \sum_i \omega_j^{(i)} (x^{(i)} - \mu_j) \right)$$

SETTING TO 0 AND USING  $\Sigma_j^{-1}$  IS FULL RANK  $\Rightarrow \sum_i \omega_j^{(i)} (x^{(i)} - \mu_j) = 0$

$$\therefore \mu_j = \frac{\sum_i \omega_j^{(i)} x^{(i)}}{\sum_i \omega_j^{(i)}} \quad (\text{AS BEFORE})$$

$\phi_j$  is CONSTRAINED,  $\sum_j \phi_j = 1$ ,  $\phi_j \geq 0$ , NEED LAGRANGIAN

$$\nabla \phi_j = \sum_{i=1}^n \omega_j^{(i)} \nabla_{\phi_j} \log \phi_j + \nabla_{\phi_j} \lambda \left( \sum_j \phi_j - 1 \right)$$

$$= \sum_{i=1}^n \frac{\omega_j^{(i)}}{\phi_j} + \lambda = 0 \Rightarrow \phi_j = -\frac{1}{\lambda} \sum_{i=1}^n \omega_j^{(i)}$$

$$\text{SINCE } \sum \phi_j = 1, \quad \sum_j \phi_j = -\frac{1}{\lambda} \sum_{i,j} \omega_j^{(i)} = -\frac{n}{\lambda}$$

$$\therefore \phi_j = \frac{1}{n} \sum_i \omega_j^{(i)}$$

MESSAGE: EM RECOVERS GMM AUTOMATICALLY.

NB: IF  $z^{(i)}$  IS CONTINUOUS, ONE CAN REPLACE SUMS w/ INTEGRALS

## Factor Analysis

MANY fewer parts than dimensions "n < d"

cf: GMMs n > d lots of neurons, few sources.

### How does this happen?

PLACE SENSORS ALL OVER CAMPUS, Record @ 1000s of locations  
 $\Rightarrow d \approx 1000s$

But Only record for 30 days ( $n < d$ )

WANT TO FIT A DENSITY, but seems hopeless.

KEY IDEA: Assume there is some latent r.v. that  
IS NOT TOO COMPLEX and explains behavior.

1<sup>st</sup> let's see problems w/ GMMs.... Even 1 GAUSSIAN

$$\mu = \frac{1}{n} \sum_{i=1}^n x^{(i)} \rightarrow \text{this is OK}$$

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

RANK( $\Sigma$ )  $\leq n < d$  - not full rank.

Problem in GAUSSIAN likelihood

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

IS NOT DEFINED.  
 $\hookrightarrow |\Sigma| = 0$

WE will fix these issues by examining three models

that are simpler. Spoiler: we'll combine these in the end!

RECALL MLE for GAUSSIAN

$$\underset{\mu, \Sigma}{\text{MAX}} \sum_{i=1}^n \log \frac{1}{2\pi|\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}-\mu)^T \Sigma^{-1} (\mathbf{x}-\mu) \right\}$$

Equivalent

$$\underset{\mu, \Sigma}{\text{MIN}} \sum_{i=1}^n (\mathbf{x}-\mu)^T \Sigma^{-1} (\mathbf{x}-\mu) + \log |\Sigma|$$

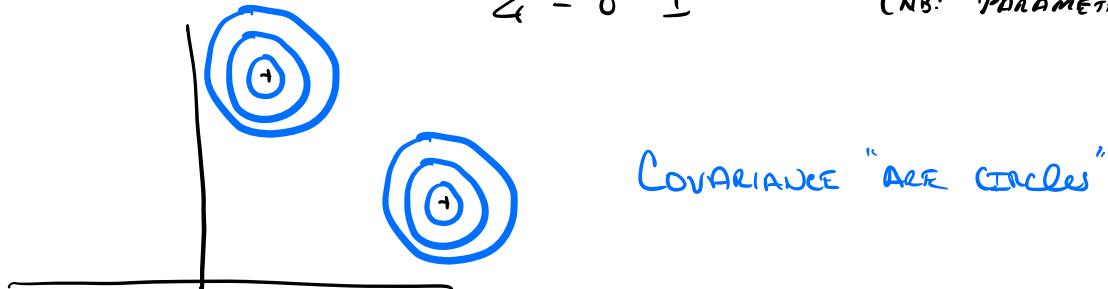
If  $\Sigma$  is full rank,  $\nabla_{\mu} = \sum_{i=1}^n \Sigma^{-1} (\mathbf{x}-\mu) = 0 \Rightarrow \mu = \frac{1}{n} \sum_i \mathbf{x}^{(i)}$

We'll use this as plugin below.

## Building Block 1

Suppose INDEPENDENT AND IDENTICAL COVARIANCE

$$\Sigma = \sigma^2 I \quad (\text{NB: PARAMETER } \sigma^2)$$



WHAT IS MLE FOR  $\Sigma$ ?

$$|\Sigma| = 2d$$

$$\underset{\sigma^2}{\text{MIN}} \sigma^{-2} \underbrace{\sum_{i=1}^n (\mathbf{x}^{(i)} - \mu)^T (\mathbf{x}^{(i)} - \mu)}_C + d \log \sigma^2$$

$$\text{let } z = \sigma^2 \quad \underset{z}{\text{MIN}} \frac{1}{z} C + d \log z$$

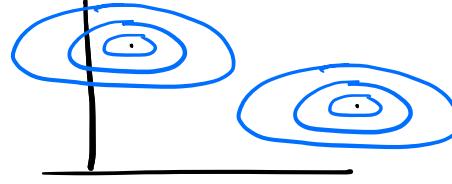
$$\Rightarrow \frac{\partial}{\partial z} = -z^{-2} C + \frac{nd}{z} = 0 \Rightarrow z = \frac{C}{nd}$$

$$\therefore \sigma^2 = \frac{1}{nd} \sum_{i=1}^n (\mathbf{x}^{(i)} - \mu)^T (\mathbf{x}^{(i)} - \mu)$$

"SUBTRACT MEAN AND SQUARE ALL ENTRIES."

## Building Block 2

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



Axis Aligned ellipse

SET  $z_i = \sigma_i^2$  (SAME IDEA AS ABOVE)

$$\min_{z_1 \dots z_d} \sum_{i=1}^n \sum_{j=1}^d z_j^{-1} (x^{(i)} - \mu_j)^2 + \log z_j$$

this is  $d$  problems for each 1 dimension

$$\Rightarrow \sum_{i=1}^n z_j^{-1} (x^{(i)} - \mu_j)^2 + \log z_j$$

$$\Rightarrow \sigma_j^2 = \frac{1}{n} \sum_i (x_j^{(i)} - \mu_j)^2$$

## Our FACTOR model

### PARAMETERS

$$\mu \in \mathbb{R}^d$$

$$\Lambda \in \mathbb{R}^{d \times s}$$

$$\Phi \in \mathbb{R}^{d \times d} \text{ - DIAGONAL MATRIX}$$

### MODEL

$$P(x, z) = P(x|z) P(z) \quad z \text{ IS LATENT}$$

$$z \sim N(0, I) \in \mathbb{R}^s \text{ for } s < d \text{ "small dim"}$$

$$x = \underbrace{\mu}_{\text{MEAN IN THE SPACE}} + \underbrace{\Lambda z}_{\text{MAPS FROM SMALL LATENT SPACE TO LARGE SPACE}} + \epsilon \quad \text{OR} \quad x \sim N(\mu + \Lambda z, \Phi)$$

$\epsilon \sim N(0, \Phi)$  Noisy

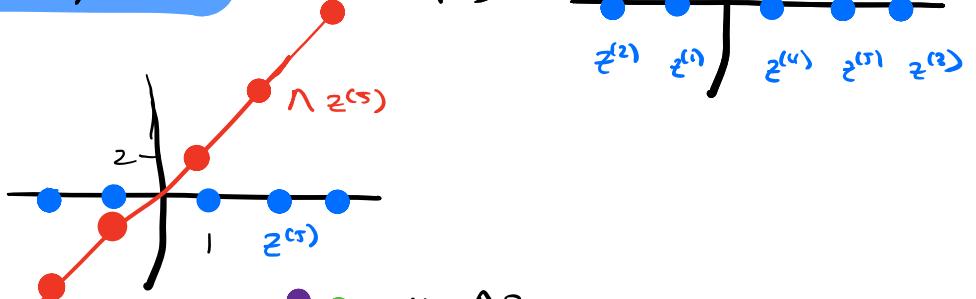
Ex:  $d=2, s=1, n=5$

$$x = \underbrace{\mu}_{\text{MEAN}} + \underbrace{\Lambda z}_{\text{MAPS FROM SMALL LATENT SPACE TO LARGE SPACE}} + \epsilon$$

1. GENERATE  $z^{(1)}, \dots, z^{(5)}$  from  $N(0, I)$

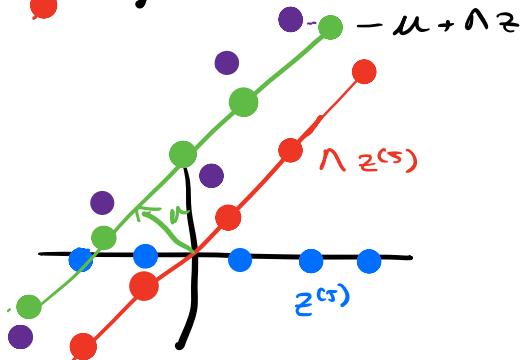


2. Suppose  $\lambda = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$



3. Add  $\mu$

4. Add  $\Sigma$   
↳ BIG SPACE



DATA WE WOULD OBSERVE ARE Purple DOTS

SO SMALL LATENT SPACE PRODUCES DATA IN HIGH DIM SPACE.

### TECHNICAL TOOLS: Block GAUSSIANS

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad x_1 \in \mathbb{R}^{d_1}, x_2 \in \mathbb{R}^{d_2} \quad x \in \mathbb{R}^d$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \quad \Sigma_{ij} \in \mathbb{R}^{d_i \times d_j} \quad i, j \in \{1, 2\}$$

NOTATION IS widely USED AND helpful.

FACT 1:  $P(x_1) = \int_{x_2} P(x_1, x_2)$

MARGINALIZATION

FOR GAUSSIANS,  $P(x_1) = N(\mu_{11}, \Sigma_{11})$  (not surprising)

FACT 2:  $P(x_1 | x_2) \sim N(\mu_{1|2}, \Sigma_{1|2})$

Conditioning

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)$$

$$\hat{\Sigma}_{12} = \hat{\Sigma}_{11} - \hat{\Sigma}_{12} \hat{\Sigma}_{22}^{-1} \hat{\Sigma}_{21} \quad (\text{matrix inversion lemma})$$

Proofs outline (happy to add)

Summary: MARGINALIZATION  $\not\equiv$  CONDITIONING GAUSSIAN  $\Rightarrow$   
 Another GAUSSIAN (CLOSED)  
 WE HAVE formula for PARAMETERS.

Back to Factor Analysis

$$x = \mu + \Lambda z + \epsilon$$

$$\begin{pmatrix} z \\ x \end{pmatrix} \sim N\left(\begin{pmatrix} 0 \\ \mu \end{pmatrix}, \Sigma\right) \quad \text{SINCE } \mathbb{E}[z] = 0$$

$$\mathbb{E}[x] = \mu$$

WHAT IS  $\Sigma$ ?

$$\hat{\Sigma}_{11} = \mathbb{E}[zz^T] = I$$

$$\begin{aligned} \hat{\Sigma}_{12} &= \mathbb{E}[z(x-\mu)^T] = \mathbb{E}[zz^T\Lambda^T] + \mathbb{E}[z\epsilon^T] \\ &= \Lambda^T \end{aligned}$$

$$\hat{\Sigma}_{21} = \hat{\Sigma}_{12}^T$$

$$\begin{aligned} \hat{\Sigma}_{22} &= \mathbb{E}[(x-\mu)(x-\mu)^T] \\ &= \mathbb{E}[(\Lambda z + \epsilon)(\Lambda z + \epsilon)^T] \\ &= \mathbb{E}[\Lambda z z^T \Lambda^T] + \mathbb{E}[\epsilon \epsilon^T] \\ &= \Lambda \Lambda^T + \Phi \end{aligned}$$

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

E-STEP :  $Q_i(z) = P(z^{(i)} | x^{(i)}; \theta)$  - USE CONDITIONAL!

M-STEP : WE HAVE CLOSED FORMS!

### Summary:

- WE SAW THAT EM CAPTURES GMM
- WE LEARNED ABOUT FACTOR ANALYSIS (LATENT LOW DIM. STRUCTURE)
- WE SAW HOW TO ESTIMATE PARAMETERS OF FA USING EM.