

# Learning From Data

## Lecture 13: Unsupervised Learning IV

Yang Li   [yangli@sz.tsinghua.edu.cn](mailto:yangli@sz.tsinghua.edu.cn)

December 17, 2021

# Today's Lecture

## Unsupervised Learning (Part IV)

- ▶ Mixture of Gaussians
- ▶ The EM Algorithm
- ▶ Factor Analysis

## Review: k-means clustering

Given input data  $\{x^{(1)}, \dots, x^{(m)}\}$ ,  $x^{(i)} \in \mathbb{R}^d$ , **k-means clustering** partition the input into  $k \leq m$  sets  $C_1, \dots, C_k$  to minimize the within-cluster sum of squares (WCSS).

$$\underset{C}{\operatorname{argmin}} \underbrace{\sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|^2}$$

## Review: k-means clustering

Given input data  $\{x^{(1)}, \dots, x^{(m)}\}$ ,  $x^{(i)} \in \mathbb{R}^d$ , **k-means clustering** partition the input into  $k \leq m$  sets  $C_1, \dots, C_k$  to minimize the within-cluster sum of squares (WCSS).

$$\operatorname{argmin}_C \sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|^2$$

↑

### Lloyd's Algorithm (1957,1982)

Let  $c^{(i)} \in \{1, \dots, k\}$  be the cluster label for  $x^{(i)}$

```

Initialize cluster centroids  $\underline{\mu_1, \dots, \mu_k} \in R^n$  randomly
Repeat until convergence{
  For every  $i$ ,
     $\underline{c^{(i)}} := \operatorname{argmin}_j \underline{\|x^{(i)} - \underline{\mu_j}\|^2}$ 

  For each  $j$ 
     $\underline{\mu_j} := \frac{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\} x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\}}$ 
}
  
```

## Review: k-means clustering

Given input data  $\{x^{(1)}, \dots, x^{(m)}\}$ ,  $x^{(i)} \in \mathbb{R}^d$ , **k-means clustering** partition the input into  $k \leq m$  sets  $C_1, \dots, C_k$  to minimize the within-cluster sum of squares (WCSS).

$$\operatorname{argmin}_C \sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|^2$$

### Lloyd's Algorithm (1957,1982)

Let  $c^{(i)} \in \{1, \dots, k\}$  be the cluster label for  $x^{(i)}$

Initialize cluster centroids  $\mu_1, \dots, \mu_k \in \mathbb{R}^n$  randomly

Repeat until convergence{

For every  $i$ ,

$c^{(i)} := \operatorname{argmin}_j \|x^{(i)} - \mu_j\|^2$  ← assign  $x^{(i)}$  to the cluster with the closest centroid

For each  $j$

$$\mu_j := \frac{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\} x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\}}$$

}

## Review: k-means clustering

Given input data  $\{x^{(1)}, \dots, x^{(m)}\}$ ,  $x^{(i)} \in \mathbb{R}^d$ , **k-means clustering** partition the input into  $k \leq m$  sets  $C_1, \dots, C_k$  to minimize the within-cluster sum of squares (WCSS).

$$\operatorname{argmin}_C \sum_{j=1}^k \sum_{x \in C_j} \|x - \mu_j\|^2$$

### Lloyd's Algorithm (1957,1982)

Let  $c^{(i)} \in \{1, \dots, k\}$  be the cluster label for  $x^{(i)}$

Initialize cluster centroids  $\mu_1, \dots, \mu_k \in \mathbb{R}^n$  randomly

Repeat until convergence{

For every  $i$ ,

$c^{(i)} := \operatorname{argmin}_j \|x^{(i)} - \mu_j\|^2$  ← assign  $x^{(i)}$  to the cluster with the closest centroid

For each  $j$

$\mu_j := \frac{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\}x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{c^{(i)}=j\}}$  ← update centroid

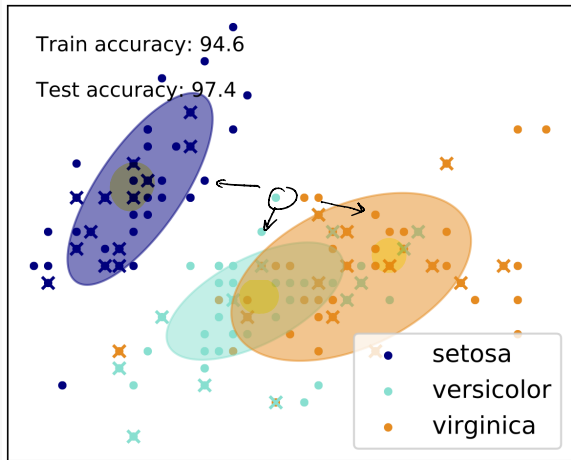
}

## Mixture of Gaussians

# Mixture of Gaussians

(GMM)

A “soft” version of k-means clustering.



Clustering results of iris dataset using *mixture of Gaussians*



# Mixture models

$$\phi \quad \text{hidden} \quad \downarrow \text{observed} \\ \phi - z \longrightarrow x.$$

## Model-based clustering

A **mixture model** assumes data are generated by the following process:

1. Sample  $\underline{z}^{(i)} \in \{1, \dots, k\}$  and  $\underline{z}^{(i)} \sim \text{Multinomial}(\underline{\phi})$   
 $\underline{p}(z^{(i)} = j) = \underline{\phi}_j$  for all  $j$ 

$\begin{bmatrix} \phi_1 \\ \vdots \\ \phi_k \end{bmatrix}$

$\underline{z}^{(i)}$  are called **latent variables**.

2. Sample observables  $x^{(i)}$  from some distribution  $p(x^{(i)}, z^{(i)})$ :

$$\underline{p}(x^{(i)}, z^{(i)}) = \underline{p}(x^{(i)} | z^{(i)}) \underline{p}(z^{(i)})$$

# Mixture models

## Model-based clustering

A **mixture model** assumes data are generated by the following process:

1. Sample  $z^{(i)} \in \{1, \dots, k\}$  and  $z^{(i)} \sim \text{Multinomial}(\phi)$

$$p(z^{(i)} = j) = \phi_j \text{ for all } j$$

$z^{(i)}$  are called **latent variables**.

2. Sample observables  $x^{(i)}$  from some distribution  $p(x^{(i)}, z^{(i)})$ :

$$p(x^{(i)}, z^{(i)}) = \underbrace{p(x^{(i)} | z^{(i)})}_{\sim \text{Bernoulli}(\theta_j)} p(z^{(i)})$$

Examples:

- ▶ Unsupervised handwriting recognition is a mixture with 10 Bernoulli distributions

# Mixture models

## Model-based clustering

A **mixture model** assumes data are generated by the following process:

1. Sample  $z^{(i)} \in \{1, \dots, k\}$  and  $z^{(i)} \sim \text{Multinomial}(\phi)$

$$p(z^{(i)} = j) = \phi_j \text{ for all } j$$

$z^{(i)}$  are called **latent variables**.

2. Sample observables  $x^{(i)}$  from some distribution  $p(x^{(i)}, z^{(i)})$ :

$$\underline{p(x^{(i)}, z^{(i)})} = \underline{p(x^{(i)}|z^{(i)})}p(z^{(i)})$$

Examples:

- ▶ Unsupervised handwriting recognition is a mixture with 10 Bernoulli distributions
- ▶ Financial return estimation uses a mixture of 2 Gaussians for normal situation and crisis time distribution

# Mixture of Gaussians

Mixture of Gaussians Model:

$$\left. \begin{array}{l} z^{(i)} \sim \text{Multinomial}(\underline{\phi}) \\ \underline{x}^{(i)} | z^{(i)} = j \sim \mathcal{N}(\underline{\mu}_j, \underline{\Sigma}_j) \end{array} \right\}$$

How to learn  $\phi_j, \mu_j$  and  $\Sigma_j$  for all  $j$  ?

$z^{(i)}$  is known:  $(\underline{x}^{(i)}, \underline{z}^{(i)})_{i=1}^n$

$z^{(i)}$  is unknown:

# Mixture of Gaussians

Mixture of Gaussians Model:

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

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

How to learn  $\phi_j, \mu_j$  and  $\Sigma_j$  for all  $j$  ?

$z^{(i)}$  is known: (supervised) use maximum likelihood estimation  
(quadratic discriminant analysis). QDA GDA

$$\phi_j = \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\}, \quad \mu_j = \frac{\sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\}}$$

$$\Sigma_j = \frac{\sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\}}$$

$z^{(i)}$  is unknown:

# Mixture of Gaussians

Mixture of Gaussians Model:

$$z^{(i)} \sim \text{Multinomial}(\phi)$$
$$x^{(i)} | z^{(i)} = j \sim \mathcal{N}(\mu_j, \Sigma_j)$$

How to learn  $\phi_j, \mu_j$  and  $\Sigma_j$  for all  $j$  ?

$z^{(i)}$  **is known**: (supervised) use maximum likelihood estimation (quadratic discriminant analysis).

$$\phi_j = \frac{1}{m} \sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\}, \quad \mu_j = \frac{\sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\}}$$
$$\Sigma_j = \frac{\sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m \mathbf{1}\{z^{(i)} = j\}}$$

$z^{(i)}$  **is unknown**: (unsupervised) use **expectation maximization**

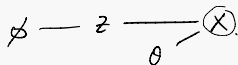
## Expectation Maximization

Overview

Algorithm Derivation

EM for mixture of Gaussians

# The EM Algorithm



The EM algorithm is an iterative method for maximum likelihood estimation when the model depends on latent (unobserved) variables.

Log-likelihood of data:

$$l(\theta) = \sum_{i=1}^m \log p(x^{(i)}; \theta) = \sum_{i=1}^m \log \sum_{z^{(i)}=1}^k p(x^{(i)}, z^{(i)}; \theta)$$

Main idea: iterate over two steps:

- ▶ Expectation (E) step : guess  $z^{(i)}$  for each  $x^{(i)}$ .
- ▶ Maximization (M) step : update  $\theta$  via maximum likelihood estimation based on guessed  $z^{(i)}$ 's



# Generalized EM Algorithm

## Listing 1: Generalized EM Algorithm

Initialize  $\theta$

Repeat until convergence {

(E-step) For each  $i$ , set *guess*  $z^i$  given  $x^i$ .

$Q_i(z^{(i)}) := p(z^{(i)}|x^{(i)}; \theta)$  ← Soft assignment:  
posterior distribution  $z|x$  under  $\theta$

(M-step) Set

$$\theta := \operatorname{argmax}_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad (*)$$

← Update parameter  $\theta$

}

# Generalized EM Algorithm

## Listing 2: Generalized EM Algorithm

Initialize  $\theta$

Repeat until convergence {

(E-step) For each  $i$ , set

$$Q_i(z^{(i)}) := \underbrace{p(z^{(i)}|x^{(i)}; \theta)} \leftarrow \text{Soft assignment:}$$

posterior distribution  $z|x$  under  $\theta$

(M-step) Set

$$\theta := \operatorname{argmax}_{\theta} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{\underbrace{Q_i(z^{(i)})}} \quad (*)$$

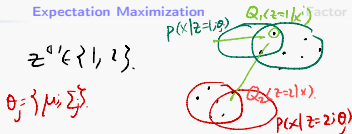
$\leftarrow$  Update parameter  $\theta$

}

We will show...

- ▶ Solving  $(*)$  is equivalent to  $\operatorname{argmax}_{\theta} \underline{l(\theta)}$   
 $\rightarrow$  Equation  $(*)$  is a (tight) lower bound on log-likelihood  $\underline{l(\theta)}$

## Generalized EM Algorithm



## Listing 3: Generalized EM Algorithm

Initialize  $\theta$ 

Repeat until convergence {

(E-step) For each  $i$ , set

$$Q_i(z^{(i)}) := p(z^{(i)}|x^{(i)}; \theta) \leftarrow \text{Soft assignment:}$$

posterior distribution  $z|x$  under  $\theta$ 

(M-step) Set

$$\theta := \operatorname{argmax}_{\theta} \sum_i \sum_{z^{(i)}} \overbrace{Q_i(z^{(i)})}^{J(Q, \theta)} \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad (*)$$

 $\leftarrow$  Update parameter  $\theta$ 

}

We will show...

- ▶ Solving (\*) is equivalent to  $\operatorname{argmax}_{\theta} l(\theta)$   
 $\rightarrow$  Equation (\*) is a (tight) lower bound on log-likelihood  $l(\theta)$
- ▶ This algorithm converges.

# Proof of Correctness: E-step

For each  $i$ , let  $Q_i(z)$  be a distribution of  $z$ :

$$\sum_{z \in \mathcal{Z}} Q_i(z) = 1, Q_i(z) \geq 0$$

Define

$$\underline{J(Q, \theta)} = \sum_{i=1}^m \sum_{z^{(i)} \in \mathcal{Z}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

## Proposition 1

1.  $\underline{J(Q, \theta)}$  is a lower bound on log-likelihood  $l(\theta)$
2. This lower bound is tight when  $Q_i(z^{(i)}) = \underline{p(z^{(i)} | x^{(i)}; \theta)}$

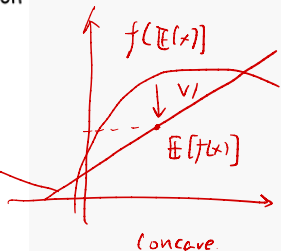
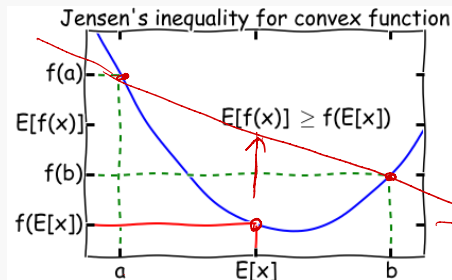
(Hint: use Jensen's inequality)

# Jensen's Inequality

## Theorem 1

Let  $f$  be a **convex** function, and let  $X$  be a random variable. Then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$

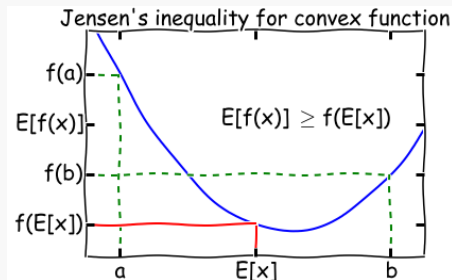


# Jensen's Inequality

## Theorem 1

Let  $f$  be a **convex** function, and let  $X$  be a random variable. Then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$



Remarks

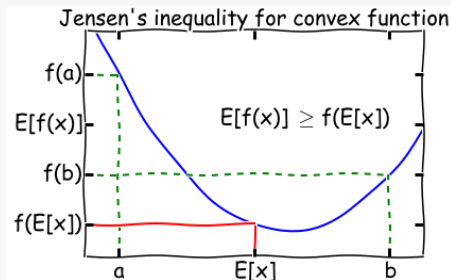
1. Let  $f$  be a **concave** function, then  $\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])$

# Jensen's Inequality

## Theorem 1

Let  $f$  be a **convex** function, and let  $X$  be a random variable. Then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$



Remarks

1. Let  $f$  be a **concave** function, then  $\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])$
2. When  $f(X)$  is a constant function,  $\mathbb{E}[f(X)] = f(\mathbb{E}[X])$

# Proof of Correctness

## Proposition 1

1.  $J(Q, \theta)$  is a lower bound on log-likelihood  $l(\theta)$

$$\begin{aligned}
 \text{proof. } \underline{l(\theta)} &= \sum_{i=1}^m \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta) \frac{Q_i(z^{(i)})}{Q_i(z^{(i)})} \\
 &= \sum_{i=1}^m \log \sum_{z^{(i)}} \underbrace{Q_i(z^{(i)})}_{=1} \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \\
 &= \sum_{i=1}^m \log \mathbb{E}_{z^{(i)} \sim Q_i} \left[ \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right].
 \end{aligned}$$

By Jensen's inequality,

$$\geq \sum_{i=1}^m \mathbb{E}_{z^{(i)} \sim Q_i} \left( \log \left[ \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right] \right) = \underbrace{\sum_{i=1}^m \sum_{z^{(i)} \in Z} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}}_{J(Q, \theta)}.$$

$J(Q, \theta)$  is a lower bound of  $\underline{l(\theta)}$



# Proof of Correctness

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

## Proposition 1

1.  $J(Q, \theta)$  is a lower bound on log-likelihood  $l(\theta)$  posterior distribution of  $z^{(i)}$  given  $x^{(i)}$
2. This lower bound is tight when  $Q_i(z^{(i)}) = p(z^{(i)} | x^{(i)}; \theta)$  (E-step)

proof. Suppose  $\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} = c$  for some constant  $c$ .

$$Q_i(z^{(i)}) = \frac{P(x^{(i)}, z^{(i)}; \theta)}{c}$$

$$\text{Since } \sum_{j=1}^K Q_i(z^{(i)}=j) = 1, \quad \sum_{j=1}^K \frac{P(x^{(i)}, z^{(i)}=j; \theta)}{c}$$

$$c = \sum_{j=1}^K P(x^{(i)}, z^{(i)}=j; \theta) = P(x^{(i)}; \theta)$$

$$\underline{Q_i(z^{(i)})} = \frac{P(x^{(i)}, z^{(i)}; \theta)}{P(x^{(i)}; \theta)} = \underline{P(z^{(i)} | x^{(i)}; \theta)}.$$

□

# Proof of Convergence

## Proposition 2

EM always monotonically improves the log likelihood, i.e. Let  $\theta^{(t)}$  be the parameter value in the  $t$ -th iteration

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

Proof. In proposition 1,  $Q^{(t)} = P(z^{(i)} | x^{(i)}; \theta^{(t)})$ ,  $J(Q^{(t)}, \theta^{(t)}) = \underline{l(\theta^{(t)})}$

In the M-step,

$$\theta^{(t+1)} := \operatorname{argmax}_{\theta} J(Q^{(t)}, \theta)$$

$$\text{Then } J(Q^{(t)}, \theta^{(t+1)}) \geq J(Q^{(t)}, \theta^{(t)})$$

By proposition 1,

$$l(\theta^{(t+1)}) \geq \underline{J(Q^{(t)}, \theta^{(t+1)})} \geq \underline{J(Q^{(t)}, \theta^{(t)})} = l(\theta^{(t)})$$

□

## EM for mixture of Gaussians

## Gaussian Mixture Model

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

$$\underbrace{\{x^{(i)} | z^{(i)}\}} \sim \mathcal{N}(\underline{\mu}_j, \underline{\Sigma}_j) \leftarrow \text{no } \phi$$

Learn parameters  $\mu, \Sigma, \phi$

**E-Step:**  $w_j^{(i)} = Q_i(z^{(i)} = j) = \underbrace{p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)}$

By Bayes's Rule  $= \frac{p(x^{(i)} | z^{(i)} = j; \mu, \Sigma) p(z^{(i)} = j; \phi)}{p(x^{(i)}; \mu, \Sigma)}$

$$= \sum_{l=1}^k \frac{p(x^{(i)} | z^{(i)} = l; \mu, \Sigma) p(z^{(i)} = l; \phi)}{p(x^{(i)}; \mu, \Sigma)}$$

## EM for mixture of Gaussians

## Gaussian Mixture Model

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

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

$$\text{Let } \nabla_{\Sigma_l} J(\phi, \mu, \Sigma) = 0.$$

$$\Rightarrow \sum_l^* = \frac{\sum_{i=1}^m w_l^{(i)} (x^{(i)} - \mu_l)(x^{(i)} - \mu_l)^T}{\sum_{i=1}^m w_j^{(i)}}$$

Learn parameters  $\mu, \Sigma, \phi$

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

**M-Step:** Maximize

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

with respect to  $\phi, \mu$  and  $\Sigma$

Assume  $w_j^{(i)} = Q_i(z^{(i)} = j)$  is given

$$J(\phi, \mu, \Sigma) = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \frac{1}{w_j^{(i)}} \left( \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right) \phi_j \right)$$

$$= \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \left( -\log((2\pi)^{\frac{D}{2}} |\Sigma_j|^{\frac{1}{2}}) - \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) - \log w_j^{(i)} + \log \phi_j \right)$$

$$\nabla_{\mu_l} J(\phi, \mu, \Sigma) = \sum_{i=1}^m w_l^{(i)} (\Sigma_l^{-1} x^{(i)} - \Sigma_l^{-1} \mu_l) = 0. \quad \mu_l^* = \frac{\sum_{i=1}^m w_l^{(i)} x^{(i)}}{\sum_{i=1}^m w_l^{(i)}}$$

$$\operatorname{argmax}_{\phi} \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j$$

$$\text{st. } \sum_{j=1}^k \phi_j = 1. \Rightarrow \underbrace{\sum_{j=1}^k \phi_j - 1 = 0.}$$

$$L(\phi, \beta) = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j + \beta \left( \sum_{j=1}^k \phi_j - 1 \right)$$

$$\frac{\partial}{\partial \phi_L} L(\phi, \beta) = \frac{\partial}{\partial \phi_L} \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j + \frac{\partial}{\partial \phi_L} \beta \left( \sum_{j=1}^k \phi_j - 1 \right)$$

$$= \sum_{i=1}^m \frac{\partial}{\partial \phi_L} w_i^{(i)} \log \phi_L + \beta \frac{\partial}{\partial \phi_L} \phi_L$$

$$\sum_{i=1}^m \frac{w_i^{(i)}}{\phi_L} + \beta = 0.$$

$$\phi_L = -\frac{1}{\beta} \sum_{i=1}^m w_i^{(i)}$$

$$\frac{\partial}{\partial \beta} L(\phi, \beta) = 0.$$

$$\Rightarrow \sum_{l=1}^k \phi_L = 1.$$

$$-\sum_{l=1}^k \sum_{i=1}^m \frac{w_i^{(i)}}{\beta} = 1.$$

$$-\underbrace{\sum_{i=1}^m \sum_{l=1}^k w_i^{(i)}}_1 = \beta.$$

$$-m = \beta.$$

$$\phi_L^* = -\frac{1}{(-m)} \sum_{i=1}^m w_i^{(i)} = \frac{1}{m} \sum_{i=1}^m w_i^{(i)}$$

# Expectation Maximization for Gaussian Mixtures

## Listing 4: EM for Gaussian Mixtures

Repeat until convergence {

(E-step) For each  $i, j$ , set

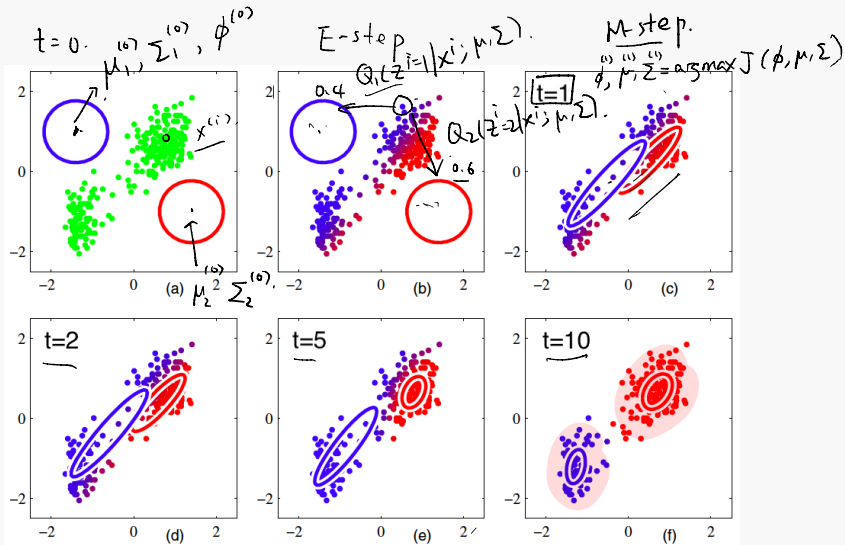
$$w_j^{(i)} := p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$$

(M-step) Update parameters: assume  $\phi_j = \mathbb{E}[w_j]$

$$\left\{ \begin{array}{l} \underline{\phi}_j := \frac{1}{m} \sum_{i=1}^m w_j^{(i)} \\ \mu_j := \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}} \\ \Sigma_j := \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}} \end{array} \right.$$

}

# Illustration of EM steps



# Comparison with k-means clustering

## Listing 4: EM Algorithm

Repeat until convergence {

(E-step) For each  $i, j$ ,

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

(M-step) Update parameters:

$$\phi_j := \frac{1}{m} \sum_{i=1}^m w_j^{(i)}$$

$$\mu_j := \frac{\sum_{i=1}^m w_j^{(i)} x_j}{\sum_{i=1}^m w_j^{(i)}}$$

$$\Sigma_j := \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^m w_j^{(i)}}$$

}

## Listing 5: (Lloyd's) k-means Alg.

Repeat until convergence {

(E-step) For every  $i$ ,

$$c^{(i)} := \underset{j}{\operatorname{argmin}} \|x^{(i)} - \mu_j\|^2$$

(M-step) Update centroids:

For each  $j$

$$\mu_j := \frac{\mathbf{1}\{c^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m \mathbf{1}\{c^{(i)} = j\}}$$

}

Similar to k-means, Gaussian mixtures are also subject to local minimums.



## Factor Analysis

Introduction

EM for Factor Analysis

Discussions

# Factor Analysis: Example

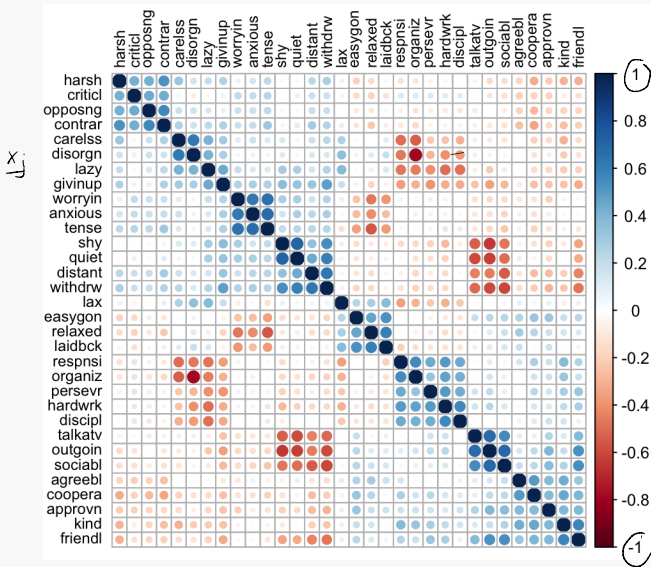
How much do you identify yourself with the following traits?

1-- the least    9 -- the most

|           | 1                     | 2                     | 3                     | 4                     | 5                     | 6                     | 7                     | 8                     | 9                     |
|-----------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| talkative | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> |
| distant   | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> |
| careless  | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> |
| hardwork  | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> |
| anxious   | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> |
| kind      | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> |

Self-ratings on 32 Personality Traits

# Factor Analysis: Example



Pairwise correlation plot of 32 variables from 240 participants

# Factor Analysis Terminology

- ▶ **observed random variables**  $x \in \mathbb{R}^n$ .

$$\underline{x} = \underline{\mu} + \underline{\Lambda} \underline{z} + \underline{\epsilon}$$

$\mathbb{R}^n$       $\mathbb{R}^{n \times k}$       $\mathbb{R}^k$       $\mathbb{R}^n$

- ▶ **factor**  $z \in \mathbb{R}^k$  is the hidden (latent) construct that “causes” the observed variables
- ▶ **factor loadings**  $\underline{\Lambda} \in \mathbb{R}^{n \times k}$  : the degree to which variable  $x_i$  is “caused” by the factors
- ▶  $\underline{\mu}, \underline{\epsilon} \in \mathbb{R}^n$  are the mean and error vectors

# Factor Analysis Terminology

- ▶ **observed random variables**  $x \in \mathbb{R}^n$

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

- ▶ **factor**  $z \in \mathbb{R}^k$  is the hidden (latent) construct that “causes” the observed variables
- ▶ **factor loadings**  $\Lambda \in \mathbb{R}^{n \times k}$  : the degree to which variable  $x_i$  is “caused” by the factors
- ▶  $\mu, \epsilon \in \mathbb{R}^n$  are the mean and error vectors

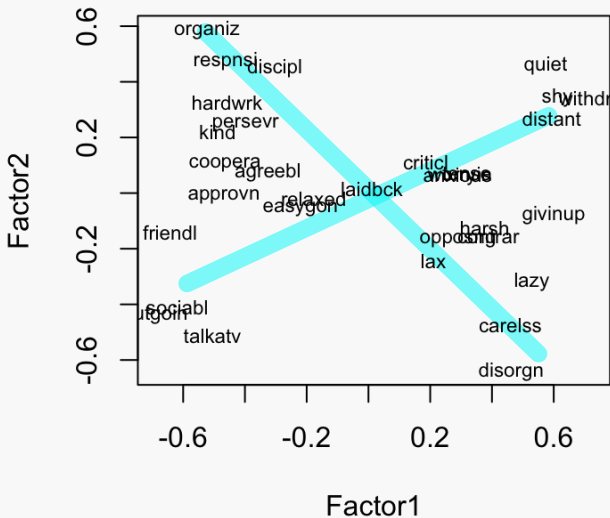
$$z \in \mathbb{R}^4$$

Matrix of factor loading  $\Lambda$  for personality test data

| variable    | factor 1 | factor 2 | factor 3 | factor 4 |
|-------------|----------|----------|----------|----------|
| distant     | 0.59     | 0.27     | 0        | 0        |
| talkative   | -0.50    | -0.51    | 0        | 0.27     |
| careless    | 0.46     | -0.47    | 0.11     | 0.14     |
| hardworking | -0.46    | 0.33     | -0.14    | 0.35     |
| kind        | -0.488   | 0.222    | 0        | 0        |
| ⋮           |          |          |          |          |

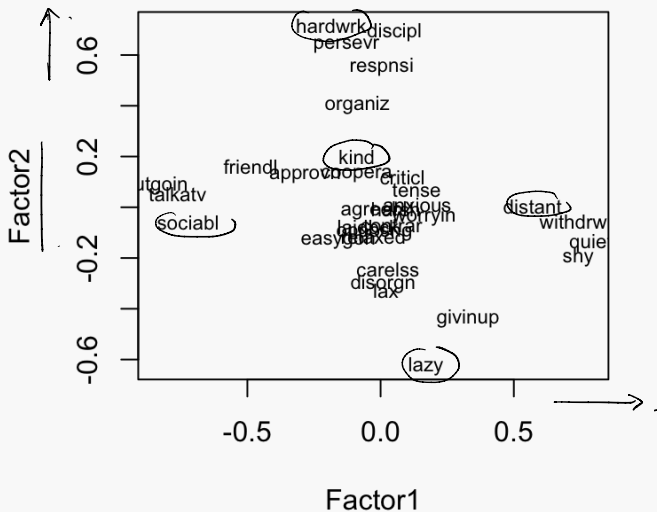
# Factor Analysis: Example

Visualize loading of the first two factors

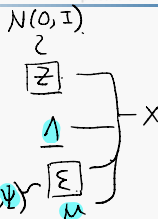


# Factor Analysis: Example

Visualize loading of the first two factors, rotated to align with axes



# Factor Analysis Model



Observed variables:  $\underline{x} \in \mathbb{R}^n$

Latent variables:  $\underline{z} \in \mathbb{R}^k$  ( $k < n$ )

The factor analysis model defines a joint distribution  $p(\underline{x}, \underline{z})$  as

$$\underline{z} \sim \mathcal{N}(0, I)$$

$$\underline{\epsilon} \sim \mathcal{N}(0, \underline{\Psi})$$

$$\underline{x} = \underline{\mu} + \underline{\Lambda} \underline{z} + \underline{\epsilon}$$

where  $\underline{\Psi} \in \mathbb{R}^{n \times n}$  is a diagonal matrix,  $\underline{\epsilon}, \underline{\mu} \in \mathbb{R}^n$ ,  $\underline{\Lambda} \in \mathbb{R}^{n \times k}$



# Factor Analysis Model

Observed variables:  $x \in \mathbb{R}^n$

Latent variables:  $z \in \mathbb{R}^k$  ( $k < n$ )

The factor analysis model defines a joint distribution  $p(x, z)$  as

$$\text{assume } z, \epsilon \text{ independent.} \left. \begin{array}{l} z \sim \mathcal{N}(0, I) \\ \epsilon \sim \mathcal{N}(0, \Psi) \end{array} \right\}$$

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

where  $\Psi \in \mathbb{R}^{n \times n}$  is a diagonal matrix,  $\epsilon, \mu \in \mathbb{R}^n$ ,  $\Lambda \in \mathbb{R}^{n \times k}$

Given observations  $x^{(1)}, \dots, x^{(m)}$ , how to fit the parameters  $\mu, \Lambda, \Psi$ ?

# The EM Algorithm

Rubin, D. and Thayer, D. (1982). *EM algorithms for ML factor analysis*. *Psychometrika*, 47(1):69-76.

## Listing 6: EM for Factor Analysis

Initialize  $\underline{\mu, \Lambda, \Psi}$

Repeat until convergence {

(E-step) For each  $i$ , set

$$Q_i(z^{(i)}) := \underline{p(z^{(i)} | x^{(i)}; \mu, \Lambda, \Psi)} \leftarrow z \text{ is a continuous variable}$$

(M-step) Set

$$\underline{\mu, \Lambda, \Psi} := \underset{\mu, \Lambda, \Psi}{\operatorname{argmax}} \underbrace{\sum_{i=1}^m \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(z^{(i)})} dz^{(i)}}_{J(Q, \theta)} \quad (*)$$

$\downarrow$   
 $\mu, \Lambda, \Psi$

# The EM Algorithm

Rubin, D. and Thayer, D. (1982). *EM algorithms for ML factor analysis*. Psychometrika, 47(1):69-76.

## Listing 7: EM for Factor Analysis

Initialize  $\mu, \Lambda, \Psi$

Repeat until convergence {

(E-step) For each  $i$ , set

$$Q_i(z^{(i)}) := p(z^{(i)}|x^{(i)}; \mu, \Lambda, \Psi) \leftarrow z \text{ is a continuous variable}$$

(M-step) Set

$$\mu, \Lambda, \Psi := \operatorname{argmax}_{\mu, \Lambda, \Psi} \sum_{i=1}^m \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(z^{(i)})} dz^{(i)} \quad (*)$$

First, we need to write  $\underline{p(z^{(i)}|x^{(i)})}$  and  $\underline{p(x^{(i)}, z^{(i)})}$  in terms of the model parameters.

## EM Derivations

$$z \sim \mathcal{N}(0, I) \quad x \sim \mathcal{N}(0, \Psi)$$

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

It can be shown that, random vector  $\begin{Bmatrix} z \\ x \end{Bmatrix} \sim \mathcal{N}(\mu_{zx}, \Sigma)$  where  $\mu_{xz} = \begin{bmatrix} 0 \\ \mu \end{bmatrix} \leftarrow \mathbb{E}[z]$

and  $\Sigma = \begin{bmatrix} 1 & \Lambda^T \\ \Lambda & \Lambda \Lambda^T + \Psi \end{bmatrix} = \begin{bmatrix} \Sigma_{zz} & \Sigma_{zx} \\ \Sigma_{xz} & \Sigma_{xx} \end{bmatrix}$

$\mathbb{E}[x] = \mu$ . Goal: find  $\Sigma_{zz}, \Sigma_{zx}, \dots$

$\mathbb{E}[z] = 0$ .  $\Sigma_{zz} = I$ .

$\Sigma_{zx} = \mathbb{E}[(z - \mathbb{E}[z])(x - \mathbb{E}[x])^T] = \mathbb{E}[z(x - \mu)^T]$ .

$\Sigma_{xx} = \Lambda \Lambda^T + \Psi$ .

$\mu = \mathbb{E}[z(\mu + \Lambda z + \xi)^T]$   
 $= \mathbb{E}[z(z^T \Lambda^T) + z \xi^T]$

$= \mathbb{E}[z z^T] \Lambda^T + \mathbb{E}[z \xi^T]$   
 $\text{cov}(z) = I \quad \mathbb{E}[z] \mathbb{E}[\xi^T] = 0$

$= I \Lambda^T = \Lambda^T$

## EM Derivations

It can be shown that, random vector  $\begin{bmatrix} z \\ x \end{bmatrix} \sim \mathcal{N}(\mu_{zx}, \Sigma)$  where  $\mu_{zx} = \begin{bmatrix} 0 \\ \mu \end{bmatrix}$

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

## E-Step

The posterior distribution  $\underline{z^{(i)}|x^{(i)}} \sim \mathcal{N}(\mu_{z^{(i)}|x^{(i)}}, \Sigma_{z^{(i)}|x^{(i)}})$

$$\mu_{z^{(i)}|x^{(i)}} = \Lambda^T(\Lambda\Lambda^T + \Psi)^{-1}(x^{(i)} - \mu)$$

$$\Sigma_{z^{(i)}|x^{(i)}} = I - \Lambda^T(\Lambda\Lambda^T + \Psi)^{-1}\Lambda$$

**Fact** For any  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim \mathcal{N}(\mu, \Sigma)$ ,  $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ ,  $\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$

$x_1|x_2 \sim \mathcal{N}(\mu_{1|2}, \Sigma_{1|2})$  where  $\mu_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2)$ .

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

$x_2|x_1$

# EM Derivations

It can be shown that, random vector  $\begin{bmatrix} z \\ x \end{bmatrix} \sim \mathcal{N}(\mu_{zx}, \Sigma)$  where  $\mu_{zx} = \begin{bmatrix} 0 \\ \mu \end{bmatrix}$

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

## E-Step

The posterior distribution  $z^{(i)}|x^{(i)} \sim \mathcal{N}(\mu_{z^{(i)}|x^{(i)}}, \Sigma_{z^{(i)}|x^{(i)}})$

$$\underline{\mu_{z^{(i)}|x^{(i)}}} = \Lambda^T (\Lambda\Lambda^T + \Psi)^{-1} (x^{(i)} - \mu)$$

$$\underline{\Sigma_{z^{(i)}|x^{(i)}}} = I - \Lambda^T (\Lambda\Lambda^T + \Psi)^{-1} \Lambda$$

$$\underline{Q_i(z^{(i)})} = p(z^{(i)}|x^{(i)}; \mu, \Lambda, \Psi)$$

$$= \frac{1}{\sqrt{(2\pi)^k |\underline{\Sigma_{z^{(i)}|x^{(i)}}}|}} \exp\left(-\frac{1}{2} (z^{(i)} - \underline{\mu_{z^{(i)}|x^{(i)}}})^T \underline{\Sigma_{z^{(i)}|x^{(i)}}}^{-1} (z^{(i)} - \underline{\mu_{z^{(i)}|x^{(i)}}})\right)$$

# EM Derivations

## M-Step

$$\operatorname{argmax}_{\mu, \Lambda, \Psi} \sum_{i=1}^m \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(z^{(i)})} dz^{(i)} \quad (*)$$

Note that

$$\begin{aligned} & \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(z^{(i)})} dz^{(i)} \\ &= \underbrace{\mathbb{E}_{z \sim Q_i}} [\log p(x^{(i)} | z^{(i)}; \mu, \Lambda, \Psi) + \log p(z^{(i)}) - \log Q_i(z^{(i)})] \end{aligned}$$

## EM Derivations

## M-Step

$$\operatorname{argmax}_{\mu, \Lambda, \Psi} \sum_{i=1}^m \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(z^{(i)})} dz^{(i)} \quad (*)$$

Note that

$$\begin{aligned} & \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(z^{(i)})} dz^{(i)} \\ &= \mathbb{E}_{z \sim Q_i} [\log p(x^{(i)} | z^{(i)}; \mu, \Lambda, \Psi) + \log p(z^{(i)}) - \log Q_i(z^{(i)})] \end{aligned}$$

(\*) is equivalent to

$$\operatorname{argmax}_{\mu, \Lambda, \Psi} \sum_{i=1}^m \mathbb{E}_{z^{(i)} \sim Q_i} [\log p(x^{(i)} | z^{(i)}; \mu, \Lambda, \Psi)]$$





# EM Derivations

## M-Step (con't)

$$\operatorname{argmax}_{\mu, \Lambda, \Psi} \sum_{i=1}^m \mathbb{E}_{z^{(i)} \sim Q_i} [\log p(x^{(i)} | z^{(i)}; \mu, \Lambda, \Psi)] \quad (**)$$

Since  $\underline{x} = \mu + \Lambda z + \epsilon$  and  $\epsilon \sim \mathcal{N}(0, \Psi)$

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

## EM Derivations

## M-Step (con't)

$$\operatorname{argmax}_{\mu, \Lambda, \Psi} \sum_{i=1}^m \mathbb{E}_{z^{(i)} \sim Q_i} [\log p(x^{(i)} | z^{(i)}; \mu, \Lambda, \Psi)] \quad (**)$$

Since  $x = \mu + \Lambda z + \epsilon$  and  $\epsilon \sim \mathcal{N}(0, \Psi)$

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

$$\begin{aligned} & p(x^{(i)} | z^{(i)}; \mu, \Lambda, \Psi) \\ &= \frac{1}{(2\pi)^{n/2} |\Psi|^{1/2}} \exp\left(-\frac{1}{2} (x^{(i)} - \mu - \Lambda z^{(i)})^T \Psi^{-1} (x^{(i)} - \mu - \Lambda z^{(i)})\right) \end{aligned}$$

We can maximize (\*\*) with respect to  $\underline{\mu}$ ,  $\underline{\Lambda}$  and  $\underline{\Psi}$

# Factor Analysis Discussions

## Comparison with Mixture of Gaussians

- ▶ Mixture of Gaussians assumes sufficient data and relative few response variables. i.e. when  $n \approx m$  or  $n > m$ ,  $\Sigma$  is singular

$x_1, \dots, x_n$

sample size  $m$   $>$   $n$

# Factor Analysis Discussions

## Comparison with Mixture of Gaussians

- ▶ Mixture of Gaussians assumes sufficient data and relative few response variables. i.e. when  $n \approx m$  or  $n > m$ ,  $\Sigma$  is singular
- ▶ Factor Analysis works when  $n > m$  by allowing model noise  $\xi$

# Factor Analysis Discussions

## Relationship to PCA

- ▶ Both PCA and factor analysis can find low dimensional latent subspace in data

# Factor Analysis Discussions

## Relationship to PCA

- ▶ Both PCA and factor analysis can find low dimensional latent subspace in data
- ▶ PCA is good for data reduction (reduce correlation among observed variables)

# Factor Analysis Discussions

## Relationship to PCA

- ▶ Both PCA and factor analysis can find low dimensional latent subspace in data
- ▶ PCA is good for data reduction (reduce correlation among observed variables)
- ▶ Factor analysis is good for data exploration (find independent, common factors in observed variables)

# Factor Analysis Discussions

## Relationship to PCA

- ▶ Both PCA and factor analysis can find low dimensional latent subspace in data
- ▶ PCA is good for data reduction (reduce correlation among observed variables)
- ▶ Factor analysis is good for data exploration (find independent, common factors in observed variables)
- ▶ Factor analysis allows the noise to have an arbitrary diagonal covariance matrix, while PCA assumes the noise is spherical.

## Additional readings

- ▶ Zoubin Ghahramani and Geoffrey E. Hinton, The EM Algorithm for Mixtures of Factor Analyzers, 1997



# Next Lecture

## Semi-Supervised Learning

- ▶ Semi-supervised SVM
- ▶ Graph-based semi-supervised learning spectral clustering
- ▶ Deep semi-supervised learning

*No WA5. Please focus on your final project!*

*The class on Dec 31st will be PA5 discussion, project Q&A and group study*