

# Intro to ML

October 25<sup>th</sup>, 2021

CHAPTER 6:

# Dimensionality Reduction

# Why Reduce Dimensionality?

- Reduces time complexity: Less computation
- Reduces space complexity: Fewer parameters
- Saves the cost of observing the feature
- Simpler models are more robust on small datasets
- More interpretable; simpler explanation
- Data visualization (structure, groups, outliers, etc) if plotted in 2 or 3 dimensions

# Factor Analysis

- Find a small number of unobservable factors  $\mathbf{z}$ , which when combined generate  $\mathbf{x}$ :

$$x_i - \mu_i = v_{i1}z_1 + v_{i2}z_2 + \dots + v_{ik}z_k + \varepsilon_i$$

where  $z_j, j=1, \dots, k$  are the latent factors with

$$E[z_j]=0, \text{Var}(z_j)=1, \text{Cov}(z_i, z_j)=0, i \neq j,$$

$\varepsilon_i$  are the **noise sources**

$$E[\varepsilon_i]=0, \text{Var}(\varepsilon_i)=\psi_i, \text{Cov}(\varepsilon_i, \varepsilon_j)=0, i \neq j, \text{Cov}(\varepsilon_i, z_j)=0$$

and  $v_{ij}$  are the **factor loadings**

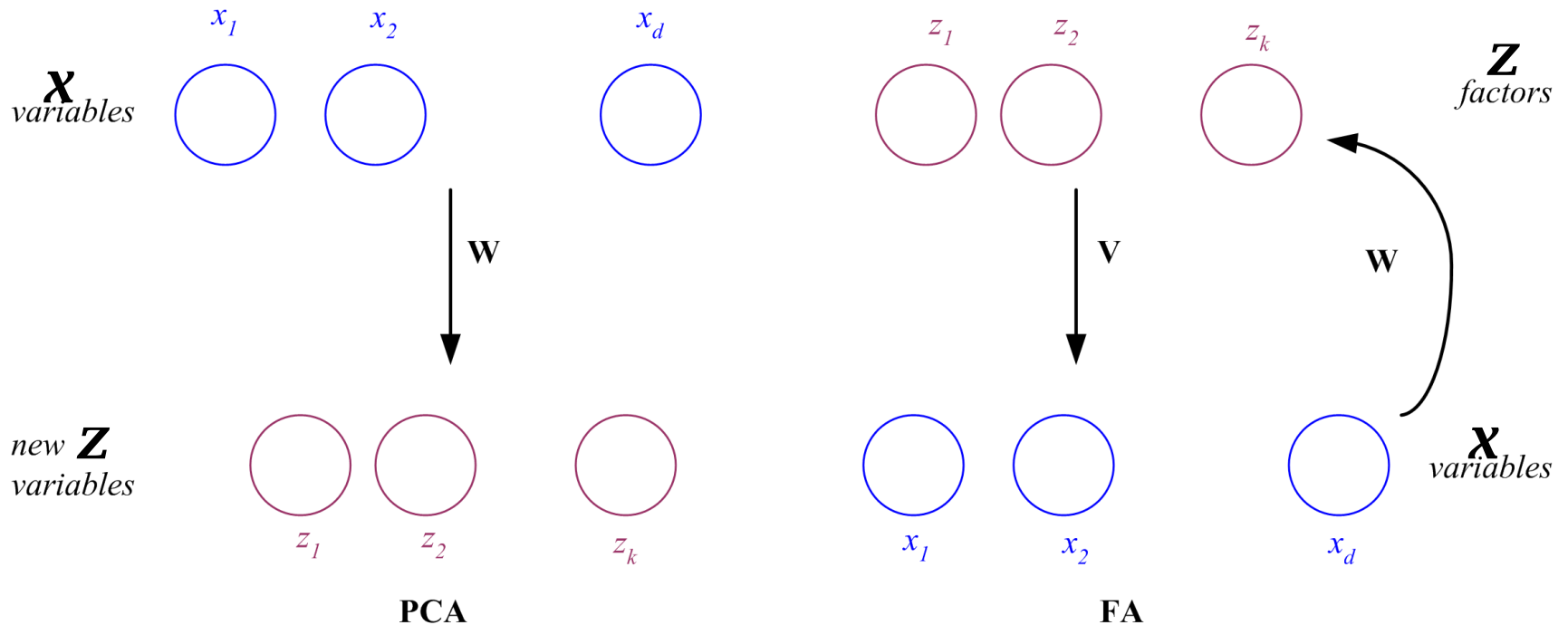
Factors are:  
Unit normal and  
uncorrelated

# PCA vs FA

- PCA From  $\mathbf{x}$  to  $\mathbf{z}$
- FA From  $\mathbf{z}$  to  $\mathbf{x}$

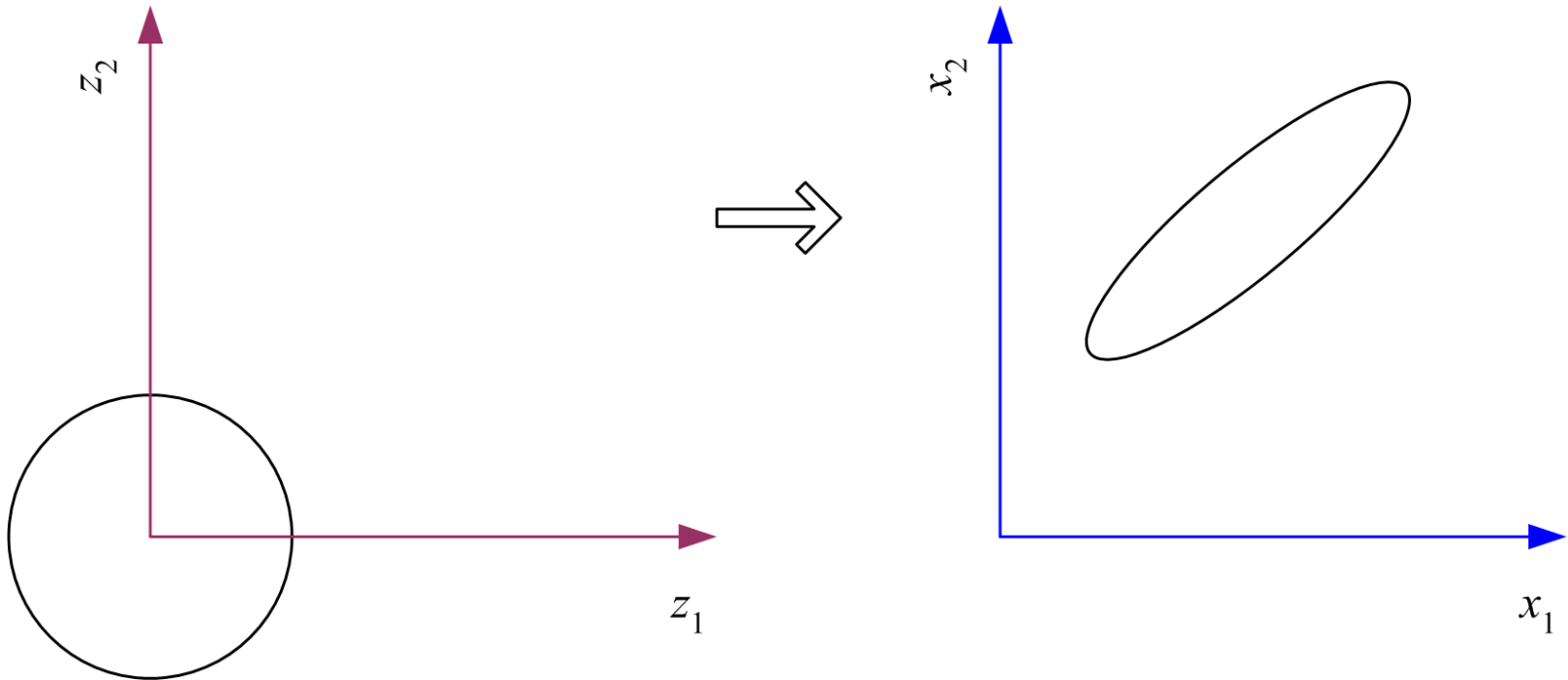
$$\mathbf{z} = \mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu})$$

$$\mathbf{x} - \boldsymbol{\mu} = \mathbf{V}\mathbf{z} + \boldsymbol{\varepsilon}$$



# Factor Analysis

- In FA, factors  $z_j$  are stretched, rotated and translated to generate  $\mathbf{x}$



- $X_i - \mu_i = v_{i1}z_1 + v_{i2}z_2 + \dots + v_{ik}z_k + \varepsilon_i$
- $Z$

$$\text{Var}(X_i) = v_{i1}^2 + v_{i2}^2 + \dots + \psi_i$$

- Variance of  $X_i$  attributed to common factor 1-k and noise

$$\begin{aligned} \Sigma &= \text{cov}(X) = \text{cov}(Vz + \varepsilon) \\ &= \text{cov}(Vz) + \text{cov}(\varepsilon) \\ &= V \boxed{\text{cov}(z)} V^T + \psi \quad \text{Identity matrix} \\ &= VV^T + \boxed{\psi} \quad \text{Diagonal matrix} \end{aligned}$$

Assume there are two 'hidden factors', and two observable features

$$\text{Cov}(x_1, x_2) = v_{11}v_{21} + v_{12}v_{22}$$

- $\text{Cov}(x_1, x_2) = v_{11}v_{21} + v_{12}v_{22}$
- If  $x_1$   $x_2$  have high covariance due to the first 'factor'
  - The  $v_{11}v_{21}$  will be high or  $v_{12}v_{22}$  will be high
  - Either way the above terms will be high
- If they depend on different factors, one term high one term low, this summation will be small

$$\text{Cov}(x_1, z_2) = \text{Cov}(v_{12}z_2, z_2) = v_{12}$$

- $\text{Cov}(x, z) = V$
- Loading represent correlation between variables and the factors
- S estimate  $\Sigma$
- Factor analysis basically solves the following equation

$$S = VV^T + \Psi$$



# Linear Discriminant Analysis

- Find a low-dimensional space s.t. when  $\mathbf{x}$  is projected classes are as well separated as possible (supervised method)
  - $\mathbf{z} = \mathbf{W}^T \mathbf{x}$ : projection of  $\mathbf{x}$  onto  $\mathbf{w}$ , a dimensional from  $d \rightarrow 1$
- $m_1$ : original sample mean,  $m_1$ : after projection
- $m_2$ : original sample mean,  $m_2$ : after projection
- For a sample,  $\mathbf{X} = \{\mathbf{x}^t, r^t\}$ ,  $r^t = 1$  is  $\mathbf{x}^t = \text{class 1}$

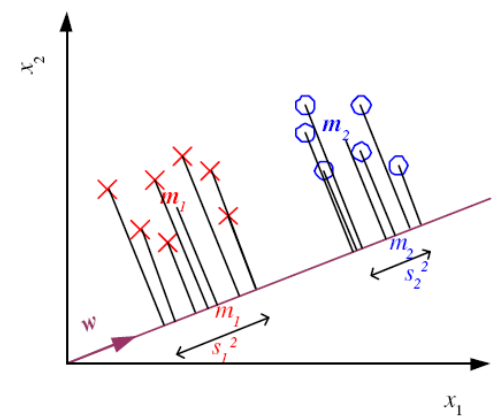
$$m_1 = \frac{\sum_t \mathbf{w}^T \mathbf{x}^t r^t}{\sum_t r^t} \quad s_1^2 = \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t$$



- find  $\mathbf{w}$  that maximizes

After projection for  
 $m_1, s_1$  (scatter)

# LDA



- The goal is to let mean as far apart as possible and let the scatter for each class to be as clustered as possible
  - $|m_1 - m_2|^2$  large,  $s_1^2 + s_2^2$  as small
- Fisher's discriminant, find  $w$ , such that

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

Is maximized

- Between-class scatter:

$$\begin{aligned}
 (m_1 - m_2)^2 &= (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\
 &= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} \\
 &= \mathbf{w}^T \mathbf{S}_B \mathbf{w} \text{ where } \mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T
 \end{aligned}$$

- Within-class scatter:

$$\begin{aligned}
 s_1^2 &= \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t \\
 &= \sum_t \mathbf{w}^T (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} r^t = \mathbf{w}^T \mathbf{S}_1 \mathbf{w}
 \end{aligned}$$


where  $\mathbf{S}_1 = \sum_t (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T r^t$

$$s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_W \mathbf{w} \text{ where } \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

# Fisher's Linear Discriminant

- Find  $\mathbf{w}$  that max

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} = \frac{|\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|^2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

- LDA soln:  $\mathbf{w} = c \cdot \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$   Only direction matters, set  $c=1$

- For more than 2 classes ( $k > 2$ )

$$\mathbf{z} = \mathbf{W}^T \mathbf{x}$$

- $\mathbf{z}$  is  $k$ -dimensional,  $\mathbf{W}$  is  $d \times k$

# K>2 Classes

- Within-class scatter:

Total scatter  $\rightarrow$

$$\mathbf{S}_W = \sum_{i=1}^K \mathbf{S}_i \quad \mathbf{S}_i = \sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T$$

For each class  $i$ ,  $r_i^t$  is 1 if  $\mathbf{x}^t$  belongs to class  $i$

- Between-class scatter:

$$\mathbf{S}_B = \sum_{i=1}^K N_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T \quad \mathbf{m} = \frac{1}{K} \sum_{i=1}^K \mathbf{m}_i$$

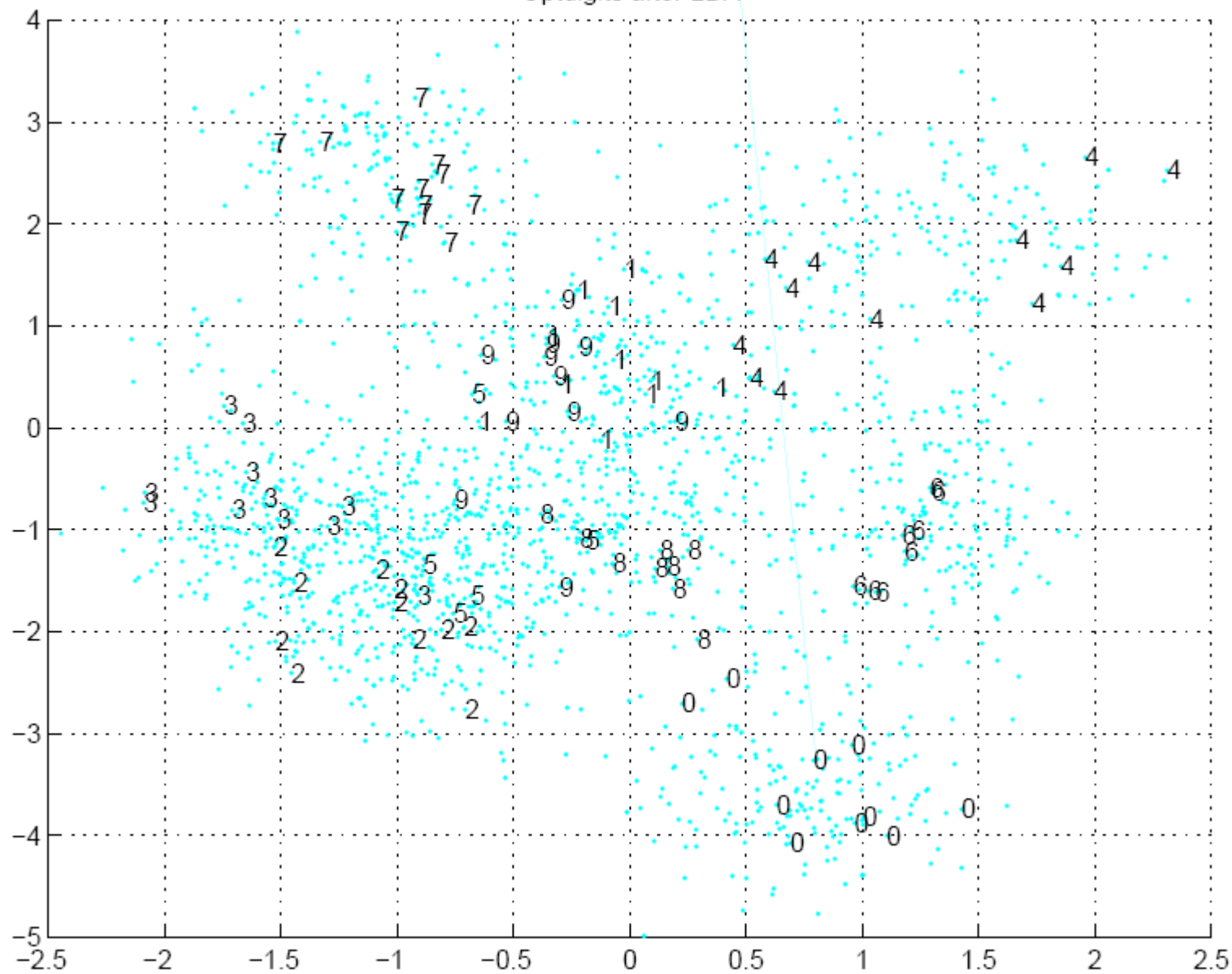
- Find  $\mathbf{W}$  matrix that max  $J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}$

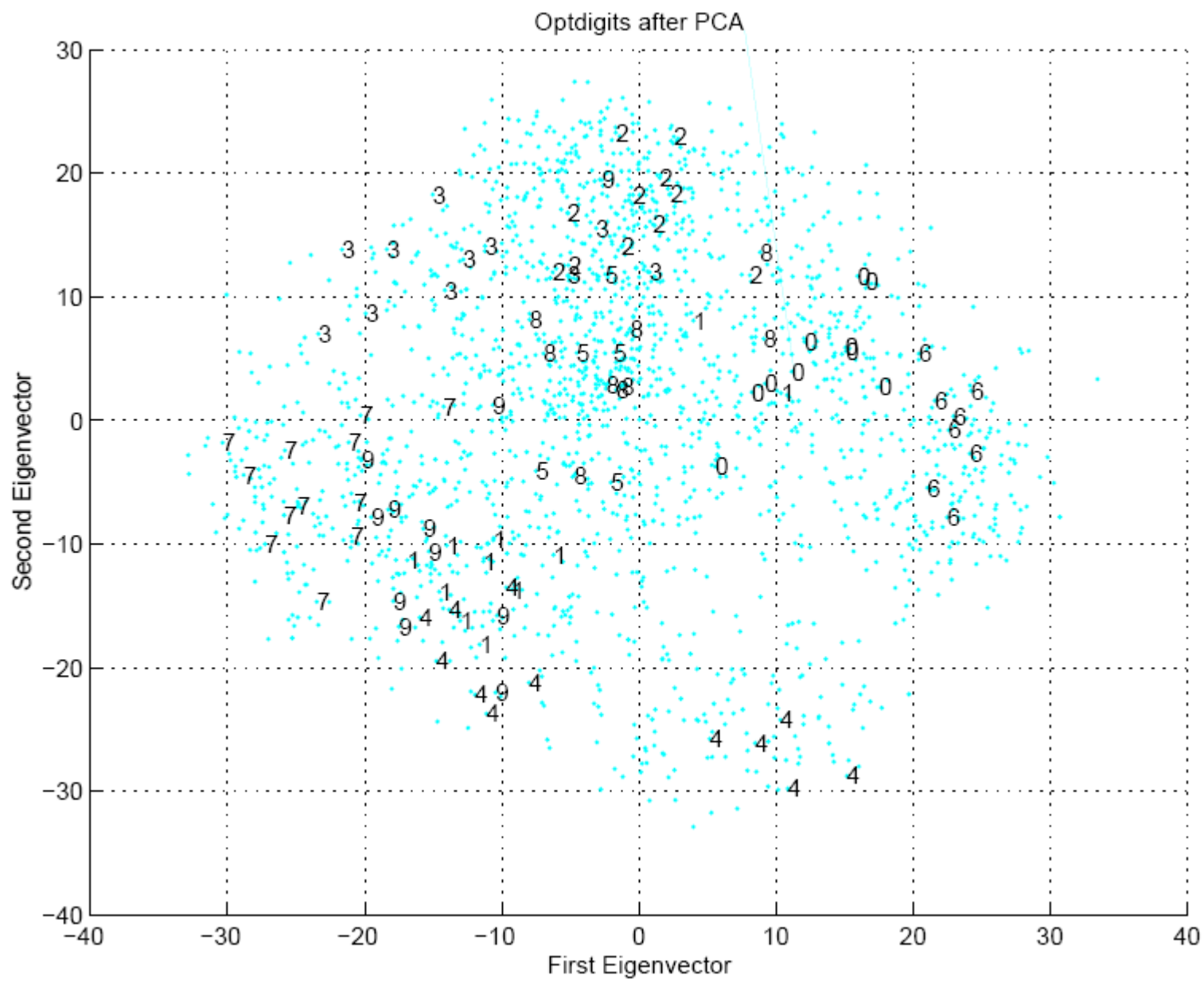
Want class means to be as far apart, and samples from the same class are close to their mean

# LDA ( $k > 2$ )

- The largest eigenvectors of  $\mathbf{S}_W^{-1}\mathbf{S}_B$  is the solution
- $\mathbf{S}_B$  is the sum of  $K$  matrices
  - Hence, maximum rank of  $K-1$  (only  $K-1$  independent)
  - $k$  is  $K-1$
- LDA only works better since it makes use of the label information (do it in the training set)
- This solution is workable if  $\mathbf{S}_W$  is invertible, if not, use PCA to get rid of singularity then LDA

Optdigits after LDA



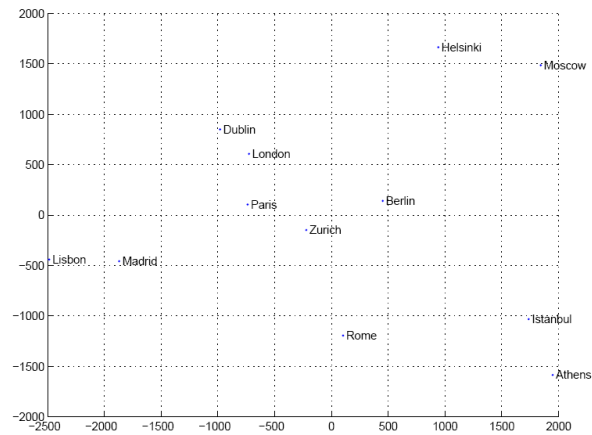


Visual display



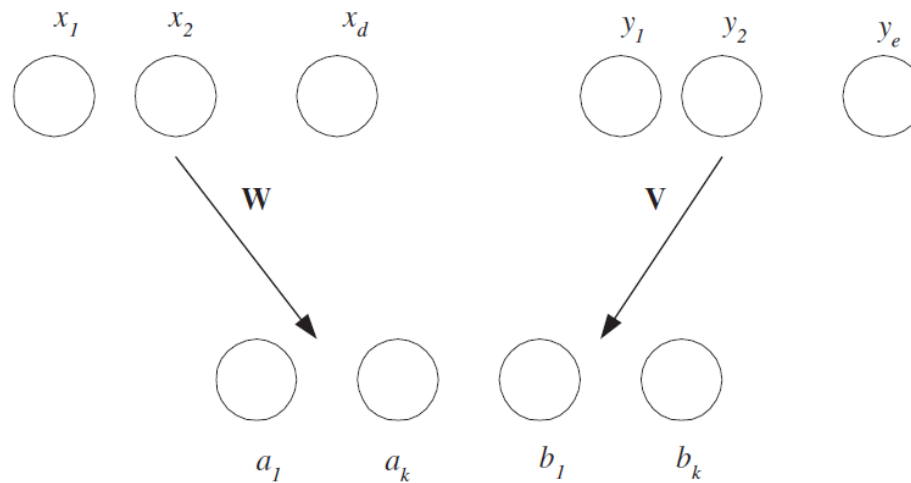
# Other important in the chapter

- **Multidimensional scaling**
  - Given pairwise distances between  $N$  points,  
 $d_{ij}, i, j = 1, \dots, N$  place on a low-dim map s.t. distances are preserved
  - Good for visualization with proper distancing
- **Canonical correlation**
  - $X = \{x^t, y^t\}_t$  ; two sets of variables  $x$  and  $y$
  - We want to find two projections  $w$  and  $v$  st when  $x$  is projected along  $w$  and  $y$  is projected along  $v$ , the correlation is maximized
  - Work in cases where  $x$ , and  $y$  are different dimensions



Map of Europe by MDS  
(from road travel distances)

## CCA



CHAPTER 7:

# Clustering

# Semiparametric Density Estimation

- **Parametric:** Assume a single model for  $p(\mathbf{x} | C_i)$  (Chapters 4 and 5)
- **Semiparametric:**  $p(\mathbf{x} | C_i)$  is a mixture of densities  
Multiple possible explanations/prototypes:  
Different handwriting styles, accents in speech  
Not a single Gaussian distribution for the group of interest
- **Nonparametric:** No model; data speaks for itself (Chapter 8)

# Mixture Densities

$$p(\mathbf{x}) = \sum_{i=1}^k p(\mathbf{x} | G_i) P(G_i)$$

where  $G_i$  the components/groups/clusters,

$P(G_i)$  mixture proportions (priors),

$p(\mathbf{x} | G_i)$  component densities

Gaussian mixture where  $p(\mathbf{x} | G_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$  parameters

$$\Phi = \{P(G_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^k$$

unlabeled sample  $X = \{\mathbf{x}^t\}_t$  (unsupervised learning)

# Classes vs. Clusters

- Supervised:  $X = \{\mathbf{x}^t, r^t\}_t$
- Classes  $C_i, i=1, \dots, K$

$$p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x} | C_i) P(C_i)$$

where  $p(\mathbf{x} | C_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

- $\Phi = \{P(C_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^K$

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N} \quad \mathbf{m}_i = \frac{\sum_t r_i^t \mathbf{x}^t}{\sum_t r_i^t}$$

$$\mathbf{S}_i = \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_t r_i^t}$$

- Unsupervised :  $X = \{\mathbf{x}^t\}_t$
- Clusters  $G_i, i=1, \dots, k$

$$p(\mathbf{x}) = \sum_{i=1}^k p(\mathbf{x} | G_i) P(G_i)$$

where  $p(\mathbf{x} | G_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

- $\Phi = \{P(G_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^k$

Labels  $r_i^t$  ?



Through  
clustering

# Imagine a case


- A image, 24 bits/pixel,  $\sim$  16 million colors
- Say we have a screen with 8 bits/pixel
  - 256 colors only
- Find the best colors among 16 million colors so that the image look as close to the original image as possible -> color quantization
- Vector quantization ->continuous value to discrete space

# k-Means Clustering

- Find  $k$  **reference vectors** (prototypes/codebook vectors/codewords) which best represent data
- Reference vectors,  $\mathbf{m}_j, j = 1, \dots, k$
- Use nearest (most similar) reference:

$$\|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\|$$

find a reference vector that is close to the  $\mathbf{x}$



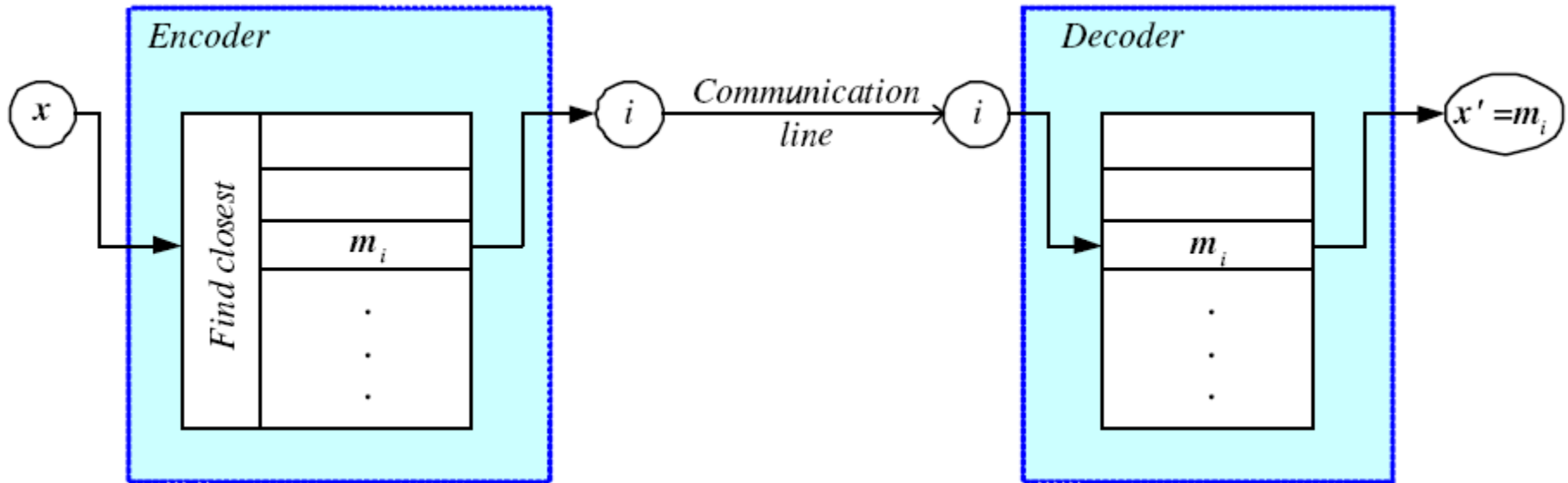
- Reconstruction error

$$E(\{\mathbf{m}_i\}_{i=1}^k | \mathcal{X}) = \sum_t \sum_i b_i^t \|\mathbf{x}^t - \mathbf{m}_i\|$$

$$b_i^t = \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$



# Encoding/Decoding



Vector quantization saves bit, at the receiving end, there is an error

# $k$ -means Clustering

No analytical solution, since it's about finding that reference vector  
Results in an iterative approach

Initialize  $\mathbf{m}_i, i = 1, \dots, k$ , for example, to  $k$  random  $\mathbf{x}^t$

Repeat

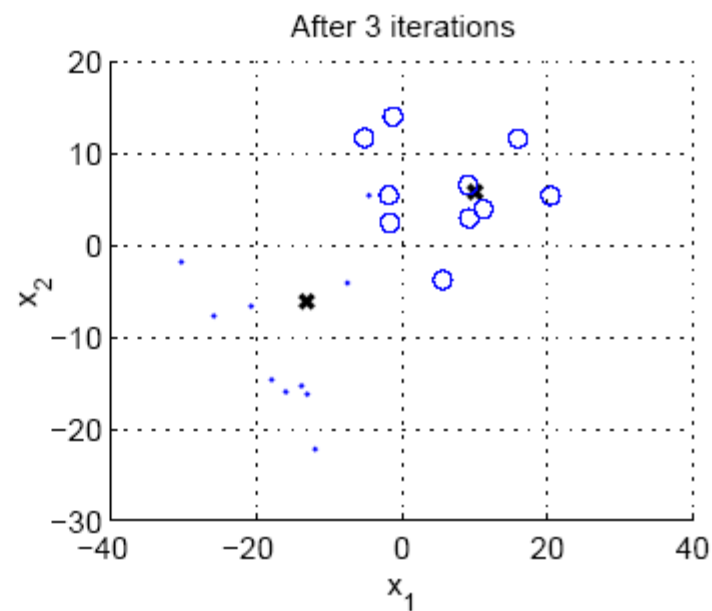
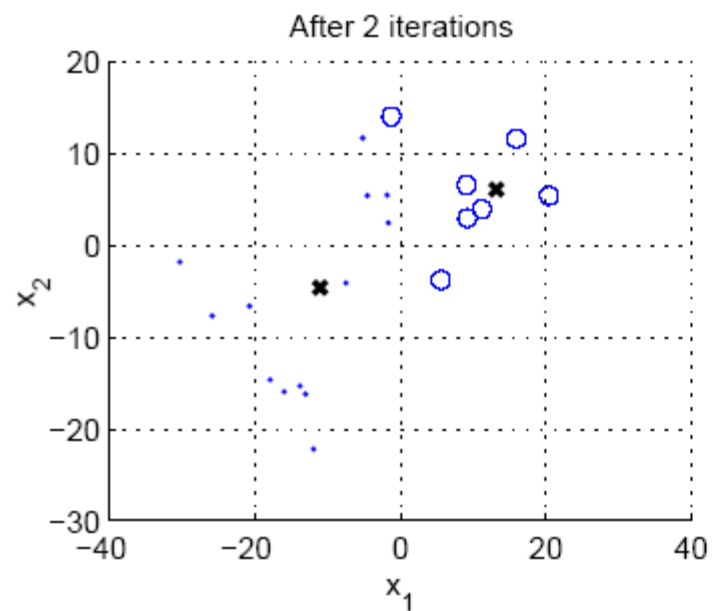
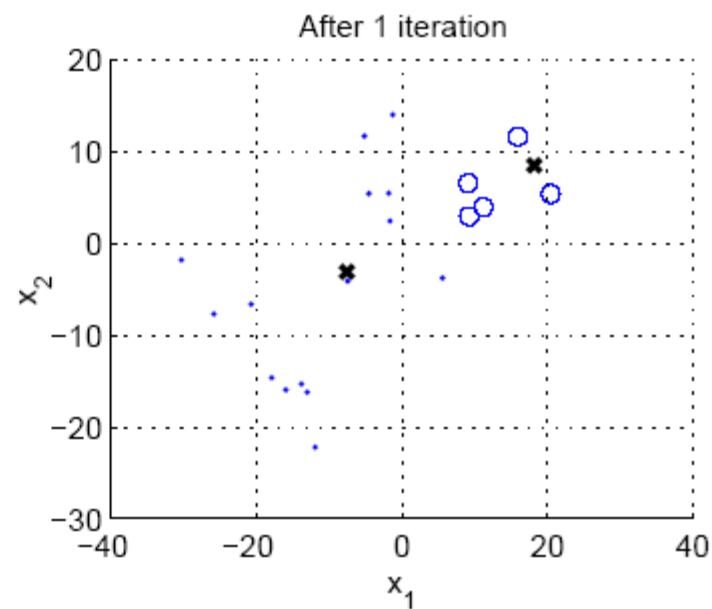
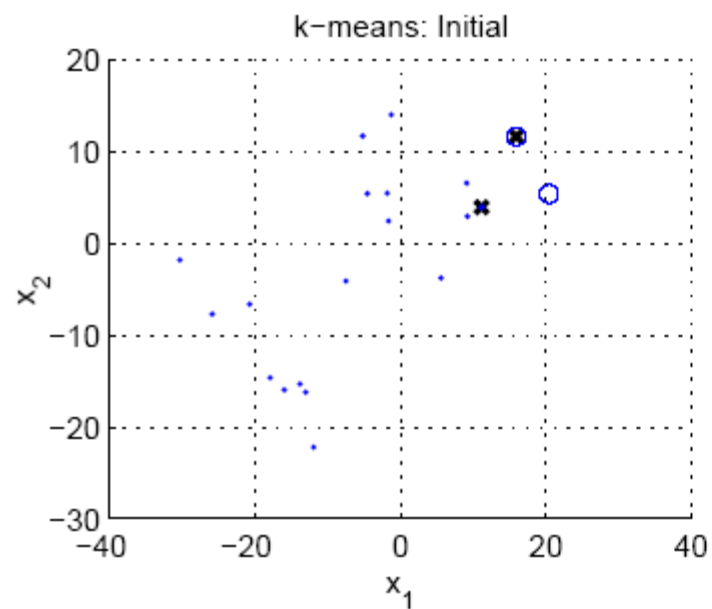
For all  $\mathbf{x}^t \in \mathcal{X}$

$$b_i^t \leftarrow \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$

For all  $\mathbf{m}_i, i = 1, \dots, k$

$$\mathbf{m}_i \leftarrow \sum_t b_i^t \mathbf{x}^t / \sum_t b_i^t$$

Until  $\mathbf{m}_i$  converge




K-means aim at finding codebook that minimizes reconstruction error

# Expectation-Maximization (EM)

- Log likelihood with a mixture model

$$\begin{aligned}\mathcal{L}(\Phi | \mathcal{X}) &= \log \prod_t p(\mathbf{x}^t | \Phi) \\ &= \sum_t \log \sum_{i=1}^k p(\mathbf{x}^t | G_i) p(G_i)\end{aligned}$$

Unknown, no  
analytical solution



EM Algorithm, core concept

- Assume there exist hidden variables  $z$ , which when known, make optimization much simpler
- Complete likelihood,  $L_c(\Phi | X, Z)$ , in terms of  $\mathbf{x}$  and  $\mathbf{z}$
- Incomplete likelihood,  $L(\Phi | X)$ , in terms of  $\mathbf{x}$

# E- and M-steps

Model parameter



Iterate the two steps

1. E-step: Estimate  $z$  given  $X$  and current  $\Phi$
2. M-step: Find new  $\Phi'$  given  $z$ ,  $X$ , and old  $\Phi$ .

$$\text{E-step: } \mathcal{Q}(\Phi | \Phi') = E[\mathcal{L}_c(\Phi | \mathcal{X}, Z) | \mathcal{X}, \Phi']$$

$$\text{M-step: } \Phi^{l+1} = \arg\max_{\Phi} \mathcal{Q}(\Phi | \Phi')$$

An increase in  $Q$  increases incomplete likelihood

$$\mathcal{L}(\Phi^{l+1} | \mathcal{X}) \geq \mathcal{L}(\Phi' | \mathcal{X})$$



There is proof  
beyond this class

# EM in Gaussian Mixtures



Mixture of Gaussian  
distribution

$$p(\mathbf{x}) = \sum_{i=1}^k p(\mathbf{x} | G_i) p(G_i)$$

- $z_j^t = 1$  if  $\mathbf{x}^t$  belongs to  $G_j$ , 0 otherwise; assume  $p(\mathbf{x} | G_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$
- E-step
  - Expectation, find the expected value of the current model parameters under the current parameter set
- M-step
  - Maximize these parameters, and iterate



$$(7.7) \quad P(\mathbf{z}^t) = \prod_{i=1}^k \pi_i^{z_i^t}$$

The likelihood of an observation  $\mathbf{x}^t$  is equal to its probability specified by the component that generated it:

$$(7.8) \quad p(\mathbf{x}^t | \mathbf{z}^t) = \prod_{i=1}^k p_i(\mathbf{x}^t)^{z_i^t}$$

$p_i(\mathbf{x}^t)$  is shorthand for  $p(\mathbf{x}^t | \mathcal{G}_i)$ . The joint density is

$$p(\mathbf{x}^t, \mathbf{z}^t) = P(\mathbf{z}^t) p(\mathbf{x}^t | \mathbf{z}^t)$$

and the complete data likelihood of the iid sample  $\mathcal{X}$  is

$$\begin{aligned} \mathcal{L}_c(\Phi | \mathcal{X}, \mathcal{Z}) &= \log \prod_t p(\mathbf{x}^t, \mathbf{z}^t | \Phi) \\ &= \sum_t \log p(\mathbf{x}^t, \mathbf{z}^t | \Phi) \\ &= \sum_t \log P(\mathbf{z}^t | \Phi) + \log p(\mathbf{x}^t | \mathbf{z}^t, \Phi) \\ &= \sum_t \sum_i z_i^t [\log \pi_i + \log p_i(\mathbf{x}^t | \Phi)] \end{aligned}$$



**E-step:** We define

$$\begin{aligned} \mathcal{Q}(\Phi|\Phi^l) &\equiv E \left[ \log P(X, Z) | \mathcal{X}, \Phi^l \right] \\ &= E \left[ \mathcal{L}_c(\Phi | \mathcal{X}, Z) | \mathcal{X}, \Phi^l \right] \\ &= \sum_t \sum_i E[z_i^t | \mathcal{X}, \Phi^l] [\log \pi_i + \log p_i(\mathbf{x}^t | \Phi)] \end{aligned}$$

where

$$\begin{aligned} E[z_i^t | \mathcal{X}, \Phi^l] &= E[z_i^t | \mathbf{x}^t, \Phi^l] \quad \mathbf{x}^t \text{ are iid} \\ &= P(z_i^t = 1 | \mathbf{x}^t, \Phi^l) \quad z_i^t \text{ is a 0/1 random variable} \\ &= \frac{p(\mathbf{x}^t | z_i^t = 1, \Phi^l) P(z_i^t = 1 | \Phi^l)}{p(\mathbf{x}^t | \Phi^l)} \quad \text{Bayes' rule} \\ &= \frac{p_i(\mathbf{x}^t | \Phi^l) \pi_i^l}{\sum_j p_j(\mathbf{x}^t | \Phi^l) \pi_j^l} \\ &= \frac{p(\mathbf{x}^t | \mathcal{G}_i, \Phi^l) P(\mathcal{G}_i)}{\sum_j p(\mathbf{x}^t | \mathcal{G}_j, \Phi^l) P(\mathcal{G}_j)} \\ &= P(\mathcal{G}_i | \mathbf{x}^t, \Phi^l) \equiv h_i^t \end{aligned}$$

(7.9)



- E-step:
  - Expected value is the posterior probability of the sample coming from that mixture (cluster)
  - It's between 0-1
  - We can also think about it as soft assignment of cluster for each sample
- In E step, given data  $X$ , compute the complete data likelihood given the current model parameters

**M-step:** We maximize  $\mathcal{Q}$  to get the next set of parameter values  $\Phi^{l+1}$ :

$$\Phi^{l+1} = \arg \max_{\Phi} \mathcal{Q}(\Phi | \Phi^l)$$

which is

$$\begin{aligned} \mathcal{Q}(\Phi | \Phi^l) &= \sum_t \sum_i h_i^t [\log \pi_i + \log p_i(\mathbf{x}^t | \Phi)] \\ &= \sum_t \sum_i h_i^t \log \pi_i + \sum_t \sum_i h_i^t \log p_i(\mathbf{x}^t | \Phi) \end{aligned}$$

The second term is independent of  $\pi_i$  and using the constraint that  $\sum_i \pi_i = 1$  as the Lagrangian, we solve for

$$\nabla_{\pi_i} \sum_t \sum_i h_i^t \log \pi_i - \lambda \left( \sum_i \pi_i - 1 \right) = 0$$

and get

$$1) \quad \pi_i^{l+1} = \frac{\sum_t h_i^t}{N}$$

which is analogous to the calculation of priors in equation 7.2.

Similarly, the first term of equation 7.10 is independent of the components and can be dropped while estimating the parameters of the components. We solve for

$$2) \quad \nabla_{\Phi} \sum_t \sum_i h_i^t \log p_i(\mathbf{x}^t | \Phi) = 0$$

# Complete steps for GMM

- If assume Gaussian component (each mixture is a Gaussian distribution)

$$P(G_i) = \frac{\sum_t h_i^t}{N} \quad \mathbf{m}_i^{l+1} = \frac{\sum_t h_i^t \mathbf{x}^t}{\sum_t h_i^t}$$

$$\mathbf{S}_i^{l+1} = \frac{\sum_t h_i^t (\mathbf{x}^t - \mathbf{m}_i^{l+1})(\mathbf{x}^t - \mathbf{m}_i^{l+1})^T}{\sum_t h_i^t}$$

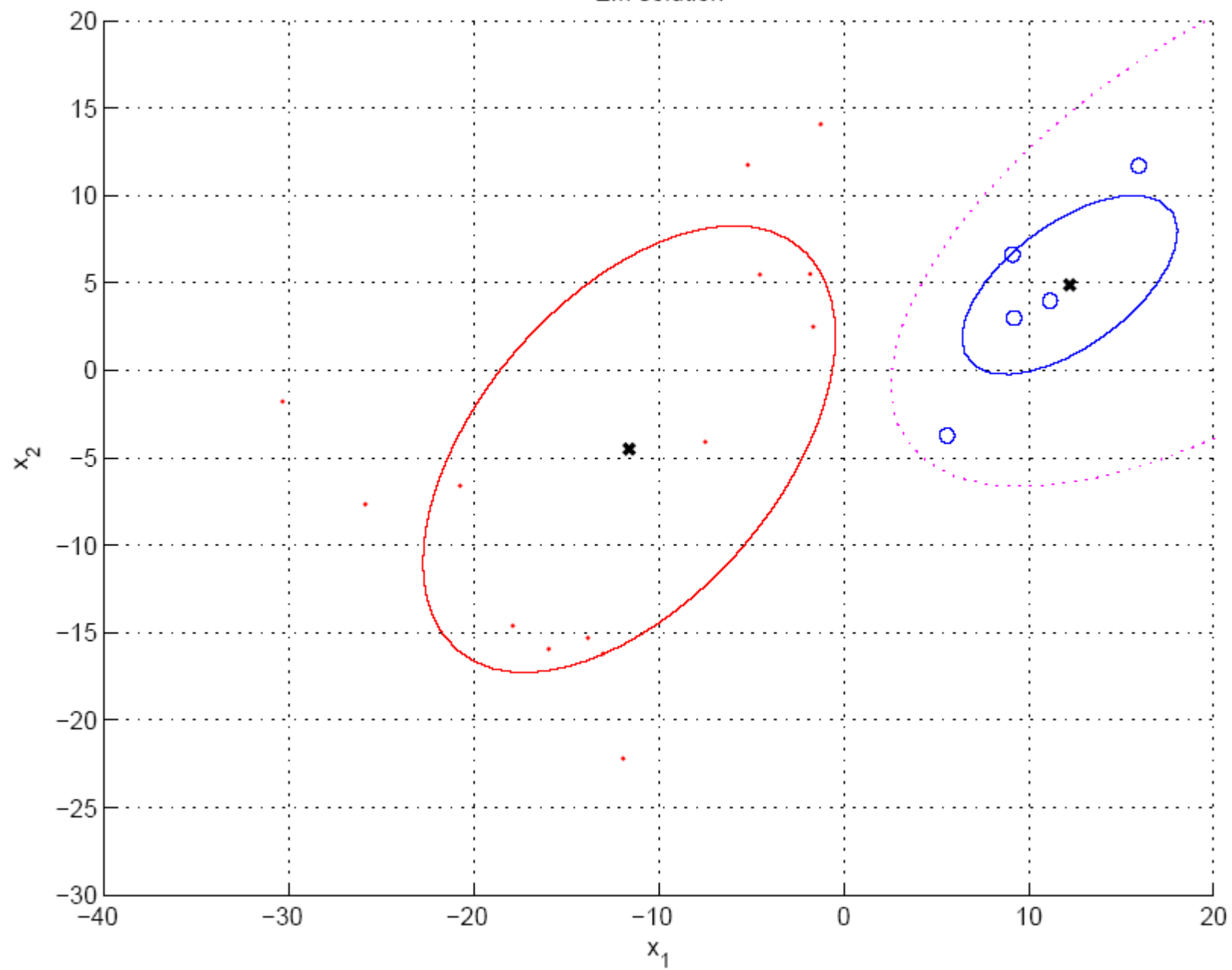
Soft assignment of a sample to a class

$$h_i^t = \frac{\pi_i |\mathbf{S}_i|^{-1/2} \exp[-(1/2)(\mathbf{x}^t - \mathbf{m}_i)^T \mathbf{S}_i^{-1} (\mathbf{x}^t - \mathbf{m}_i)]}{\sum_j \pi_j |\mathbf{S}_j|^{-1/2} \exp[-(1/2)(\mathbf{x}^t - \mathbf{m}_j)^T \mathbf{S}_j^{-1} (\mathbf{x}^t - \mathbf{m}_j)]}$$

# Practice implementation of GMM

- EM is initialized by k-means
- Once done, use  $m_i$  and samples associated with each cluster as to estimate the initial parameters used for mixture of Gaussian distributions
- Once done-learning, the GMM model can be used for 'clustering' of samples

EM solution



# After Clustering

- Dimensionality reduction methods find correlations between features and group features
- Clustering methods find similarities between instances and group instances
- Allows knowledge extraction through
  - number of clusters,
  - prior probabilities,
  - cluster parameters, i.e., center, range of features.

Example: CRM, customer segmentation