

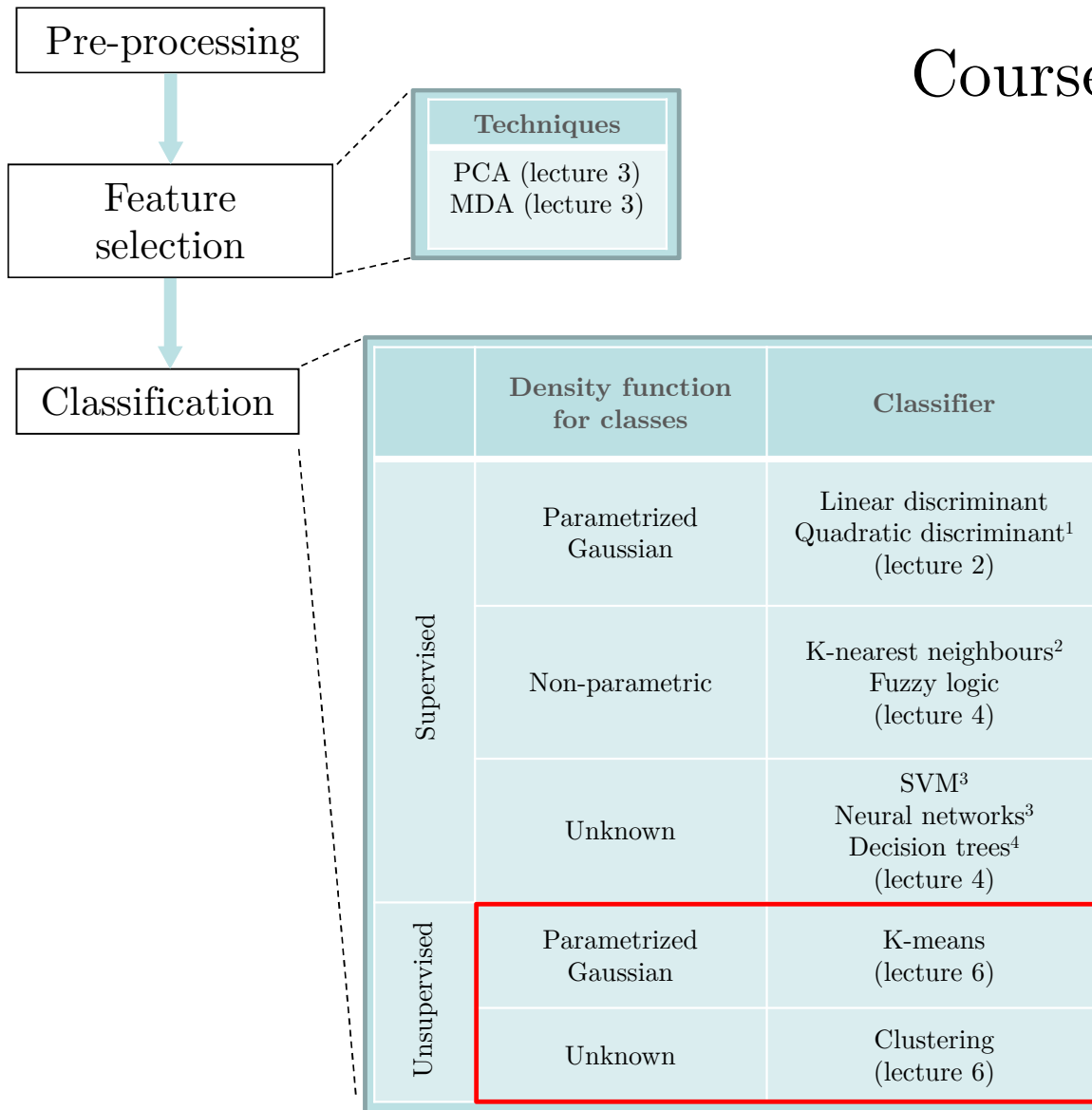
# Chapter 6

## Unsupervised learning

**Recommended bibliography:** *Pattern Classification (2nd ed)* by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000, Chapter 4

**Credits:** Some figures are taken from *Pattern Classification (2nd ed)* by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000 with the permission of the authors

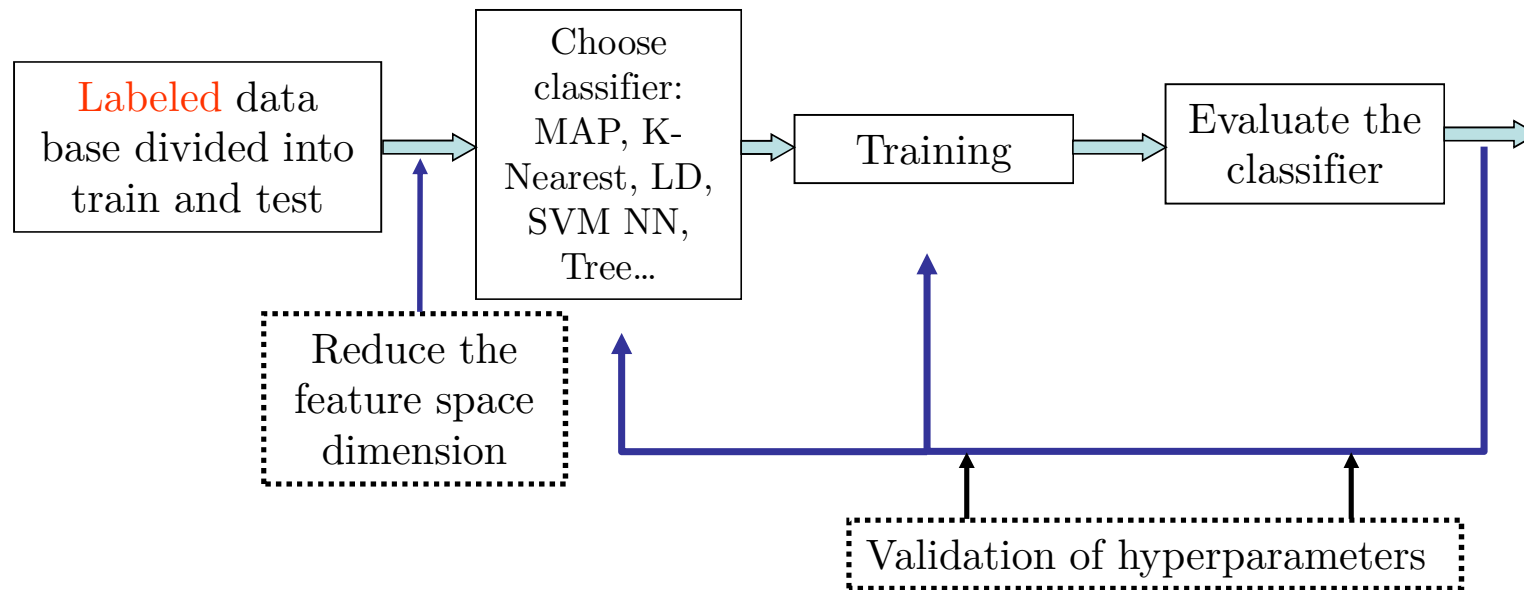
# Course overview



1. Useful only if covariance matrices are not rank deficient.
2. Useful with the number of features is very large, even larger than the number of training vectors.
3. Imposes a structure to the classifier irrespective of the training data base.
4. Useful when non-numeric features are present.

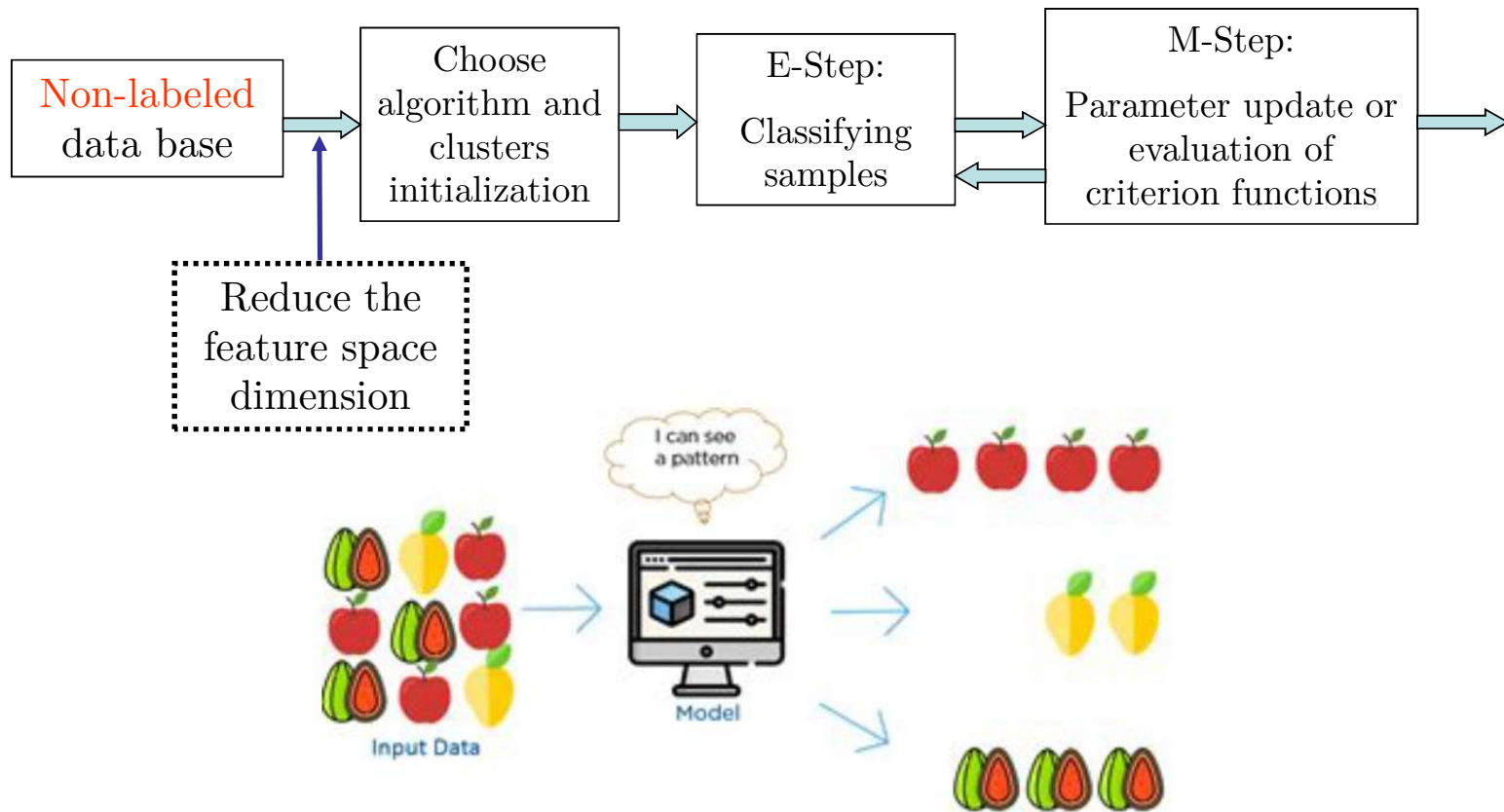
# SUPERVISED METHODS

## Labeled database

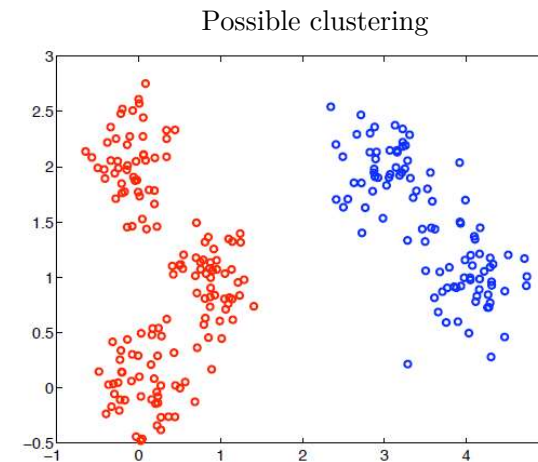
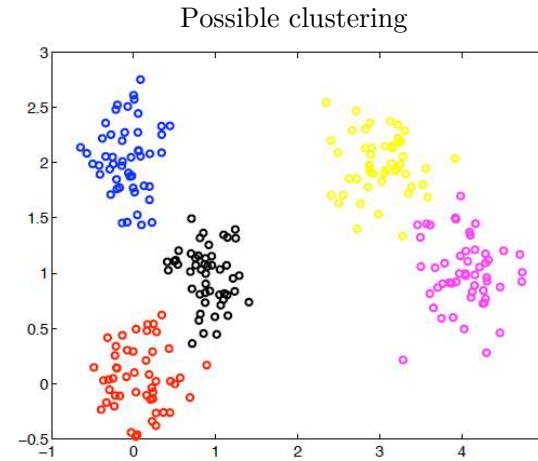
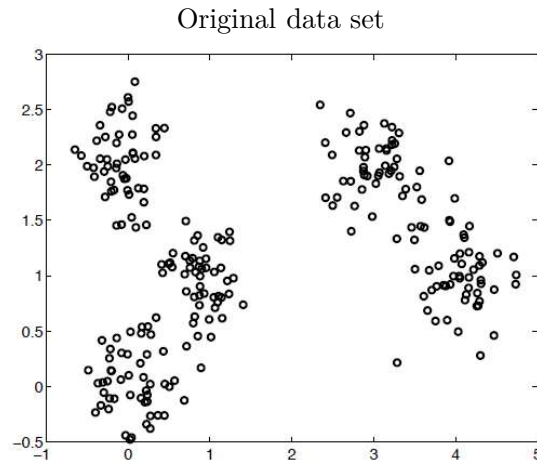


# UNSUPERVISED METHODS

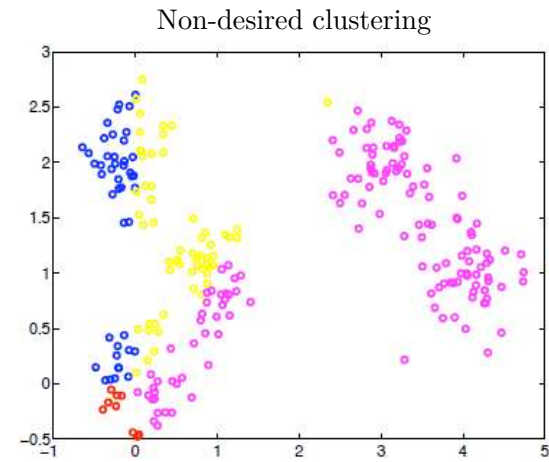
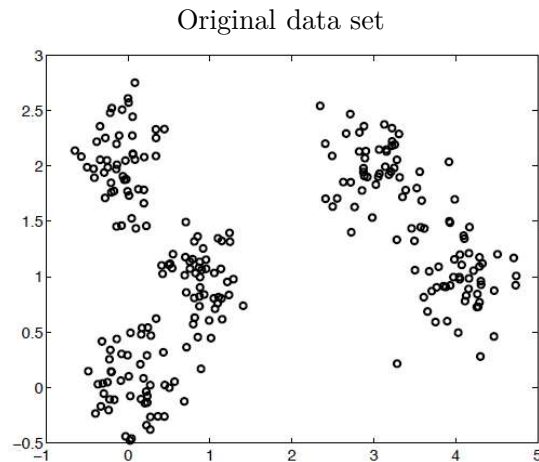
## Non-labeled database



**Goal: Find structure in the data by isolating groups of vectors that are similar in some well-defined sense...**



**Goal: Find structure in the data by isolating groups of vectors that are similar in some well-defined sense...**



## Some good reasons to learn from unlabeled data...

- Collecting and labeling a large data base can be very costly
- In a supervised problem, model each class as a mixture of clusters  $f(\mathbf{x}|\omega)$ , and find posterior probabilities that can be used for categorisation by another classifier (e.g. MAP, MBR, NN,...)
- Gain some insight on the structure of the data beyond what can be obtained from scatter plots
- Dimensionality reduction
- Anomaly detection (detection of outlier vectors)
- Tracking the characteristics of the patterns in time: some patterns may jump from one class to another

## Two approaches...

- **Parametric methods** assume a certain density function for the clusters
- **Non-parametric methods**, based on some similarity measure of the vectors to be clustered



# CONTENTS

## **6.1 Parametric methods**

6.1.1 Mixture densities and identifiability

6.1.2 ML estimates and the EM-algorithm

6.1.3 K-means clustering

6.1.4 Fuzzy K-means clustering

## **6.2 Clustering**

## 1.1 MIXTURE DENSITIES AND IDENTIFIABILITY

Assumptions:

1. The number  $c$  of clusters
2. Prior probabilities for each cluster are unknown  $\Pr(\omega_j) \quad j = 1, \dots, c$   
(mixing parameters).
3. The form of the class-conditional probabilities densities are known  $f_{\mathbf{x}|\omega_j, \theta_j}(\mathbf{x}|\omega_j, \theta_j)$
4. The values for parameters are unknown  $\theta = [\theta_1, \dots, \theta_c]$
5. The category labels are unknown: **unsupervised**

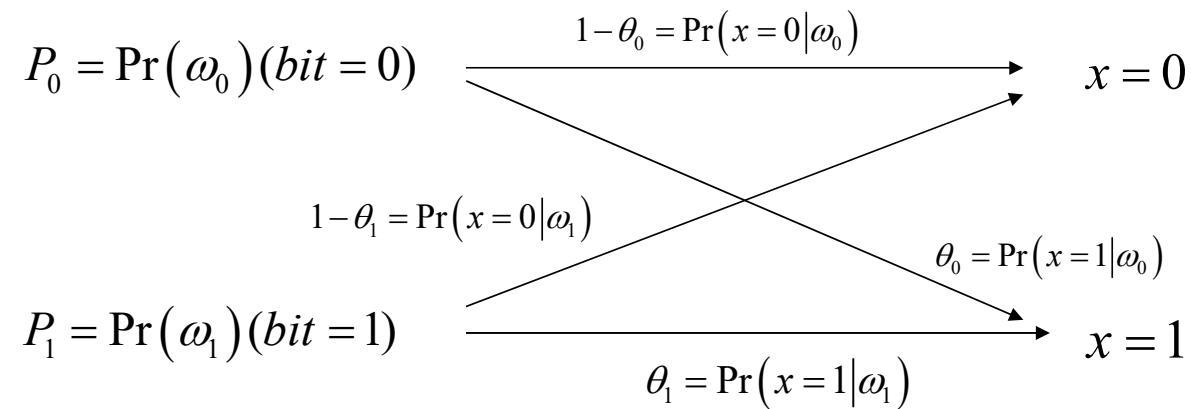
Mixture density:

$$f_{\mathbf{x}|\boldsymbol{\theta}}(\mathbf{x}|\boldsymbol{\theta}) = \sum_{j=1}^c f_{\mathbf{x}|\omega_j, \boldsymbol{\theta}_j}(\mathbf{x}|\omega_j, \boldsymbol{\theta}_j) \Pr(\omega_j)$$

1. For the moment, assumed that only the parameter vector  $\boldsymbol{\theta}$  is unknown (priors are known).
2. Necessary conditions for identifiability:

$$f_{\mathbf{x}|\boldsymbol{\theta}}(\mathbf{x}|\boldsymbol{\theta}) \neq f_{\mathbf{x}|\boldsymbol{\theta}'}(\mathbf{x}|\boldsymbol{\theta}') \quad \boldsymbol{\theta} = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_c]$$

## Example 1: Identifiability in the binary symmetric channel



Vector of parameters:  $\boldsymbol{\theta} = \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix}$

Mixture density (probability):

$$\Pr(x|\boldsymbol{\theta}) = P_0 \theta_0^x (1 - \theta_0)^{1-x} + P_1 \theta_1^x (1 - \theta_1)^{1-x}$$

## 1.2 ML ESTIMATES AND THE EM-ALGORITHM

Likelihood of the statistical independent observed samples:

$$D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$$

$$f_D(D|\boldsymbol{\theta}) = \prod_{k=1}^n f_{\mathbf{x}}(\mathbf{x}_k|\boldsymbol{\theta}) \quad l = \sum_{k=1}^n \ln f_{\mathbf{x}}(\mathbf{x}_k|\boldsymbol{\theta})$$

$$f_{\mathbf{x}}(\mathbf{x}_k|\boldsymbol{\theta}) = \sum_{j=1}^c f_{\mathbf{x}}(\mathbf{x}_k|\omega_j, \boldsymbol{\theta}_j) \Pr(\omega_j)$$

Assume the unknowns  $\boldsymbol{\theta}_i, \boldsymbol{\theta}_j$  are independent variables, then the ML criterion provides  $c$  equations:

$$\nabla_{\boldsymbol{\theta}_i} l = \sum_{k=1}^n \Pr(\omega_i|\mathbf{x}_k, \hat{\boldsymbol{\theta}}) \nabla_{\boldsymbol{\theta}_i} \ln f_{\mathbf{x}}(\mathbf{x}_k|\omega_i, \hat{\boldsymbol{\theta}}_i) = \mathbf{0} \quad i = 1, \dots, c$$



Proof:

$$\begin{aligned}\nabla_{\boldsymbol{\theta}_i} l &= \sum_{k=1}^n \frac{1}{f_{\mathbf{x}}(\mathbf{x}_k | \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}_i} f_{\mathbf{x}}(\mathbf{x}_k | \boldsymbol{\theta}) \\&= \sum_{k=1}^n \frac{1}{f_{\mathbf{x}}(\mathbf{x}_k | \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}_i} \left( \sum_{j=1}^c f_{\mathbf{x}}(\mathbf{x}_k | \omega_j, \boldsymbol{\theta}_j) \Pr(\omega_j) \right) \\&= \sum_{k=1}^n \frac{1}{f_{\mathbf{x}}(\mathbf{x}_k | \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}_i} \left( f_{\mathbf{x}}(\mathbf{x}_k | \omega_i, \boldsymbol{\theta}_i) \Pr(\omega_i) \right) \\&= \sum_{k=1}^n \frac{\Pr(\omega_i)}{f_{\mathbf{x}}(\mathbf{x}_k | \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}_i} f_{\mathbf{x}}(\mathbf{x}_k | \omega_i, \boldsymbol{\theta}_i) = \quad \leftarrow \text{Bayes theorem} \\&= \sum_{k=1}^n \Pr(\omega_i | \mathbf{x}_k, \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}_i} \ln f_{\mathbf{x}}(\mathbf{x}_k | \omega_i, \boldsymbol{\theta}_i) = \mathbf{0} \quad i = 1, \dots, c\end{aligned}$$

**The EM algorithm, including estimation of the unknown prior probabilities:** (convergence guaranteed, without proof)

1. Compute conditioned probability for classes (Bayes rule)

$$\hat{\Pr}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}}) = \frac{f_{\mathbf{x}}(\mathbf{x}_k | \omega_i, \hat{\boldsymbol{\theta}}_i) \hat{\Pr}(\omega_i)}{\sum_{j=1}^c f_{\mathbf{x}}(\mathbf{x}_k | \omega_j, \hat{\boldsymbol{\theta}}_j) \hat{\Pr}(\omega_j)}$$

Only depends on parameters of class  $i$

2. Estimate vector parameter estimates (gradient equal to zero)

$$\sum_{k=1}^n \hat{\Pr}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}}) \nabla_{\boldsymbol{\theta}_i} \ln f_{\mathbf{x}}(\mathbf{x}_k | \omega_i, \hat{\boldsymbol{\theta}}_i) = \mathbf{0} \quad i = 1, \dots, c$$

ML criterion, averaged with estimated a posteriori probabilities

3. Estimate prior probabilities

$$\hat{\Pr}(\omega_i) = \frac{1}{n} \sum_{k=1}^n \hat{\Pr}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}})$$

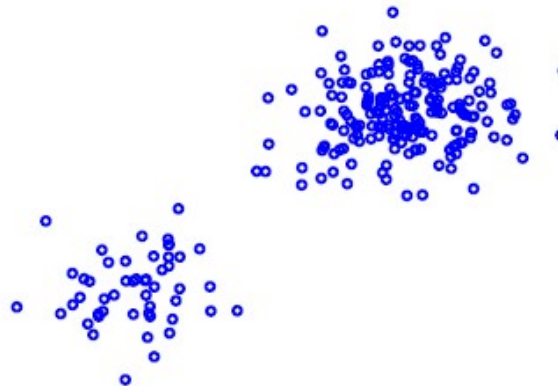
This is a chicken-and-egg problem, solved in an iterative way...

## Application to Gaussian distributions:

$$\ln f_{\mathbf{x}}(\mathbf{x}_k | \omega_i, \boldsymbol{\theta}_i) = \frac{1}{2} \ln |\mathbf{C}_i^{-1}| - \frac{d}{2} \ln(2\pi) - \frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu}_i)^T \mathbf{C}_i^{-1} (\mathbf{x}_k - \boldsymbol{\mu}_i)$$

Parameters to estimate:  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_c)$        $\boldsymbol{\theta}_i = (\boldsymbol{\mu}_i, \mathbf{C}_i)$   
and the prior probabilities  $\Pr\{\omega_i\} \quad i = 1, \dots, c$

An example for  $c = 2$







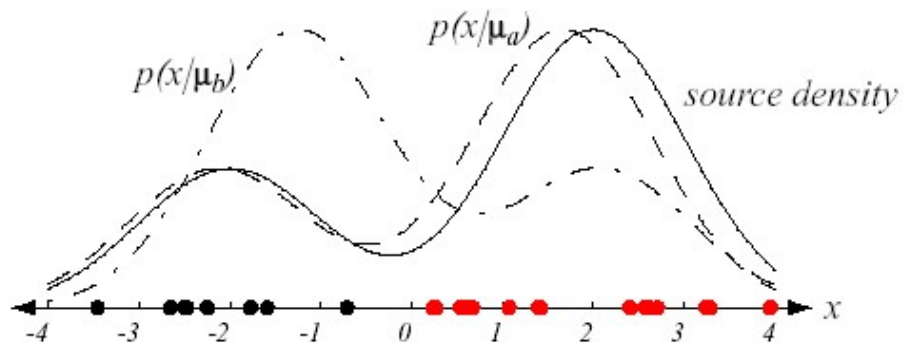
## Case 1: Unknown mean vectors and prior probabilities, covariance known

$$\sum_{k=1}^n \Pr(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}}) \nabla_{\boldsymbol{\mu}_i} \ln f_{\mathbf{x}}(\mathbf{x}_k | \omega_i, \boldsymbol{\mu}_i) = \sum_{k=1}^n \Pr(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}}) \mathbf{C}_i^{-1} (\mathbf{x}_k - \boldsymbol{\mu}_i) = 0 \quad i = 1, \dots, c$$

After multiplying by  $\mathbf{C}_i$  and rearranging...

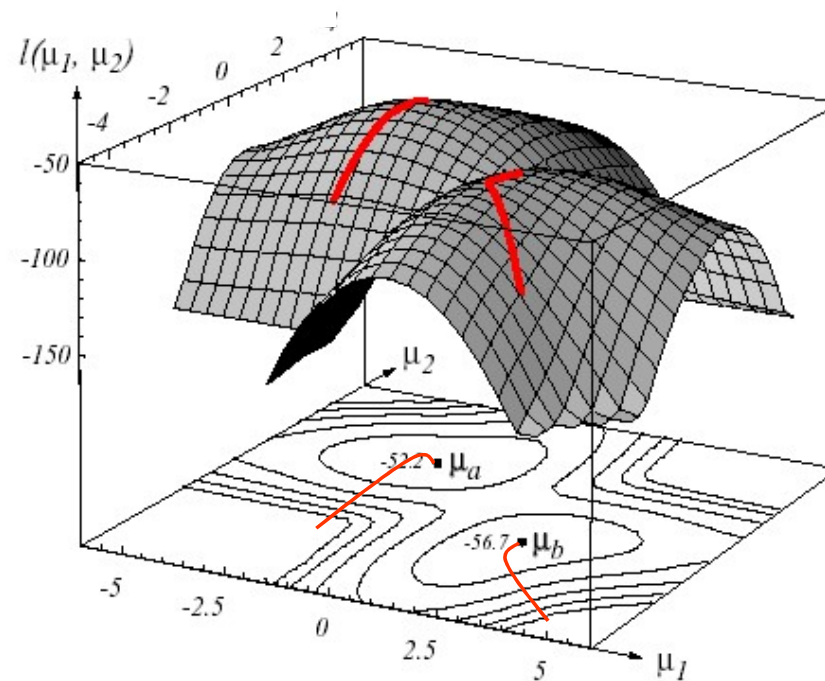
$$\hat{\boldsymbol{\mu}}_i = \frac{\sum_{k=1}^n \Pr(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\mu}}) \mathbf{x}_k}{\sum_{k=1}^n \Pr(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\mu}})} \longrightarrow \hat{\boldsymbol{\mu}}_i(j+1) = \frac{\sum_{k=1}^n \Pr(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\mu}}(j)) \mathbf{x}_k}{\sum_{k=1}^n \Pr(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\mu}}(j))}$$

This is a hill-climbing procedure that may converge to local maxima, also called Expectation-Maximization (EM) algorithm.



On each iteration, the EM-algorithm monotonically increases the (log-) likelihood of the  $n$  samples:

$$\ln f_{\mathbf{x}}(D|\boldsymbol{\theta})$$



## Case 2: Unknown mean vectors, covariance matrices and prior probabilities



1. Expectation (E-Step): estimate the probability of belonging to each class

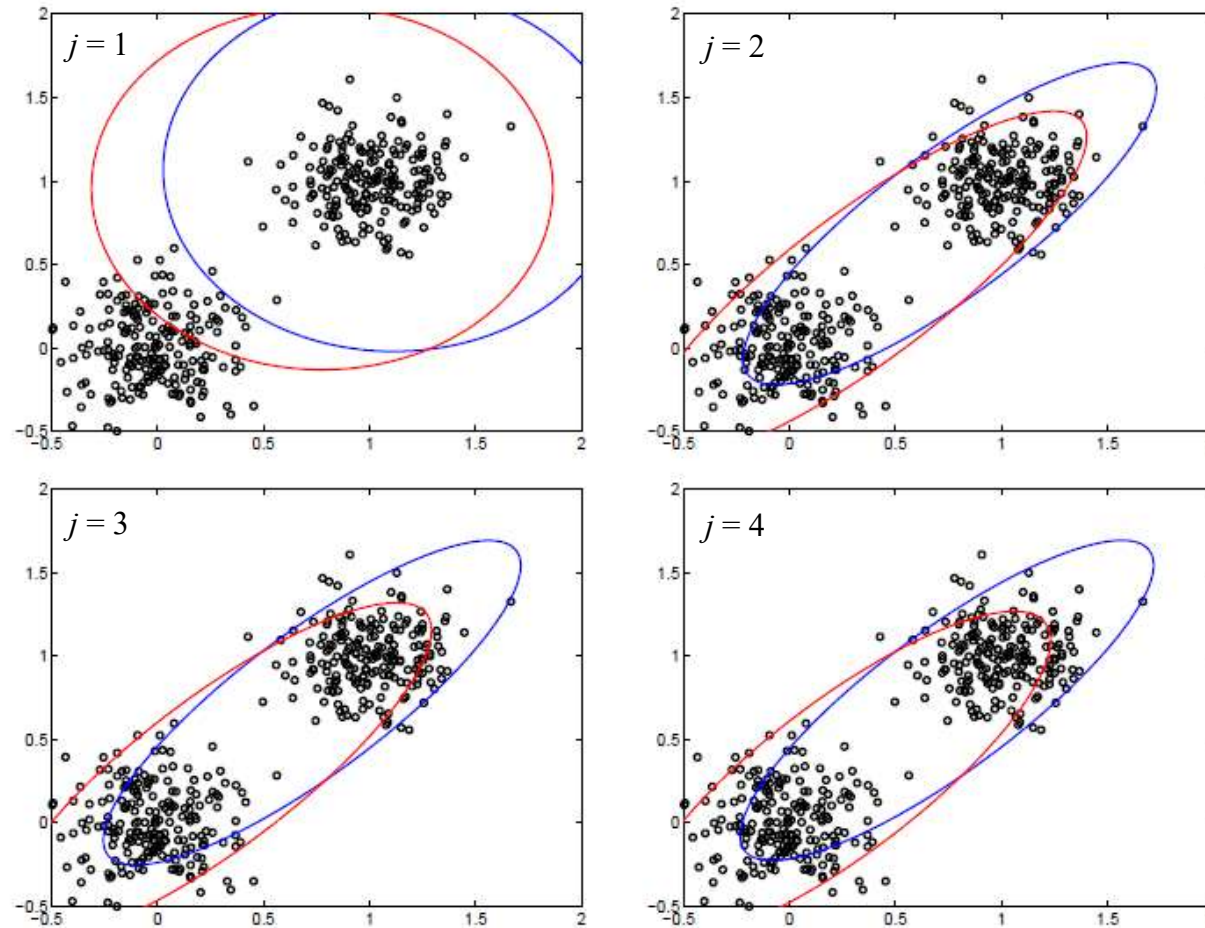
$$\hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}}) = \frac{|\hat{\mathbf{C}}_i^{-1}|^{\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x}_k - \hat{\boldsymbol{\mu}}_i)^T \hat{\mathbf{C}}_i^{-1} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_i)\right) \hat{\text{Pr}}(\omega_i)}{\sum_{j=1}^c |\hat{\mathbf{C}}_j^{-1}|^{\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x}_k - \hat{\boldsymbol{\mu}}_j)^T \hat{\mathbf{C}}_j^{-1} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_j)\right) \hat{\text{Pr}}(\omega_j)}$$

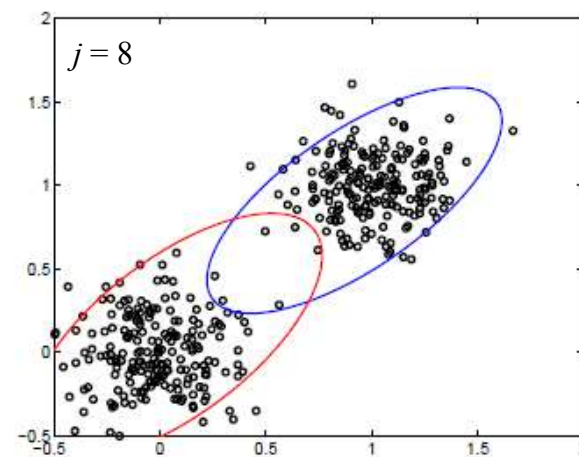
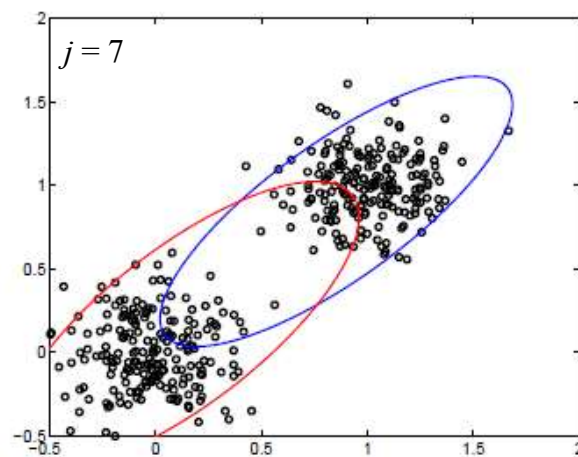
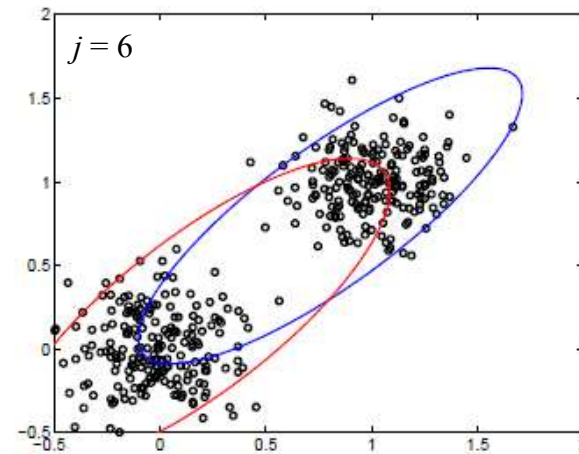
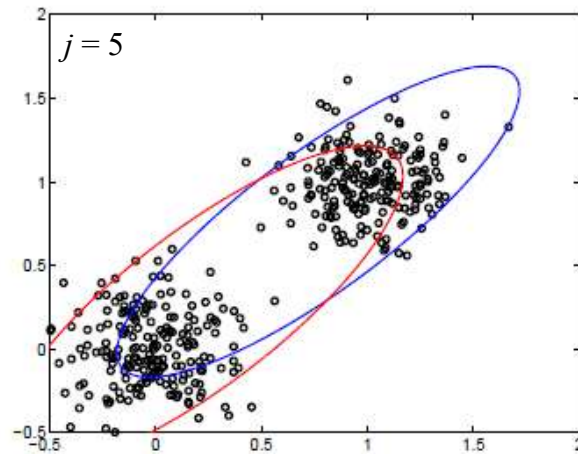
2. Maximization (M-Step): estimate the parameters of each class via ML

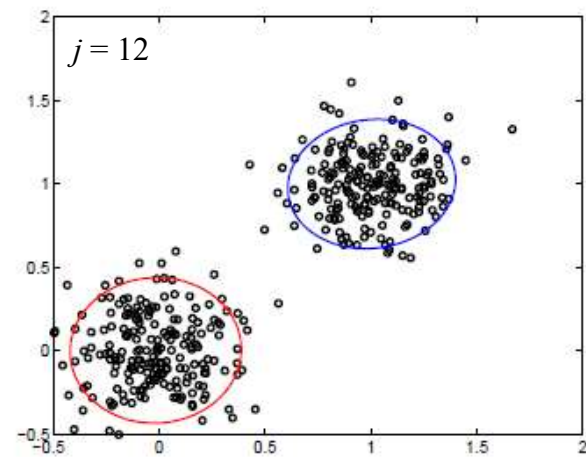
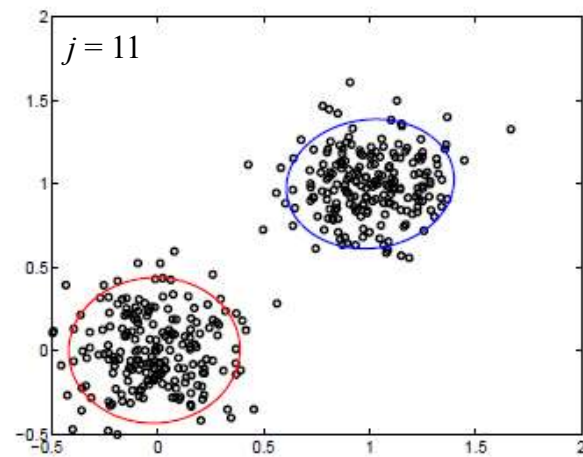
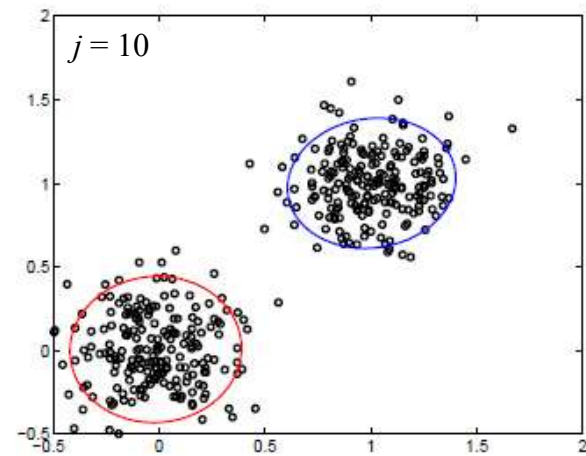
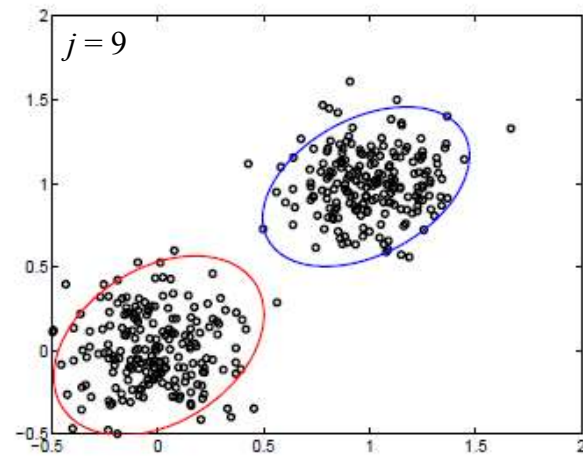
$$\hat{\boldsymbol{\mu}}_i = \frac{\sum_{k=1}^n \hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}}) \mathbf{x}_k}{\sum_{k=1}^n \hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}})} \quad \hat{\mathbf{C}}_i = \frac{\sum_{k=1}^n \hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}}) (\mathbf{x}_k - \hat{\boldsymbol{\mu}}_i)(\mathbf{x}_k - \hat{\boldsymbol{\mu}}_i)^T}{\sum_{k=1}^n \hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}})}$$
$$\hat{\text{Pr}}\{\omega_i\} = \frac{1}{n} \sum_{k=1}^n \hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}})$$

Iterations between E and M steps stop when estimated parameters do not change.

## Example 2: sequence of 12 EM iterations for two Gaussian classes







## How to select the appropriate number of clusters $c$ ?

We can find  $m$  that minimize the following asymptotic approximation to the description length criterion:

$$c = \arg \min_m DL(m)$$

$$DL(m) \simeq -\ln f_{\mathbf{x}} \left( D \middle| \hat{\boldsymbol{\theta}}_{m,ML} \right) + \frac{d_m}{2} \ln(n)$$

Number of parameters in the  $m$ -clusters mixture

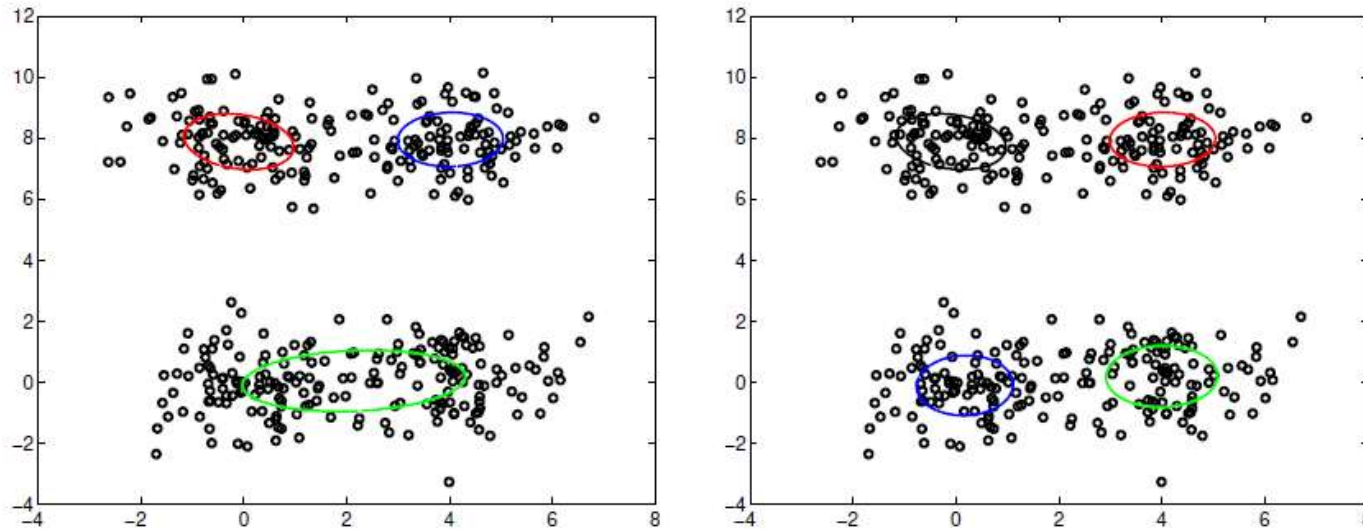
Number of vectors in  $D$

ML estimates of the  $m$ -clusters mixture

This is called the Bayesian Information Criterion (BIC).



### Example 3: number of mixture components



$m=1$ ,  $-\log P(\text{data})=2017.38$ ,  $\text{penalty}=14.98$ ,  $DL=2032.36$   
 $m=2$ ,  $-\log P(\text{data})=1712.69$ ,  $\text{penalty}=32.95$ ,  $DL=1745.65$   
 $m=3$ ,  $-\log P(\text{data})=1678.56$ ,  $\text{penalty}=50.93$ ,  $DL=1729.49$  ←  
 $m=4$ ,  $-\log P(\text{data})=1649.08$ ,  $\text{penalty}=68.90$ ,  $DL=1717.98$  ←

$$\frac{\partial}{\partial m} \ln(n)$$



## 1.3 K-MEANS CLUSTERING

We can reduce convergence time by doing a hard classification of samples at the E-step. In a simple case of equal diagonal covariance matrices...

$$\boldsymbol{\theta}_i = \boldsymbol{\mu}_i \quad \mathbf{C}_i = \sigma^2 \mathbf{I}$$

Euclidean distance of  $\mathbf{x}_k$  to each estimated  $\boldsymbol{\mu}_j$ , obtained from posterior probability in slide 15 (if ignoring  $\hat{\Pr}(\omega_i)$ )

$$\hat{\Pr}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\mu}}_i) = \begin{cases} 1 & d_e(\mathbf{x}_k, \hat{\boldsymbol{\mu}}_i) < d_e(\mathbf{x}_k, \hat{\boldsymbol{\mu}}_j) \quad \forall j \neq i \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{\Pr}\{\omega_i\} = \frac{1}{n} \sum_{k=1}^n \hat{\Pr}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\mu}}_i) = \frac{n_i}{n}$$

$$\hat{\boldsymbol{\mu}}_i = \frac{1}{n_i} \sum_{k=1}^n \hat{\Pr}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\mu}}_i) \mathbf{x}_k$$

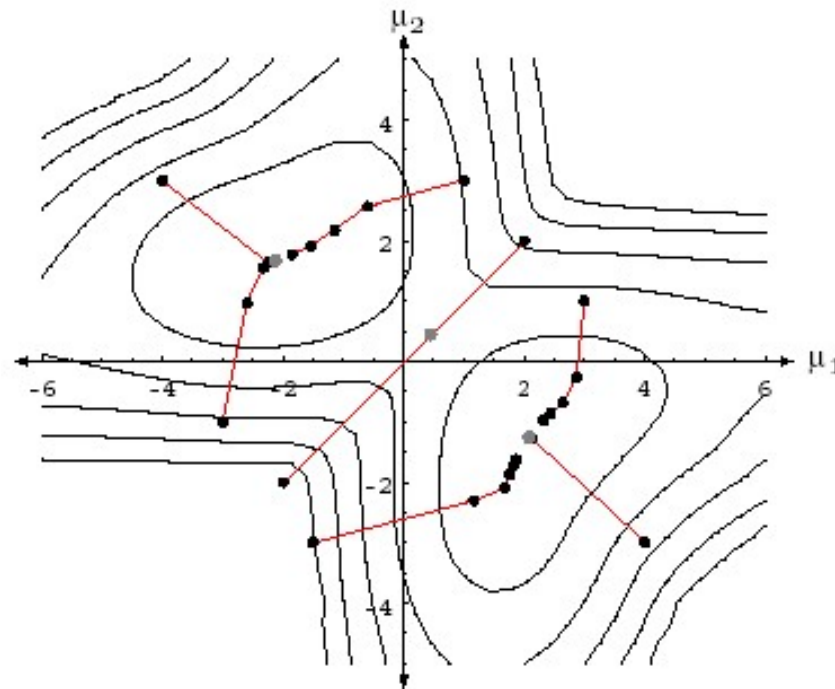
- K-Means clustering algorithm

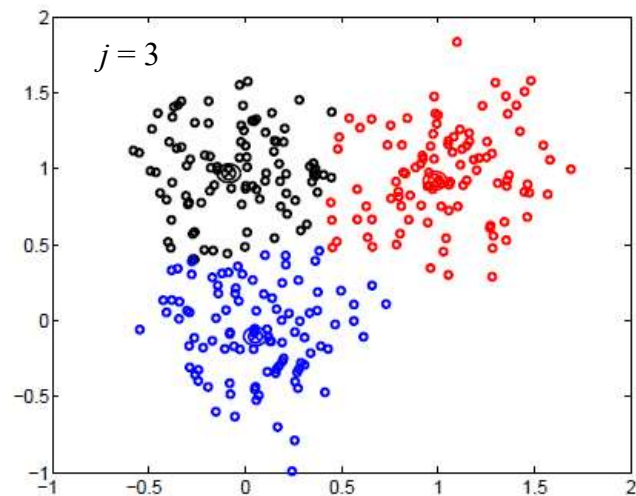
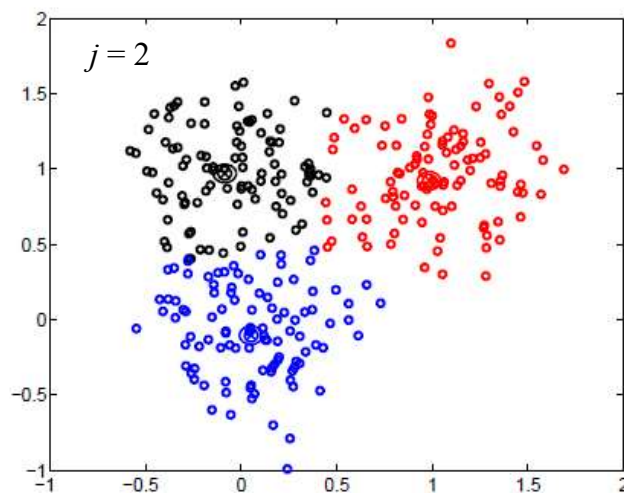
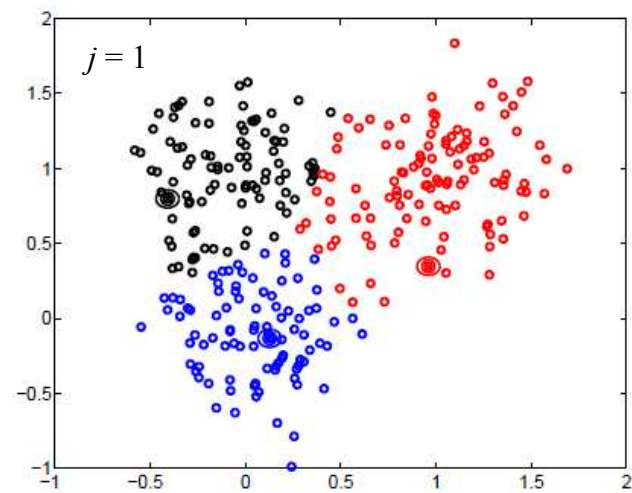
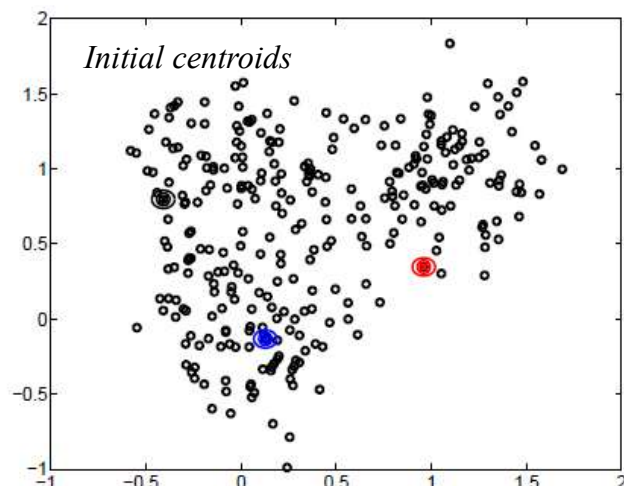
```

1 begin initialize  $n, c, \mu_1, \mu_2, \dots, \mu_c$ 
2   do classify  $n$  samples according to nearest  $\mu_i$       % E-step
3   recompute  $\mu_i$       % M-step
4   until no change in  $\mu_i$ 
5   return  $\mu_1, \mu_2, \dots, \mu_c$ 
6 end

```

Convergence is achieved  
in a few iterations





## 1.4 FUZZY K-MEANS CLUSTERING



We can relax the condition of every sample belonging to just one cluster, and assume some fuzzy membership, related to the a posteriori probabilities  $\hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}})$

We seek to minimize a heuristic function:

integer value larger than 1  
to adjust blending of  
clusters

$$J_{\text{Fuzzy}} = \sum_{i=1}^c \sum_{k=1}^n \hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}})^b d_e(\mathbf{x}_j, \hat{\boldsymbol{\mu}}_i)$$

leading to the following optimum (mean vector estimated only):

$$\hat{\boldsymbol{\mu}}_i = \frac{\sum_{k=1}^n \hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}})^b \mathbf{x}_k}{\sum_{k=1}^n \hat{\text{Pr}}(\omega_i | \mathbf{x}_k, \hat{\boldsymbol{\theta}})^b}$$



...for other values of  $b$

For  $b = 1$  we get the k-means, what happens for other values of  $b$ ?

- $b = 2$  Higher weight when computing  $\hat{\mu}_i$  for those  $\mathbf{x}_i$  with higher probability

- $b \rightarrow \infty \quad \hat{\mu}_i = \arg \max_{\mathbf{x}_k} \frac{\sum_{k=1}^n \hat{\Pr}(\omega_i | \mathbf{x}_k, \hat{\theta})^b \mathbf{x}_k}{\sum_{k=1}^n \hat{\Pr}(\omega_i | \mathbf{x}_k, \hat{\theta})^b}$

Each mean is computed using a single vector  $\mathbf{x}_k$

The value of  $b$  can be used to robustify the determination of centroids in the presence of outliers.

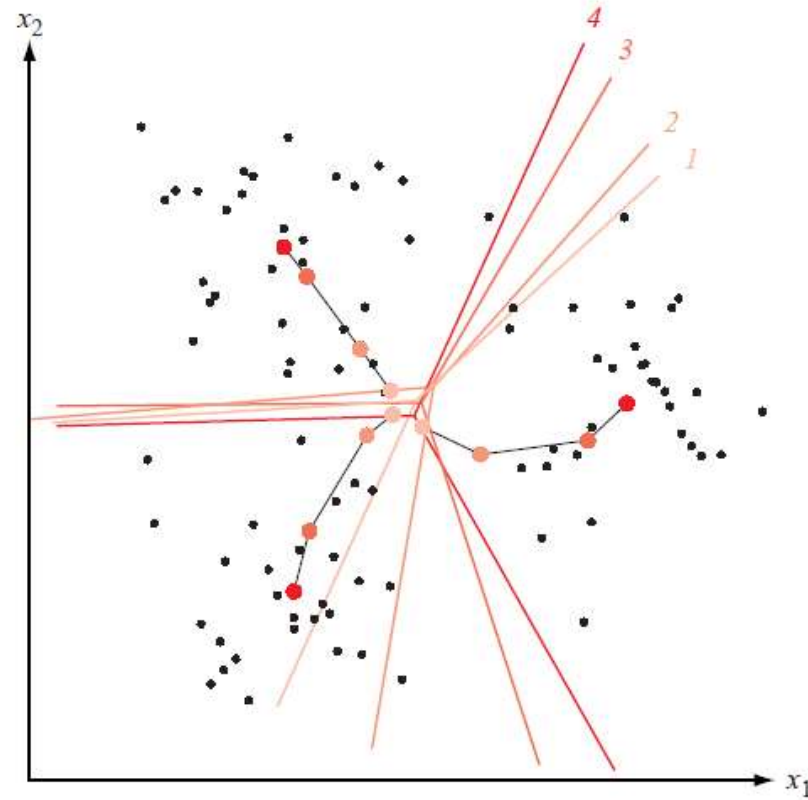
**Exercise:** Derive the estimate of  $\hat{\Pr}(\omega_i | \mathbf{x}_k, \hat{\theta})$ , have in mind the sum from  $i = 1$  to  $c$  must be 1.



```
1 begin initialize  $n, \mu_1, \dots, \mu_c, P(\omega_i | \mathbf{x}_j), i = 1 \dots, c; j = 1, \dots, n$   
2     normalize probabilities of cluster memberships  
3     do classify  $n$  samples according to nearest  $\mu_i$   
4         recompute  $\mu_i$   
5         recompute  $P(\omega_i | \mathbf{x}_j)$   
6     until no change in  $\mu_i$  and  $P(\omega_i | \mathbf{x}_j)$   
7     return  $\mu_1, \mu_2, \dots, \mu_c$   
8 end
```

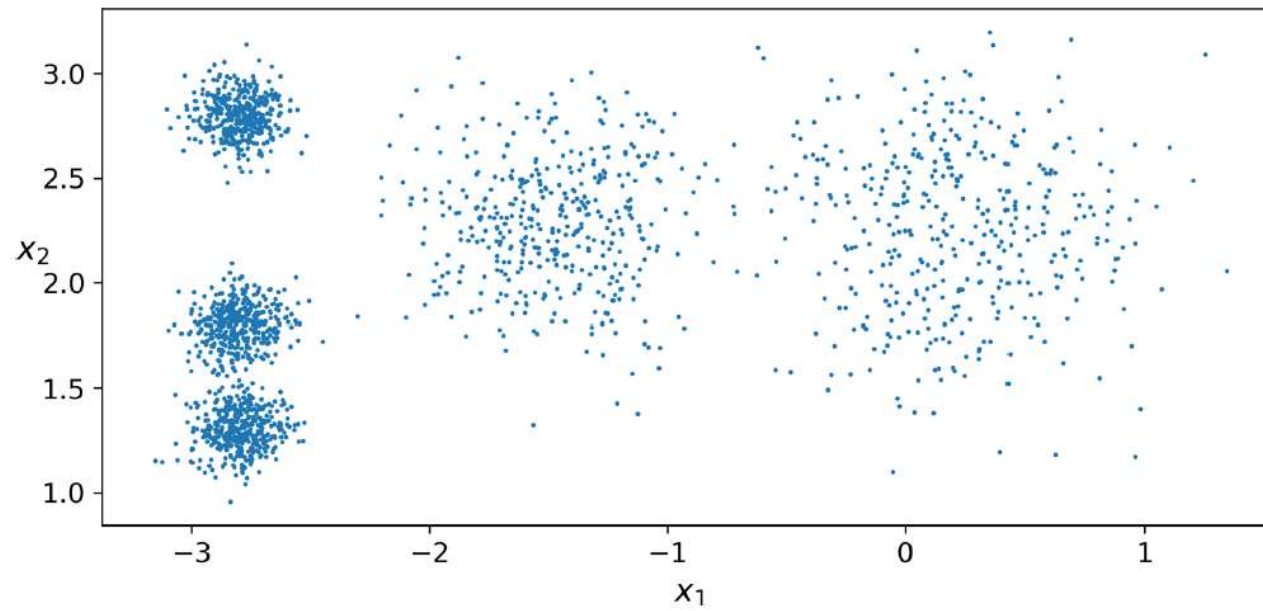
Warning! If  $c$  is incorrectly specified, serious problems may arise

For an application in vector quantization of a  $n$ -dimensional real valued vector: *J. Proakis: "Digital Communications", Chapter 3: Source Coding.*



Evolution of the regions after 4 iterations of the fuzzy k-means ( $b=2$ ).

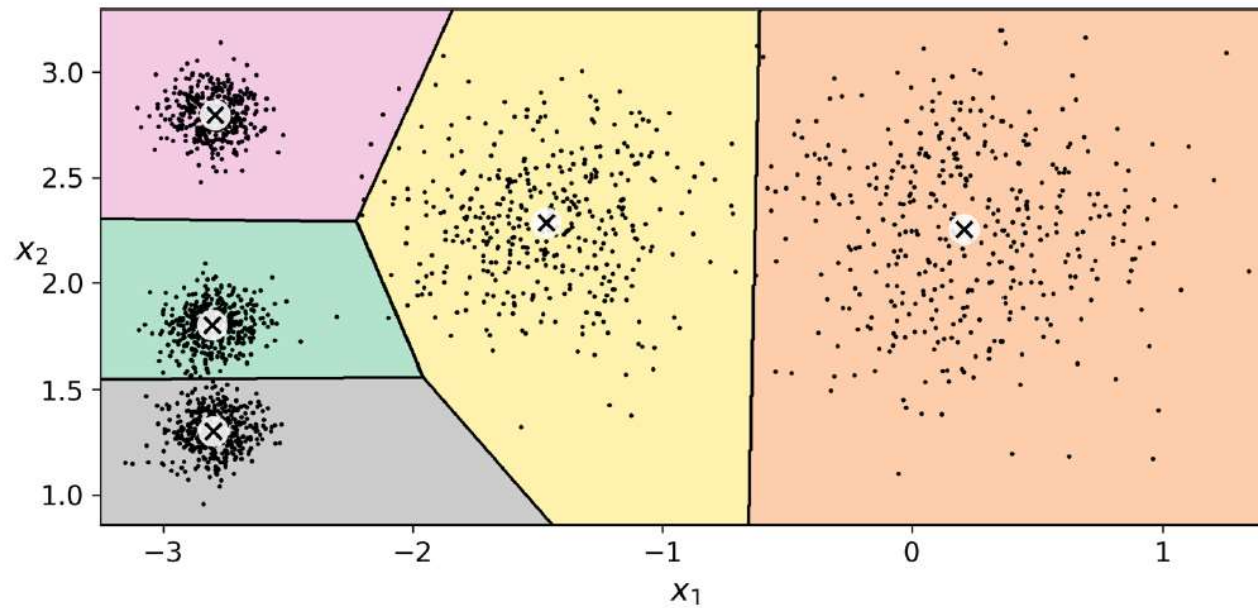
# K-means limitations





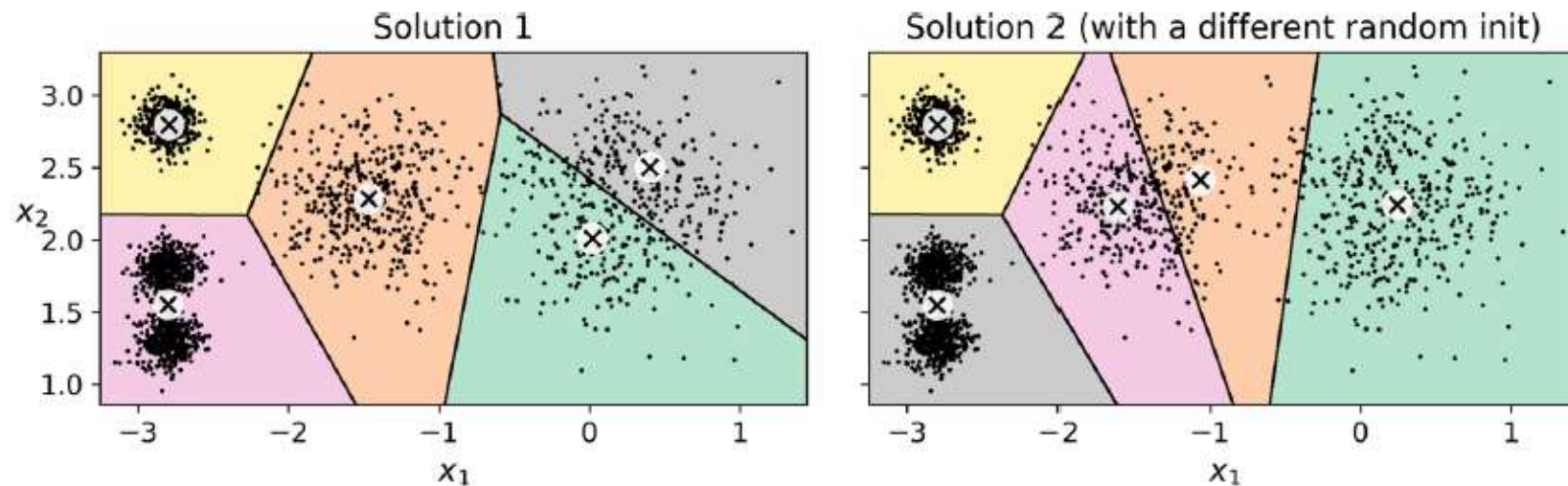
# K-means limitations

Cluster size: difficulties if clusters of different sizes and densities.



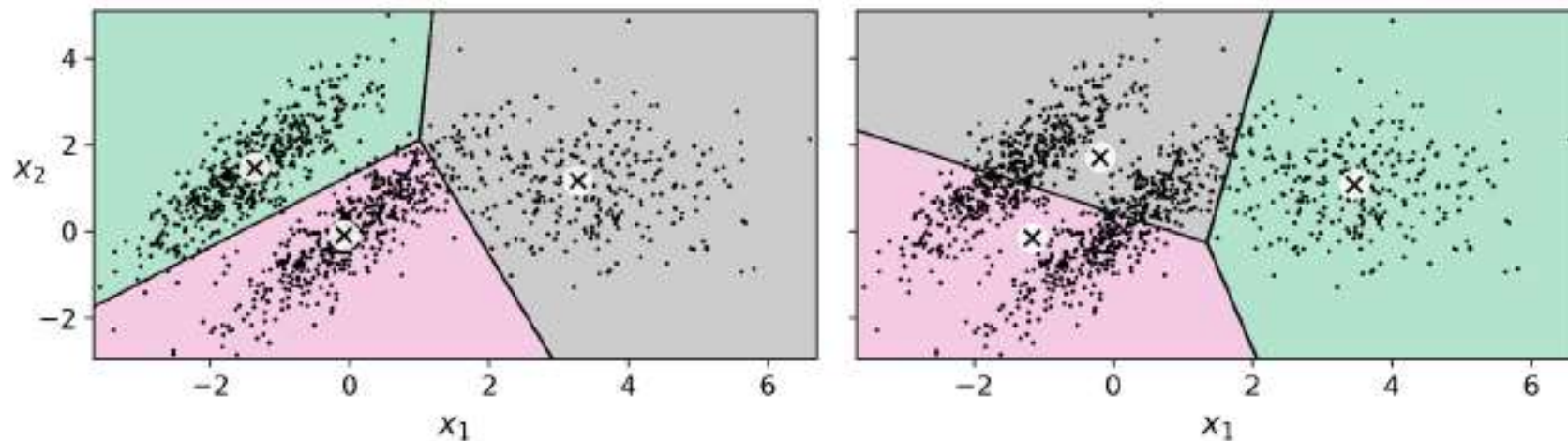
# K-means limitations

Convergence is guaranteed, but possibly to local minima, depending on the initialization.



# K-means limitations

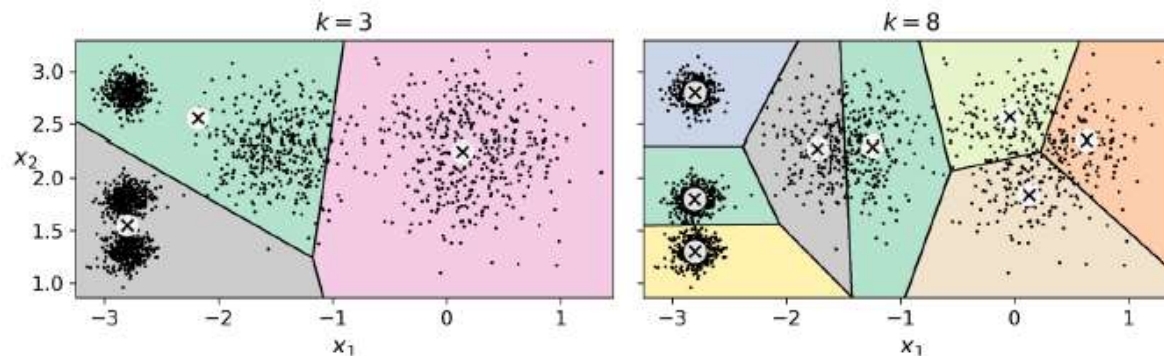
Cluster shape: difficulties if clusters are not of spherical shape.



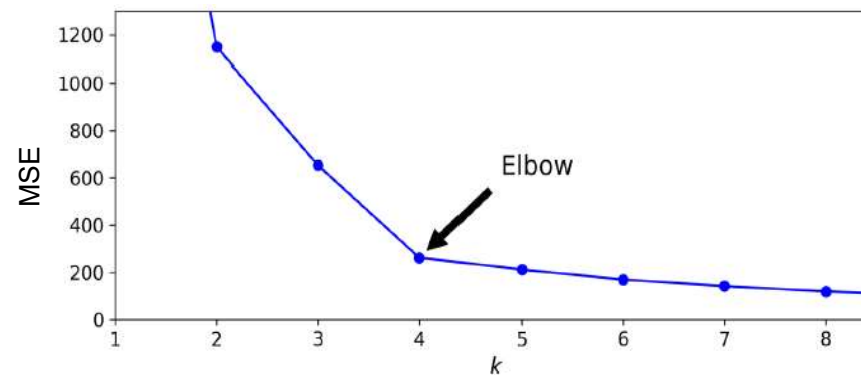
Gaussian mixture models with arbitrary  $\mathbf{C}$  matrices will perform much better in this case (check the parametric models).

## 2.4 OPTIMUM NUMBER OF CLUSTERS

Let us revisit the estimation of the number of clusters. Examples of two bad selections:



The overall MSE as a function of the number of clusters is not a good criterion, as it is a decreasing function with  $k$ :



A more precise approach is to use the *mean silhouette score*:

1. For  $h = 2$  to  $K$
2.     Compute the k-means clustering for  $h$  clusters
3.     For  $l = 1$  to  $h$
4.         Compute the silhouette score of each cluster as

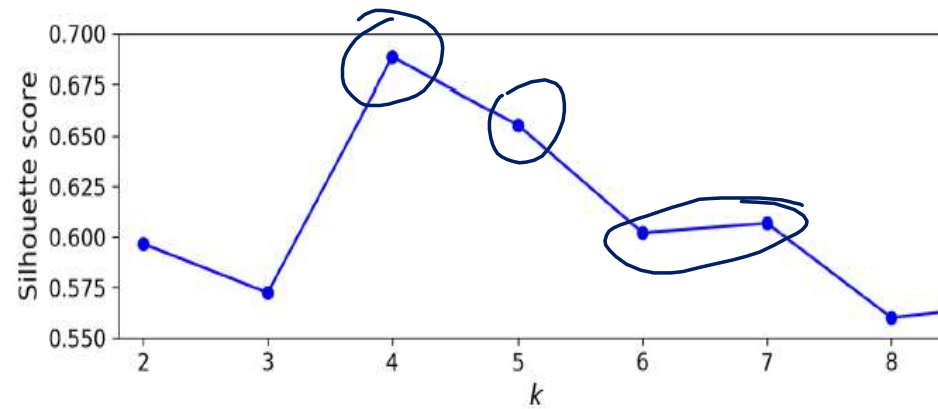
$$a_i = \frac{1}{N_l} \sum_{\mathbf{x}_j \in \omega_l} d(\mathbf{x}_i, \mathbf{x}_j) \quad \leftarrow \text{Average distance of vector } \mathbf{x}_i \text{ to other vectors within its own cluster}$$

$$b_i = \frac{1}{N_m} \sum_{\mathbf{x}_j \in \omega_m} d(\mathbf{x}_i, \mathbf{x}_j) \quad \leftarrow \text{Average distance of vector } \mathbf{x}_i \text{ to other vectors of the nearest cluster } m \text{ (the one minimizing } b, \text{ which may be different for every vector } \mathbf{x}_i)$$

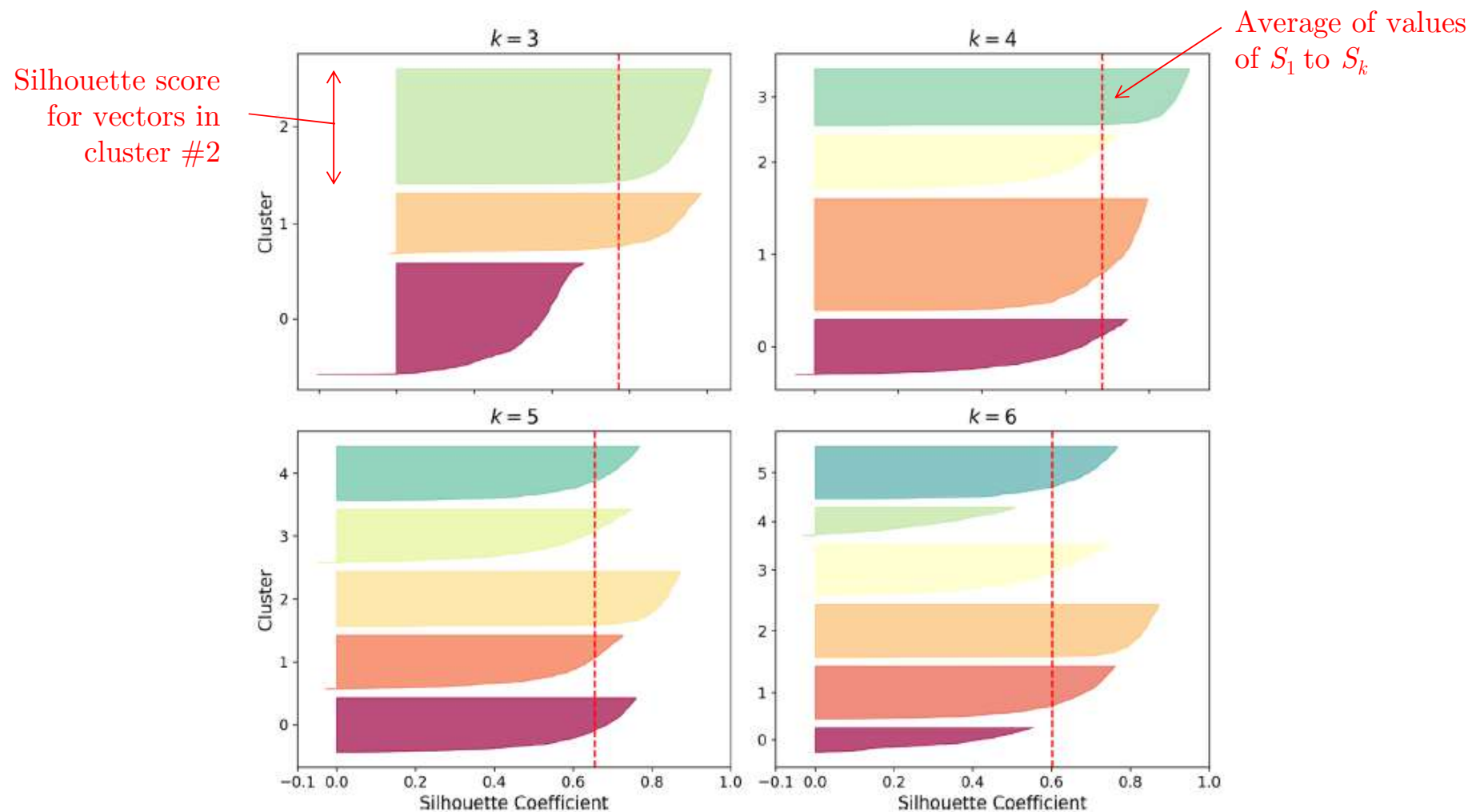
$$S_l = \sum_{\mathbf{x}_i \in \omega_l} (b_i - a_i) / \max(a_i, b_i) \quad \leftarrow \text{Silhouette score } \in (-1, 1) \text{ of vector } \mathbf{x}_i$$

5.     Average the values  $S_1$  to  $S_h$
6. The estimated number of clusters is the value in  $\{2, K\}$  that maximizes this average

For the example above, it is confirmed that the good number of clusters is 5:



A more informative visualization is obtained when plotting every vector *silhouette coefficient sorted by the cluster*:



When many vectors in a cluster have lower silhouette coefficient than the average, the clustering is bad: those vectors are too close to other clusters.

# CONTENTS

## **6.1 Parametric methods**

## **6.2 Clustering**

### 6.2.1 Introduction

### 6.2.2 Criterion functions

### 6.2.3 Iterative optimization

### 6.2.4 HDBSCAN

### 6.2.6 Conclusions

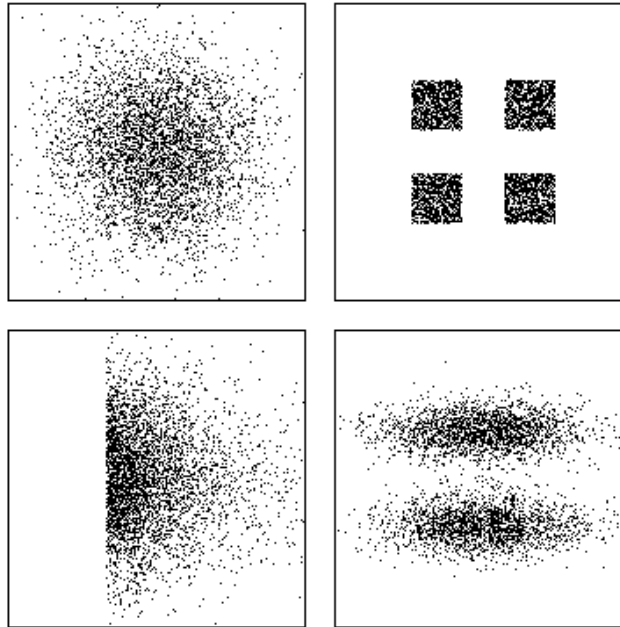


## 2.1 INTRODUCTION

- *Clusters* are clouds of points in a d-dimensional space.
- *Normal Distribution*, sufficient statistics are the
  - *Mean sample  $\mathbf{m}$* : Locates the center of gravity of the cloud and it best represents all of the data in the sense of minimizing the sum of squared distances from  $\mathbf{m}$  to the samples.
  - *Sample covariance matrix  $\mathbf{C}$* : denotes the amount the data scatters along various directions around  $\mathbf{m}$ .

Sample mean vector and sample covariance matrix are not sufficient statistics in a general case.

Distributions with identical mean and covariance:



Let us drop the Gaussian assumption for the clusters and adopt a non-parametric clustering approach:

Three key concepts in clustering procedures:

1. Data are grouped in clusters or groups of data points that possess strong internal **similarity**.
2. A **criterion function** is used to seek the grouping that extremizes it. To evaluate the partitioning of a set of samples into clusters, the similarity is measured between samples.
3. An **iterative algorithm** whereby vector is moved from cluster  $i$  to a cluster  $j$  only if the global criterion function decreases.

## SIMILARITY MEASURES

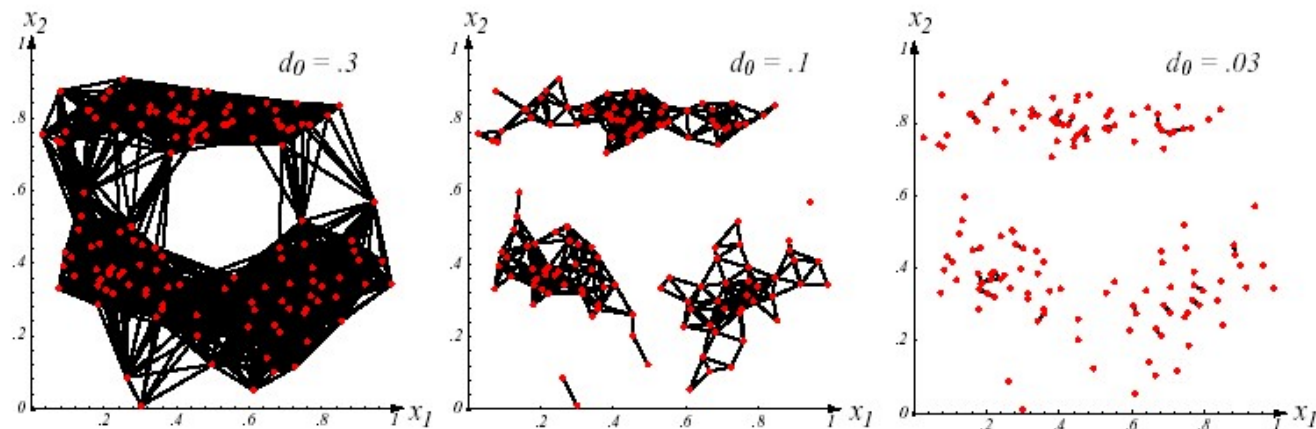
Two samples belong to the same cluster if  $d(\mathbf{x}_i, \mathbf{x}_j) < d_o$

The choice of threshold  $d_o$  is critical!

- Minkowski Distance  $d_q(\mathbf{x}_i, \mathbf{x}_j) = \left( \sum_{n=1}^d (x_i[n] - x_j[n])^q \right)^{\frac{1}{q}}$
- Mahalanobis Distance  $d_M^2(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{\Sigma}^{-1} (\mathbf{x}_i - \mathbf{x}_j)$

Advantages of Euclidean distance: clusters are invariant to rotation and translation.

Distance threshold affects the number and size of clusters:



typical intra-cluster distance  $< d_0 <$  typical inter-clusters distance

Data normalization prior to clustering:

- Each feature is translated to have zero mean
- Each feature is scaled to have unit variance
- PCA: axis coincide with the eigenvectors of the sample covariance matrix

**After normalización and PCA, clusters are invariant to displacements, scale change and rotations**

## 2.2 CRITERION FUNCTIONS

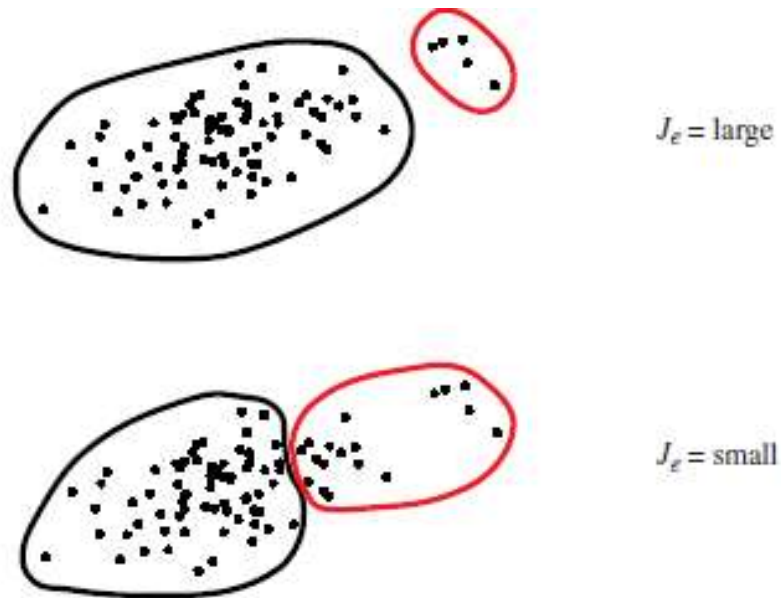
Criterion functions for clustering:

- Initial set  $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$
- Partition into exactly  $c$  subsets  $D_1, D_2, \dots, D_c$
- Objective: To find the partition that extremizes the criterion function, e.g. sum of squared errors criterion

$$J_e = \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} \|\mathbf{x} - \mathbf{m}_i\|^2$$

- $\mathbf{m}_i$  is the centroid, as the best representative of the samples in  $D_i$ .

It is appropriate when the clusters form compact clouds and uniform number of samples per cluster, otherwise it may divide natural clusters artificially:





## DISTANCE-BASED CRITERIA

Other criterion function:

$$J_e = \frac{1}{2} \sum_{i=1}^c n_i \bar{s}_i$$

$$\bar{s}_i = \frac{1}{n_i^2} \sum_{\mathbf{x} \in D_i} \sum_{\mathbf{x}' \in D_i} \|\mathbf{x} - \mathbf{x}'\|^2 \quad \text{Average distance in class } i$$

$$\bar{s}_i = \max_{\mathbf{x}, \mathbf{x}' \in D_i} d_e(\mathbf{x}, \mathbf{x}') \quad \text{Maximum distance in class } i$$

$$\bar{s}_i = \min_{\mathbf{x}, \mathbf{x}' \in D_i} d_e(\mathbf{x}, \mathbf{x}') \quad \text{Minimum distance in class } i$$

any similarity measure could be used

## SCATTER-BASED CRITERIA

Sensible criterion: maximize distance among clusters and minimize in-cluster scatter

**Scatter criteria:**

*Definitions used in clustering criteria...*

- Mean vector for the  $i$  cluster  $\mathbf{m}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in D_i} \mathbf{x}$
- Total mean vector  $\mathbf{m} = \frac{1}{n} \sum_{\mathbf{x} \in D} \mathbf{x} = \frac{1}{n} \sum_{i=1}^c n_i \mathbf{m}_i$
- Scatter matrix for the  $i$  cluster  $\mathbf{S}_i = \sum_{\mathbf{x} \in D_i} (\mathbf{x} - \mathbf{m}_i)(\mathbf{x} - \mathbf{m}_i)^T$
- Within-cluster scatter matrix  $\mathbf{S}_W = \sum_{i=1}^c \mathbf{S}_i$
- Between-cluster scatter matrix  $\mathbf{S}_B = \sum_{i=1}^c n_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T$
- Total scatter matrix  $\mathbf{S}_T = \sum_{\mathbf{x} \in D} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T = \mathbf{S}_W + \mathbf{S}_B$

## 1. Scatter criteria based on the trace...

- It measures the square of the scattering radius
- Minimize the trace of the Within-Cluster Scatter Matrix

$$Tr(\mathbf{S}_W) = \sum_{i=1}^c Tr(\mathbf{S}_i) = \sum_{i=1}^c \sum_{\mathbf{x} \in D_i} \|\mathbf{x} - \mathbf{m}_i\|^2 = J_e$$

**It results in function  $J_e$ !**

- It is equivalent to maximize between cluster scattering matrix trace.

$$Tr(\mathbf{S}_W) = Tr(\mathbf{S}_T) - Tr(\mathbf{S}_B) \qquad Tr(\mathbf{S}_B) = \sum_{i=1}^c n_i \|\mathbf{m}_i - \mathbf{m}\|^2$$

## 2. Scatter criteria based on the determinant...

- It measures the square of the scattering volume.
- $\mathbf{S}_B$  is singular if  $c \leq d$ ;  $\text{rank}(\mathbf{S}_B) \leq \min(d, c-1)$
- $\mathbf{S}_W$  is singular if  $n + c < d$
- Assuming  $n > d - c$

$$J_d = |\mathbf{S}_W| = \left| \sum_{i=1}^c \mathbf{S}_i \right|$$

- It does not change if the axes are scaled

### 3. Invariant scatter criteria...

Eigenvalues of  $\mathbf{S}_W^{-1}\mathbf{S}_B$  are **invariant to nonsingular linear transformations** of the data (see chapter 3). We can invent criteria based on the eigenvalues:

$$Tr(\mathbf{S}_W^{-1}\mathbf{S}_B) = \sum_{i=1}^d \lambda_i \quad \frac{|\mathbf{S}_W|}{|\mathbf{S}_T|} = \prod_{i=1}^d \frac{1}{1 + \lambda_i}$$
$$J_f = trace(\mathbf{S}_T^{-1}\mathbf{S}_W) = \sum_{i=1}^d \frac{1}{1 + \lambda_i}$$

← It is the sum of squared errors after normalization with matrix  $\mathbf{S}_T$

relying on the fact that good partitions are those for which the non-zero eigenvalues are large.

- Other proposed criteria (equivalent for  $c = 2$ ):
$$\min Tr(\mathbf{S}_T^{-1}\mathbf{S}_W)$$
$$\max Tr(\mathbf{S}_W^{-1}\mathbf{S}_B)$$

Property...  $\lambda_1, \dots, \lambda_i, \dots, \lambda_d = \text{eigenvalues}(\mathbf{S}_W^{-1} \mathbf{S}_B) \Rightarrow$   
 $\frac{1}{1+\lambda_1}, \dots, \frac{1}{1+\lambda_i}, \dots, \frac{1}{1+\lambda_d} = \text{eigenvalues}(\mathbf{S}_T^{-1} \mathbf{S}_W)$

- Proof:

$$\mathbf{S}_B \mathbf{v}_i = \lambda_i \mathbf{S}_W \mathbf{v}_i \Rightarrow$$

$$\mathbf{S}_T \mathbf{v}_i = \mathbf{S}_B \mathbf{v}_i + \mathbf{S}_W \mathbf{v}_i = \lambda_i \mathbf{S}_W \mathbf{v}_i + \mathbf{S}_W \mathbf{v}_i = (1 + \lambda_i) \mathbf{S}_W \mathbf{v}_i \Rightarrow$$

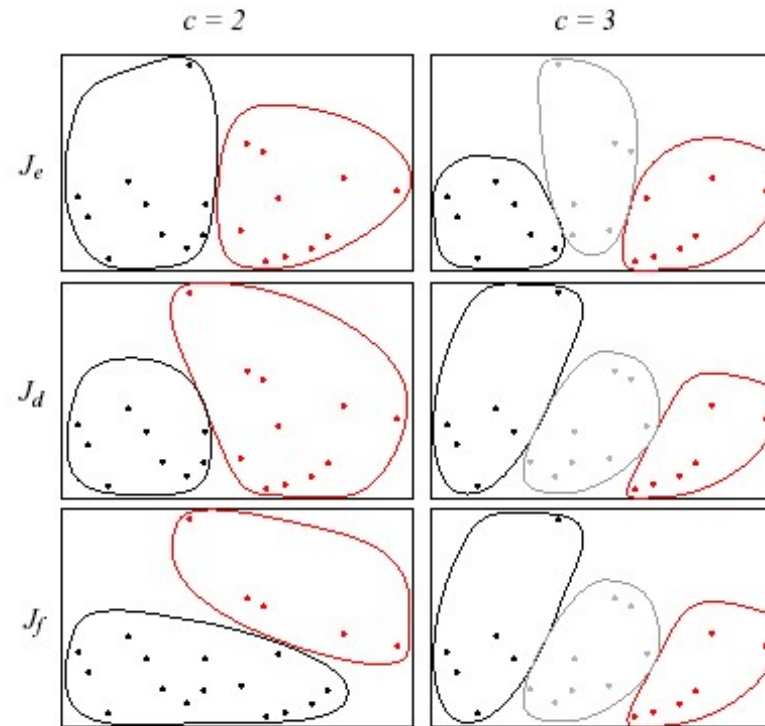
$$\mathbf{v}_i = (1 + \lambda_i) \mathbf{S}_T^{-1} \mathbf{S}_W \mathbf{v}_i \Rightarrow$$

$$\mathbf{S}_T^{-1} \mathbf{S}_W \mathbf{v}_i = \frac{1}{1+\lambda_i} \mathbf{v}_i$$

### Example:

- Trace criteria
- Determinant criteria
- Invariant criteria

It turns out that all criteria are equivalent as  $c$  increases



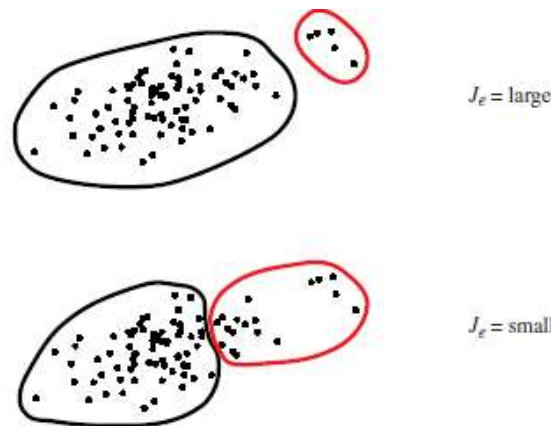
## Conclusions on clustering procedures...

- Underlying model: assumes that samples form  $c$  fairly well separated clouds of points.
- $S_w$  measures the compactness of these clouds, so it may not work in certain situations:
  - If there is a dense cluster inside a diffuse cluster
  - Separate intertwined line-like clusters
- **Problem:** Computational complexity to evaluate the overall number of possibilities in partitioning is impractical. We need some iterative procedure...



## 2.3 ITERATIVE OPTIMIZATION

- We adopt the euclidean distance criterion in the sequel.
- Direct partitioning by full enumeration of clusters has large complexity:  $c^n/c!$ , which is completely infeasible
- Practical solution: Initiate with some reasonable partition and move samples from one group to another if such move improves the value of the criterion function.
- It guarantees local but not global optimization...





- Iterative improvement to minimize the sum of squared error criterion  $J_e$ .
- Effective error per cluster  $J_i$ .

$$J_e = \sum_{i=1}^c J_i; \quad J_i = \sum_{\mathbf{x} \in D_i} \|\mathbf{x} - \mathbf{m}_i\|^2$$

- A sample is moved from cluster  $i$  to cluster  $j$ .

$$\hat{\mathbf{x}} \in D_i \Rightarrow \hat{\mathbf{x}} \in D_j$$

$$\mathbf{m}_j^* = \mathbf{m}_j + \frac{\hat{\mathbf{x}} - \mathbf{m}_j}{n_j + 1}; \quad \mathbf{m}_i^* = \mathbf{m}_i - \frac{\hat{\mathbf{x}} - \mathbf{m}_i}{n_i - 1}$$

$$n_j = n_j + 1; \quad n_i = n_i - 1$$



- Increasing / Decreasing “Effective error per cluster” (**Prove it as exercise**)

Contribution to error criterion due to cluster  $j$

$$\begin{aligned} J_j^* &= \sum_{\mathbf{x} \in D_j} \|\mathbf{x} - \mathbf{m}_j^*\|^2 + \|\hat{\mathbf{x}} - \mathbf{m}_j^*\|^2 = \\ &= \left( \sum_{\mathbf{x} \in D_j} \left\| \mathbf{x} - \mathbf{m}_j - \frac{\hat{\mathbf{x}} - \mathbf{m}_j}{n_j + 1} \right\|^2 \right) + \left\| \frac{n_j}{n_j + 1} (\hat{\mathbf{x}} - \mathbf{m}_j) \right\|^2 = J_j + \frac{n_j}{n_j + 1} \|\hat{\mathbf{x}} - \mathbf{m}_j\|^2 \end{aligned}$$

Contribution to error criterion due to cluster  $i$

$$\begin{aligned} J_i^* &= \sum_{\mathbf{x} \in D_i} \|\mathbf{x} - \mathbf{m}_i^*\|^2 - \|\hat{\mathbf{x}} - \mathbf{m}_i^*\|^2 = \\ &= \left( \sum_{\mathbf{x} \in D_i} \left\| \mathbf{x} - \mathbf{m}_i + \frac{\hat{\mathbf{x}} - \mathbf{m}_i}{n_i - 1} \right\|^2 \right) - \left\| \frac{n_i}{n_i - 1} (\hat{\mathbf{x}} - \mathbf{m}_i) \right\|^2 = J_i - \frac{n_i}{n_i - 1} \|\hat{\mathbf{x}} - \mathbf{m}_i\|^2 \end{aligned}$$

- The sample move from cluster  $i$  to cluster  $j$  is advantageous if

$$\frac{n_i}{n_i-1} \left\| \hat{\mathbf{x}} - \mathbf{m}_i \right\|^2 > \frac{n_j}{n_j+1} \left\| \hat{\mathbf{x}} - \mathbf{m}_j \right\|^2$$

## BASIC ITERATIVE MSE CLUSTERING

```

1 begin initialize  $n, c, \mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_c$ 
2   do randomly select a sample  $\hat{\mathbf{x}}$ ;
3      $i \leftarrow \arg \min_{i'} \|\mathbf{m}_{i'} - \hat{\mathbf{x}}\|$  (classify  $\hat{\mathbf{x}}$ )
      update  $\mathbf{m}_i$ 
4     if  $n_i \neq 1$  then compute
5       
$$\rho_j = \begin{cases} \frac{n_j}{n_j+1} \|\hat{\mathbf{x}} - \mathbf{m}_j\|^2 & j \neq i \\ \frac{n_j}{n_j-1} \|\hat{\mathbf{x}} - \mathbf{m}_i\|^2 & j = i \end{cases}$$

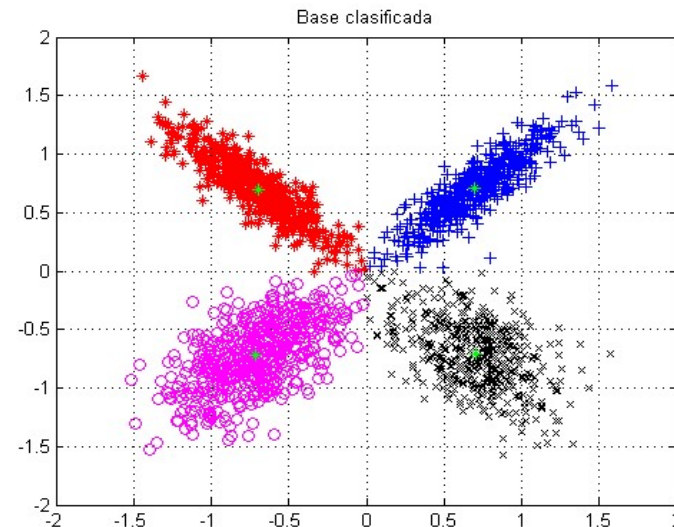
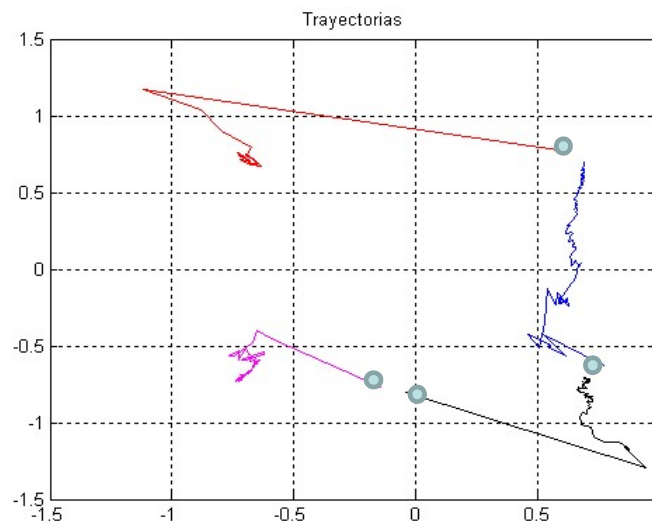
6       if  $\rho_k \leq \rho_j$  for all  $j$  then transfer  $\hat{\mathbf{x}}$  to  $\mathcal{D}_k$ 
7       recompute  $J_e, \mathbf{m}_i, \mathbf{m}_k$ 
8     until no change in  $J_e$  in  $n$  attempts
9   return  $\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_c$ 
10 end

```

If the number of samples in the cluster  $i$  is  $>1$

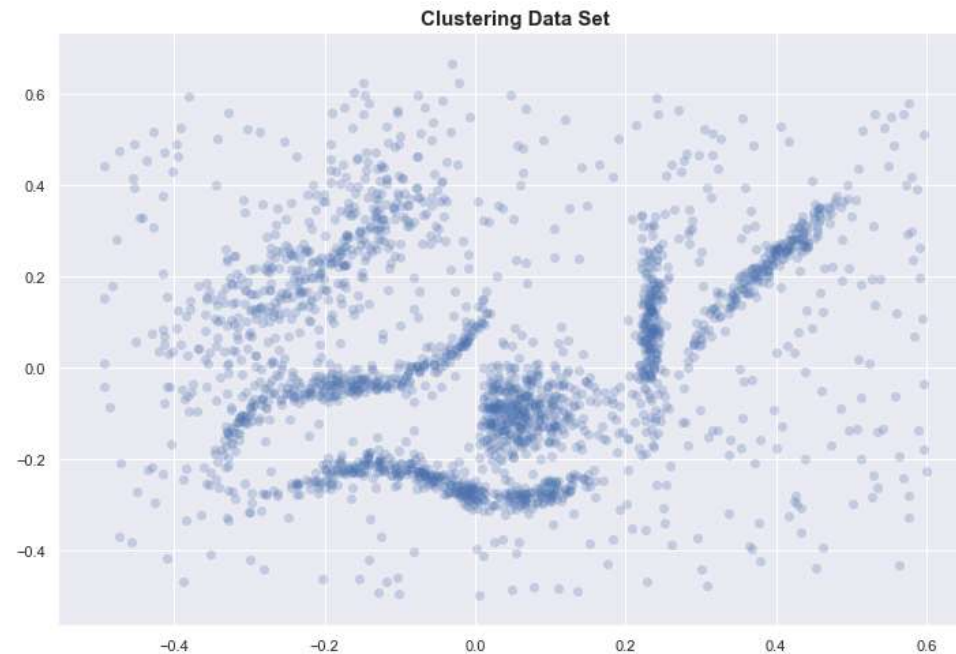
## Features of the algorithm...

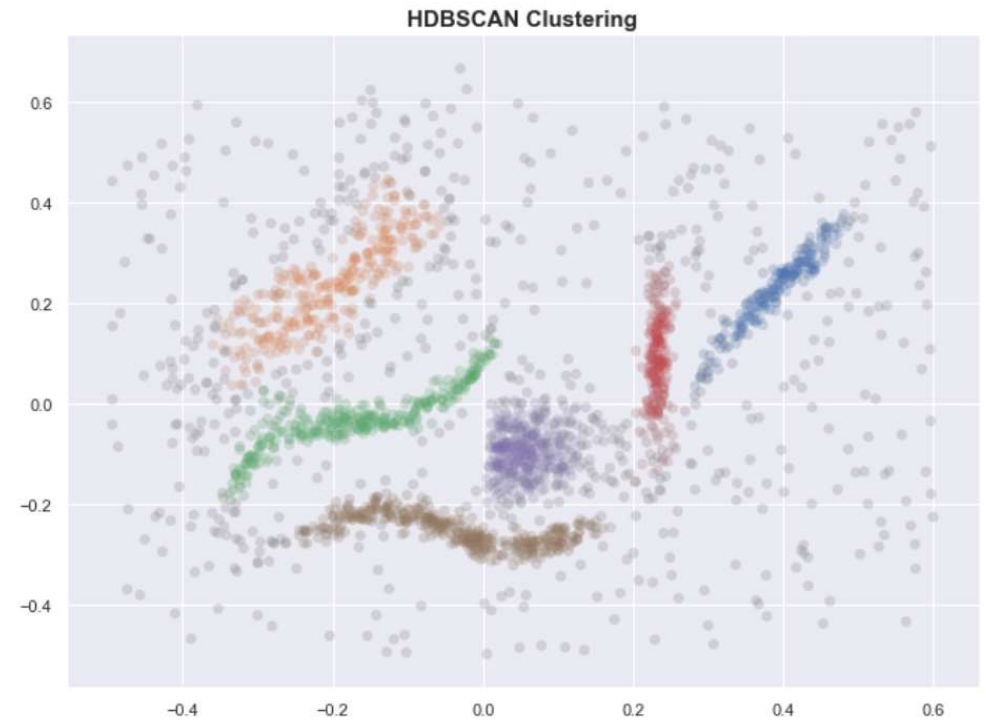
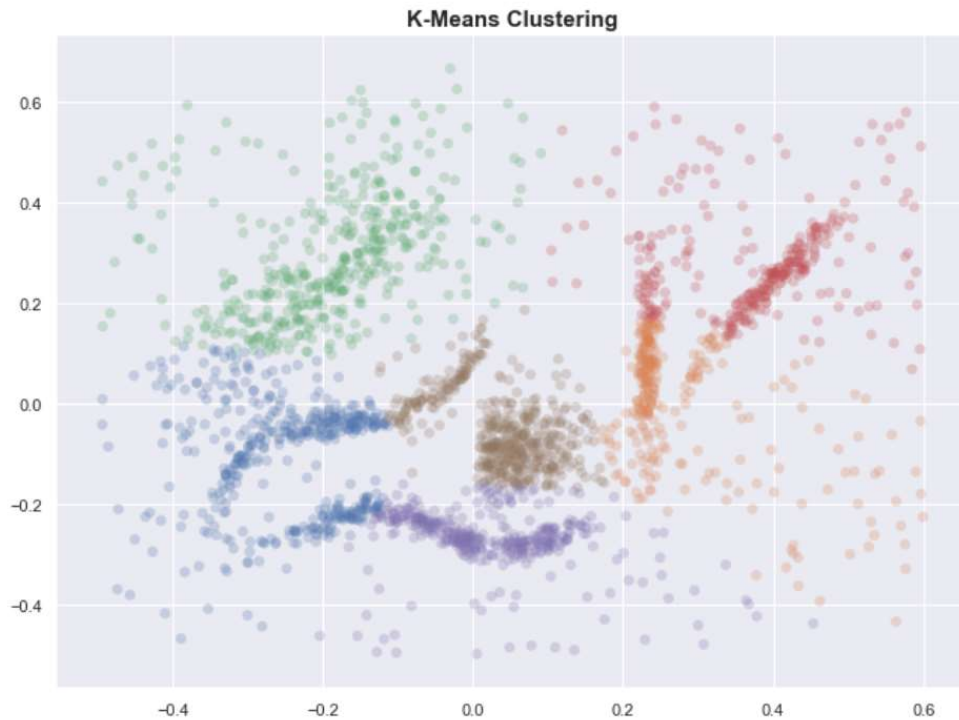
- It is a sequential version of the  $k$ -means procedure ( $k$ -means waits until all  $n$  samples are re-classified before updating means)
- This method is more susceptible to being trapped in local minima.
- Initial starting points can be obtained from the centroids obtained when seeking for  $c - 1$  clusters
- Suitable for problems where samples are acquired on-line.



## 2.4 DBSCAN

The *Density-Based Spatial Clustering of Applications with Noise*<sup>1</sup> is a clustering algorithm that defines clusters as continuous regions of high density. An example of a data base where k-means fails:





K-means fail because clusters are not spherical, not equally sized, not equally dense and contaminated by noise.



# DBSCAN

The principles of the algorithm:

1. For each vector  $\mathbf{x}_i$ , count how many samples are within a small distance  $\varepsilon$ .
2. If a vector has at least  $\text{min\_samp}$  other vectors in its  $\varepsilon$ -neighborhood, it is considered a *core instance*.
3. All vectors in the neighborhood belong to the same cluster. This include other possible *core instances*.
4. Any vector that is not a core instance and do not have one in its neighborhood is considered an anomaly.

# Further reading

The DBSCAN algorithm

Campello, Ricardo JGB, Davoud Moulavi, and Jörg Sander. “Density-based clustering based on hierarchical density estimates.” *Pacific-Asia conference on knowledge discovery and data mining*. Springer, Berlin, Heidelberg, 2013.

<https://towardsdatascience.com/understanding-hdbscan-and-density-based-clustering-121dbec1320e>

The hierarchical DBSCAN

John Healy. [HDBSCAN, Fast Density Based Clustering, the How and the Why](#). PyData NYC. 2018

## 2.5 CONCLUSIONS

- When underlying distribution comes from a mixture of component densities described by a set of unknown parameters, these parameters can be estimated by Bayesian or ML (EM algorithm) methods.
- Clustering is more general approach that does not rely on models. K-means and DSCAN are popular approaches.
- Selecting the number of clusters is crucial to have good performance.

