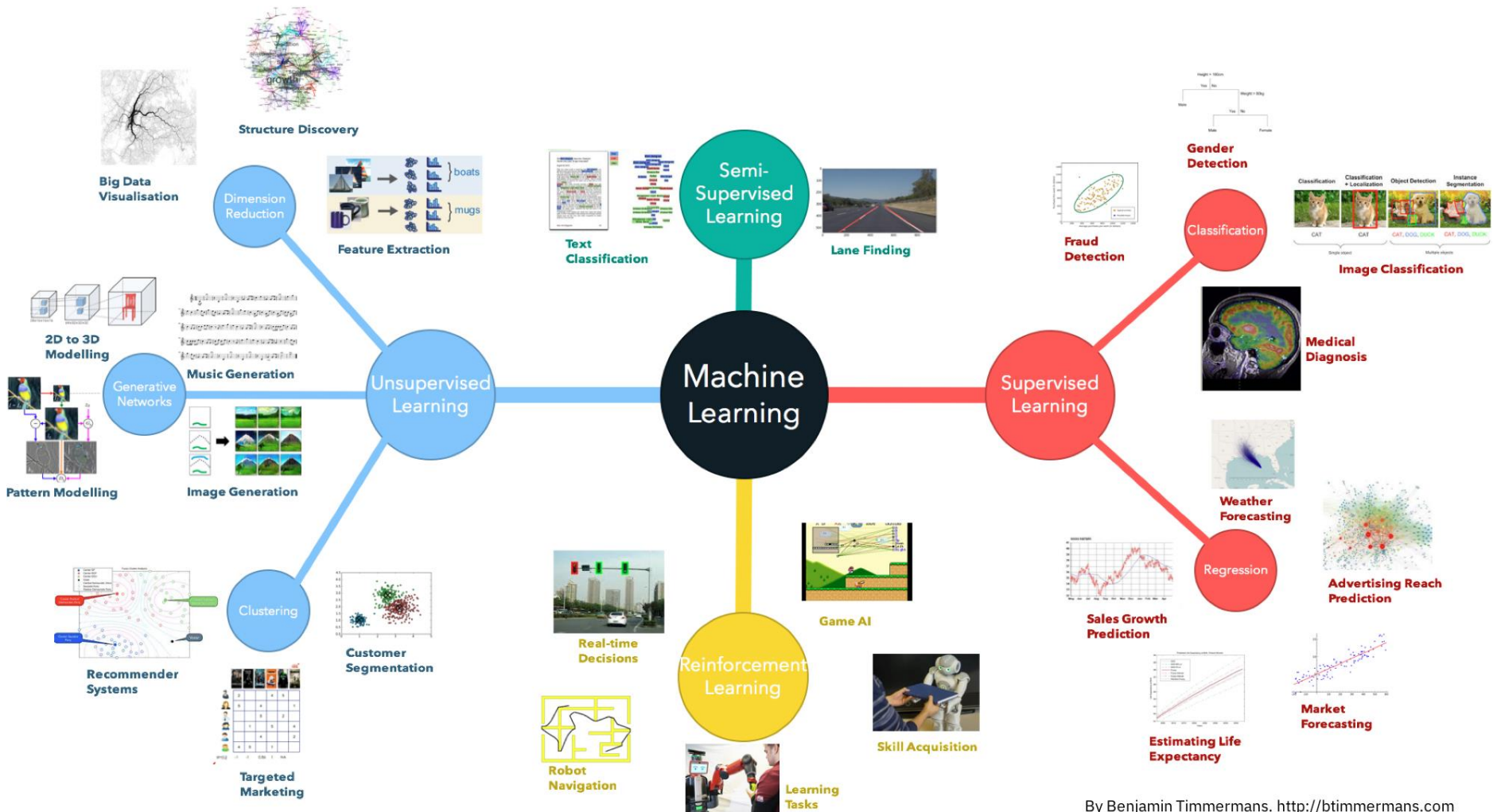


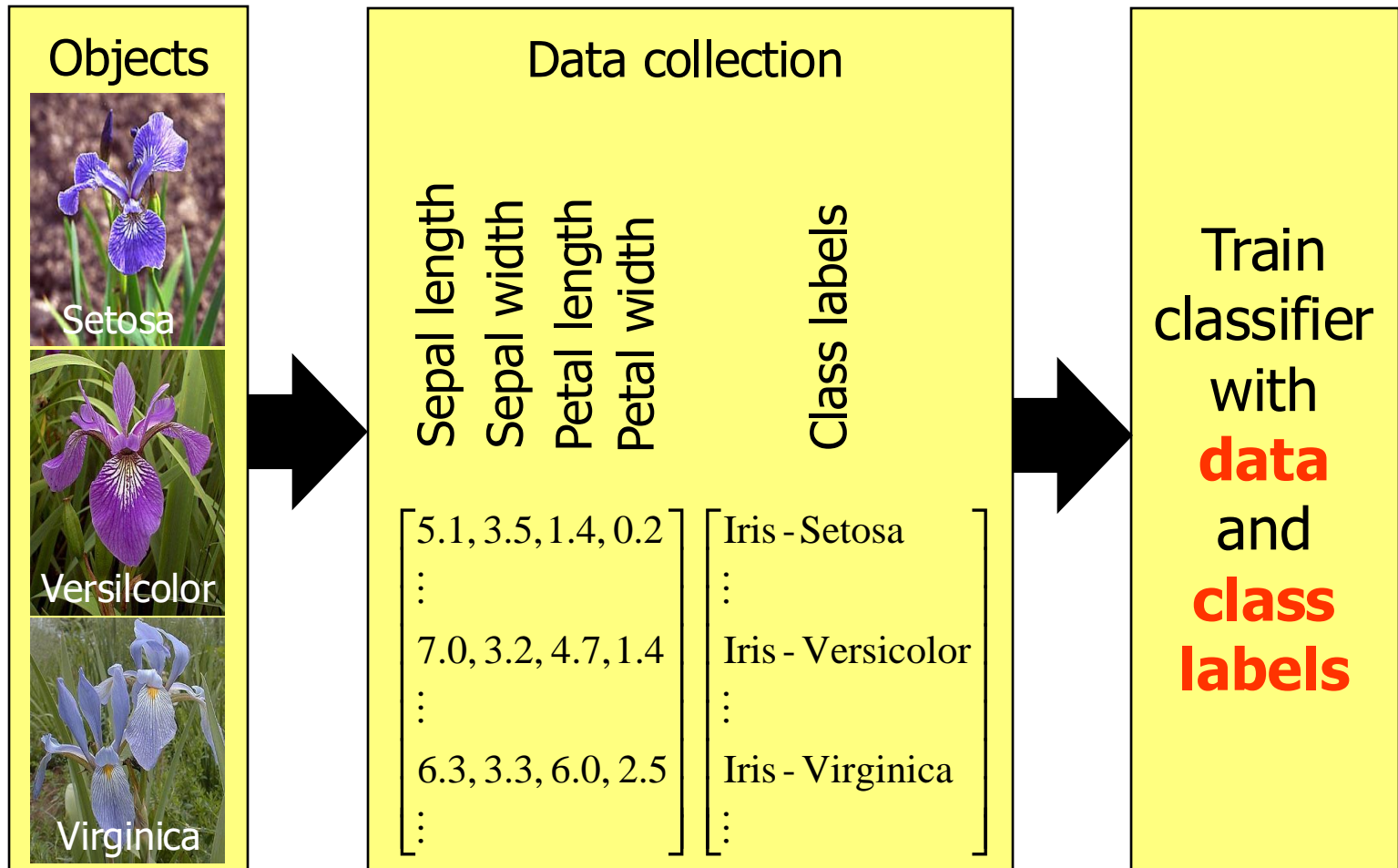
Clustering

DSAIT4020 Elements of Statistical Learning

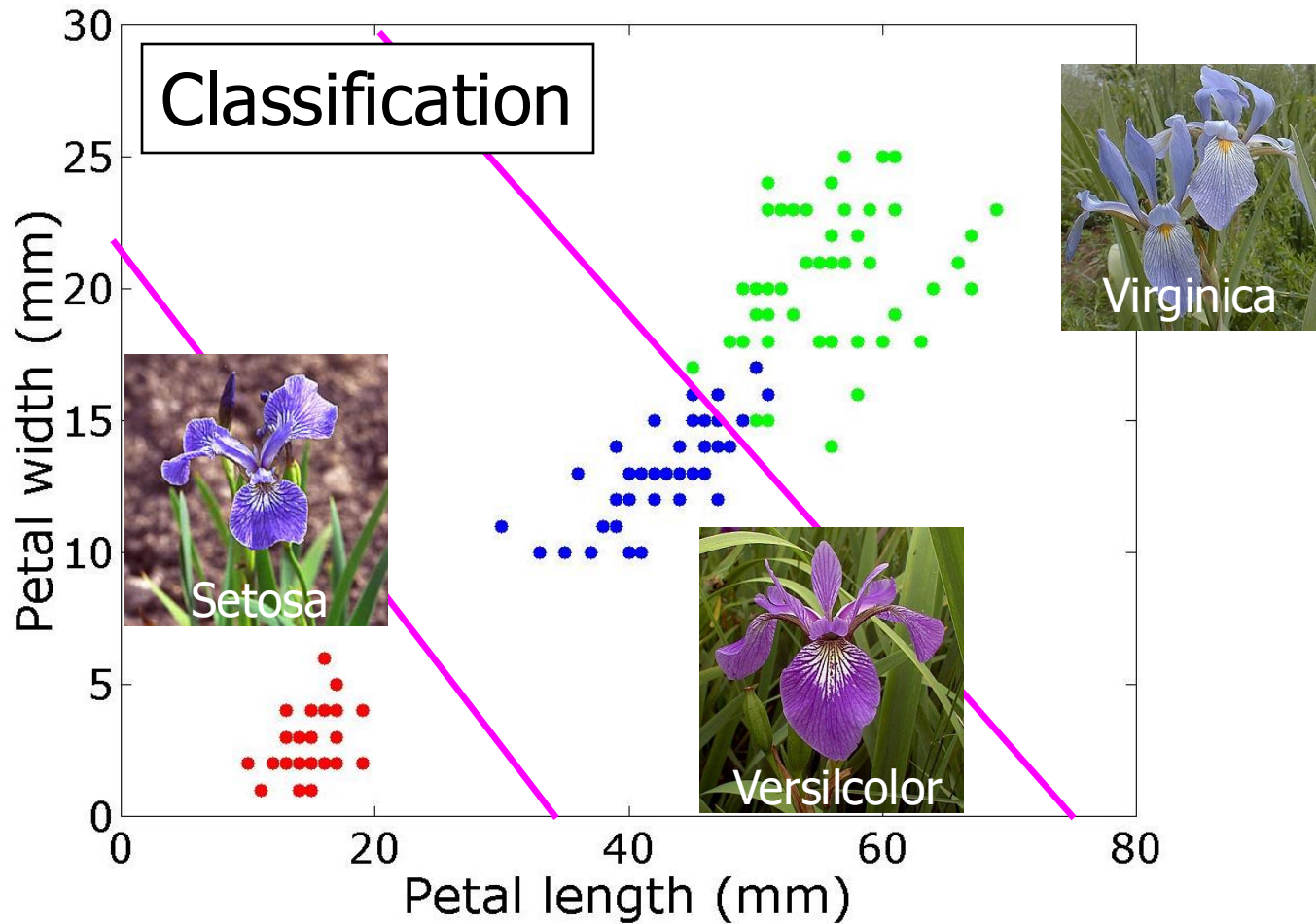


By Benjamin Timmermans. <http://btimmermans.com>

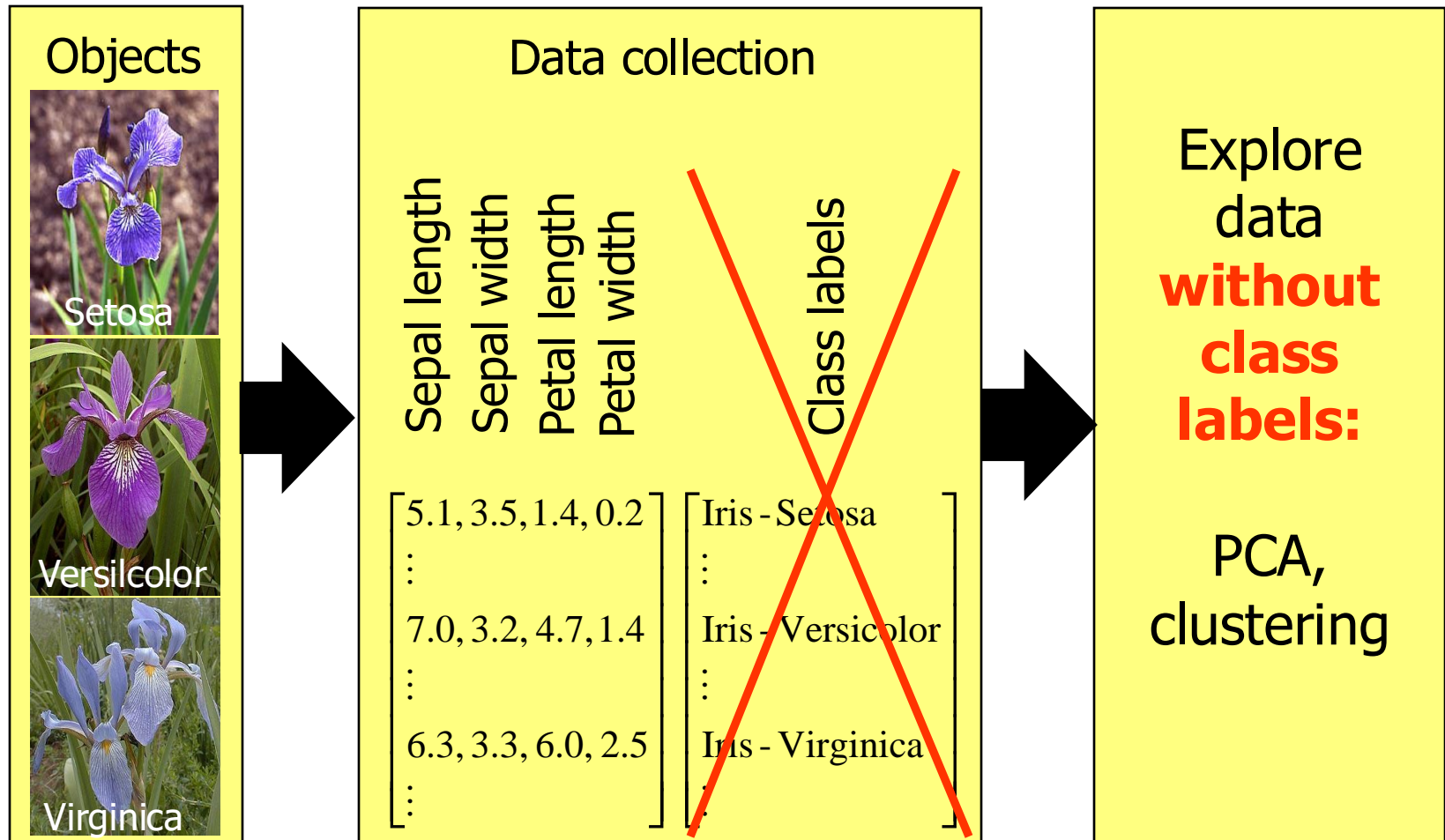
Supervised learning



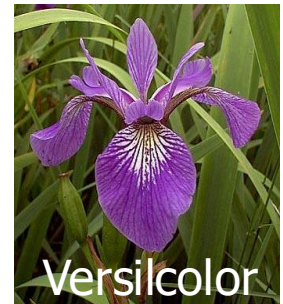
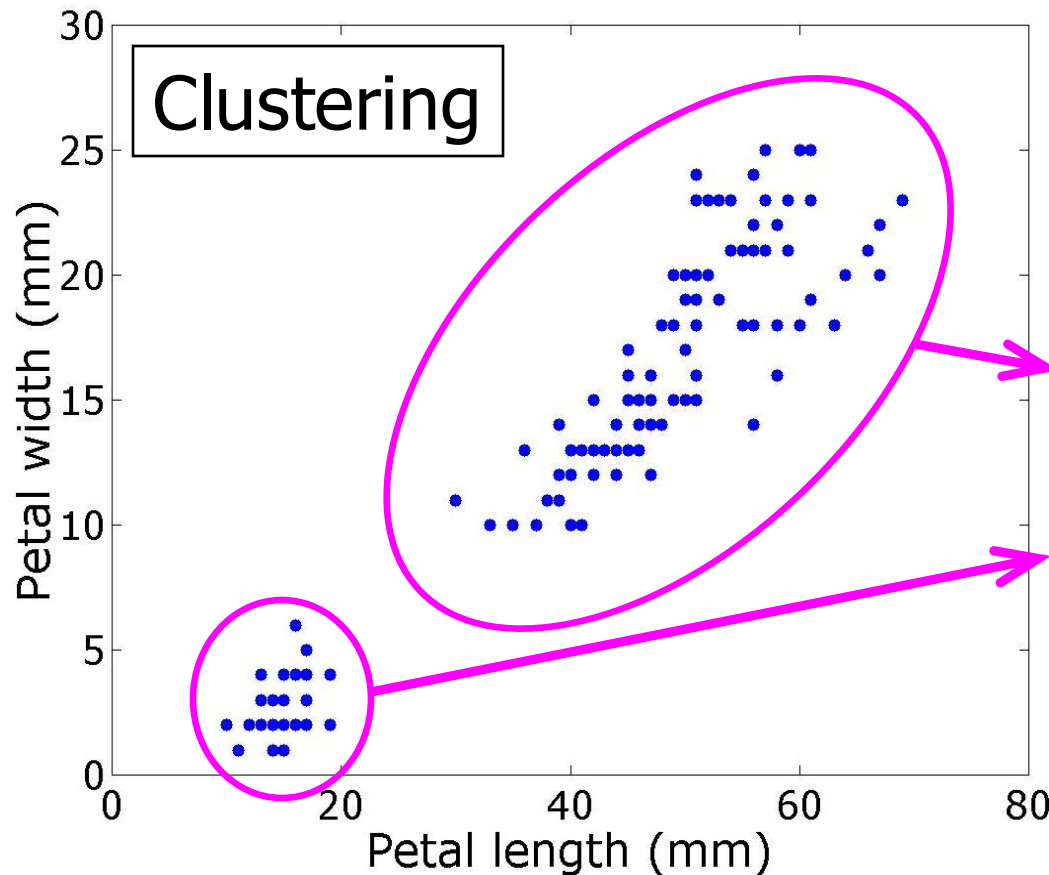
Supervised learning (2)



Unsupervised learning



Unsupervised learning (2)

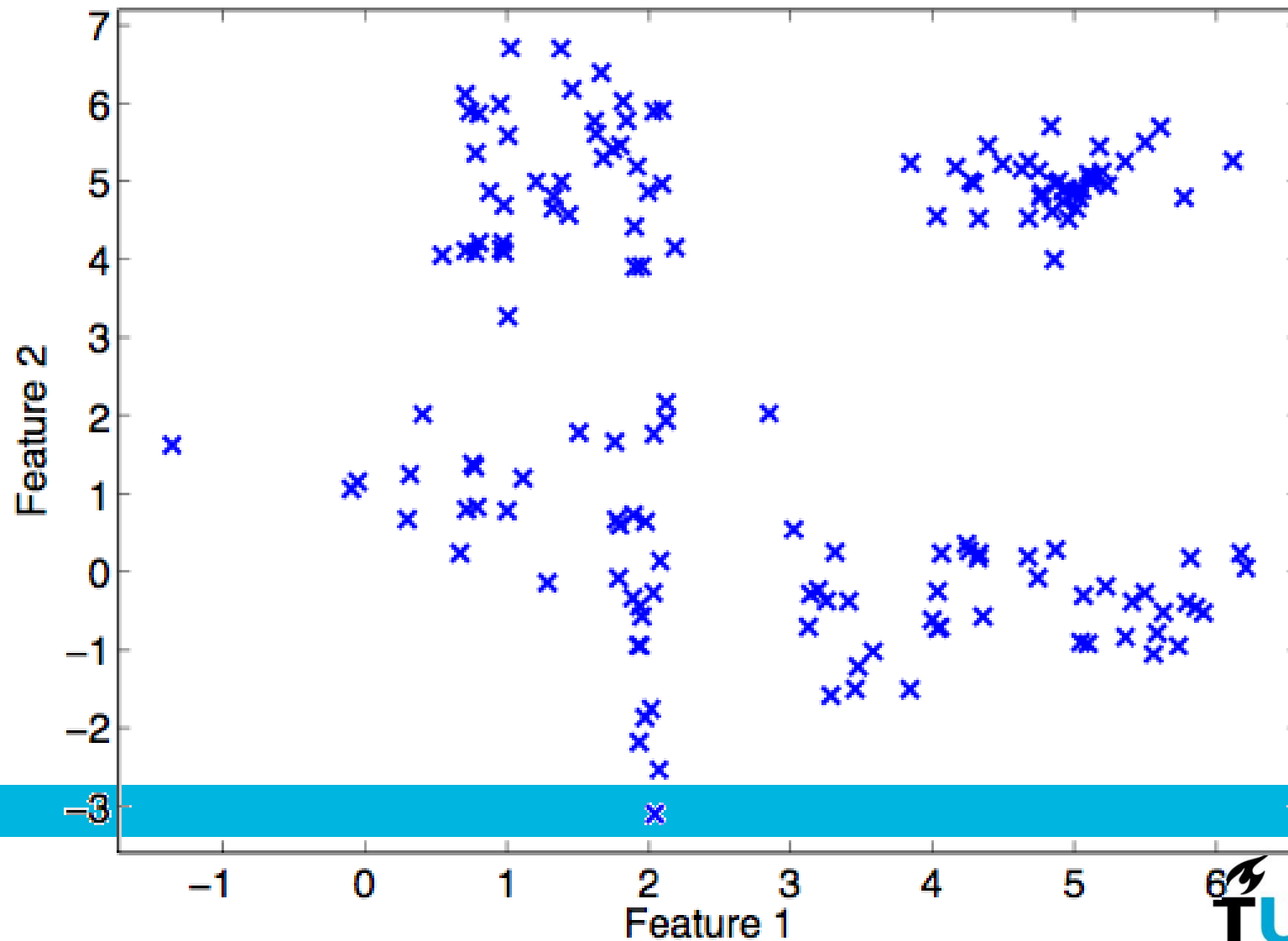


- What salient structures exist in the data?

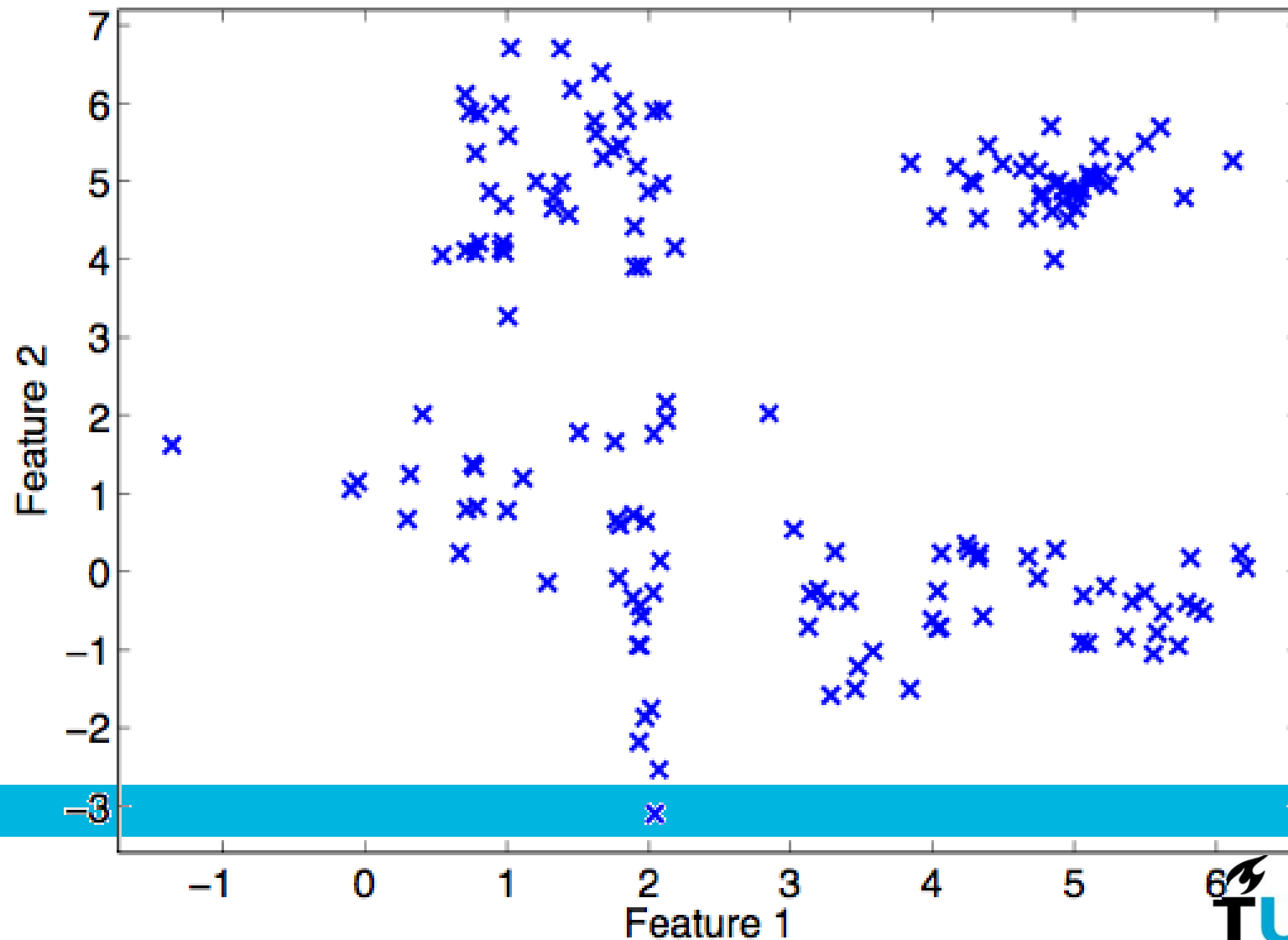
Cluster Analysis

- Grouping observations based on [dis]similarity
- E.g. data mining [exploration, searching for concepts in data]
 - Clustering species based on [genetic] similarity
 - Reducing amount of data to be analysed, helps defining concept / class
- Data reduction: selecting typical class examples
 - Multi-modal classes may be represented using typical examples
- Predicting characteristics for new data

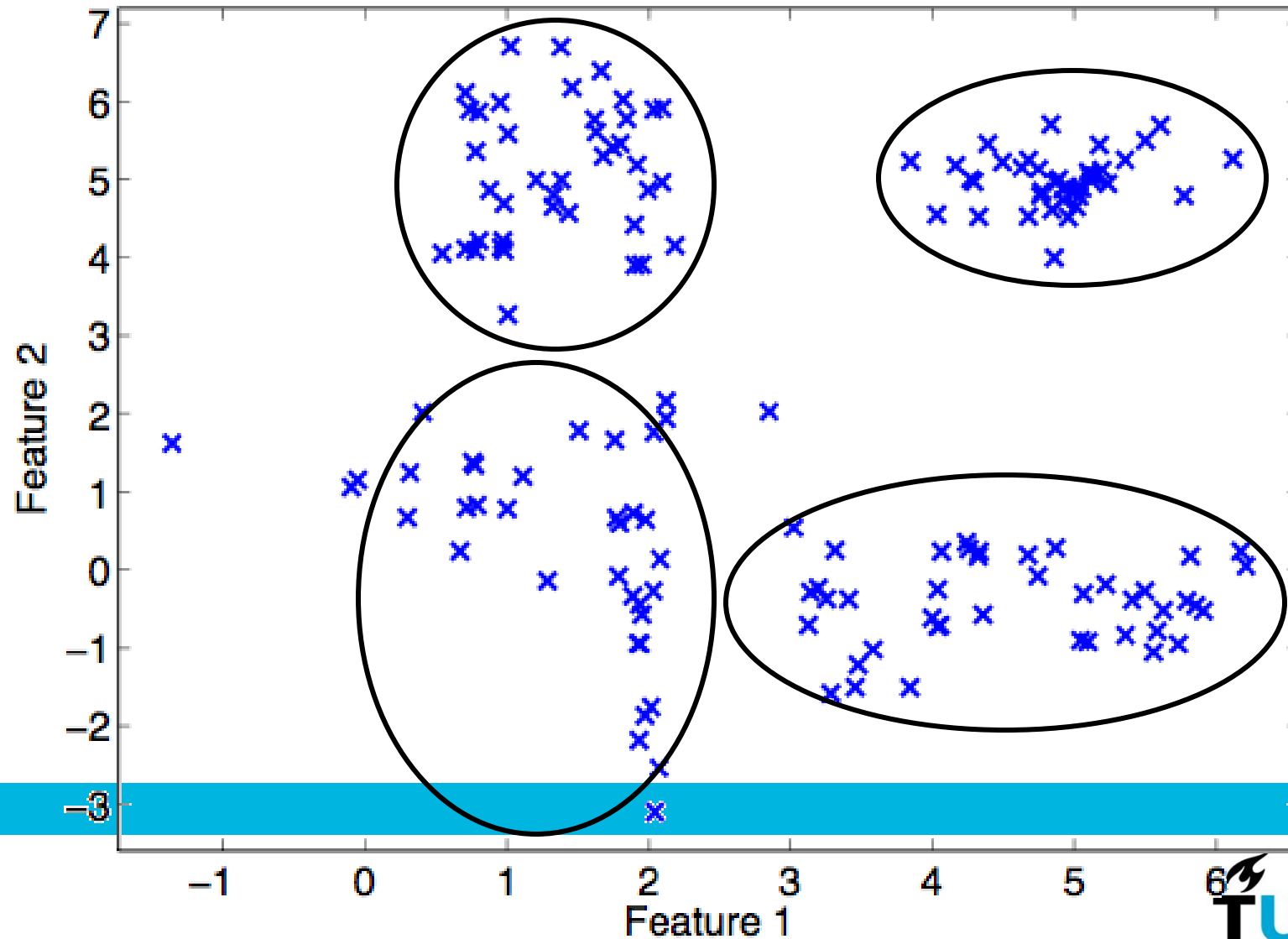
Unlabeled data: what now?



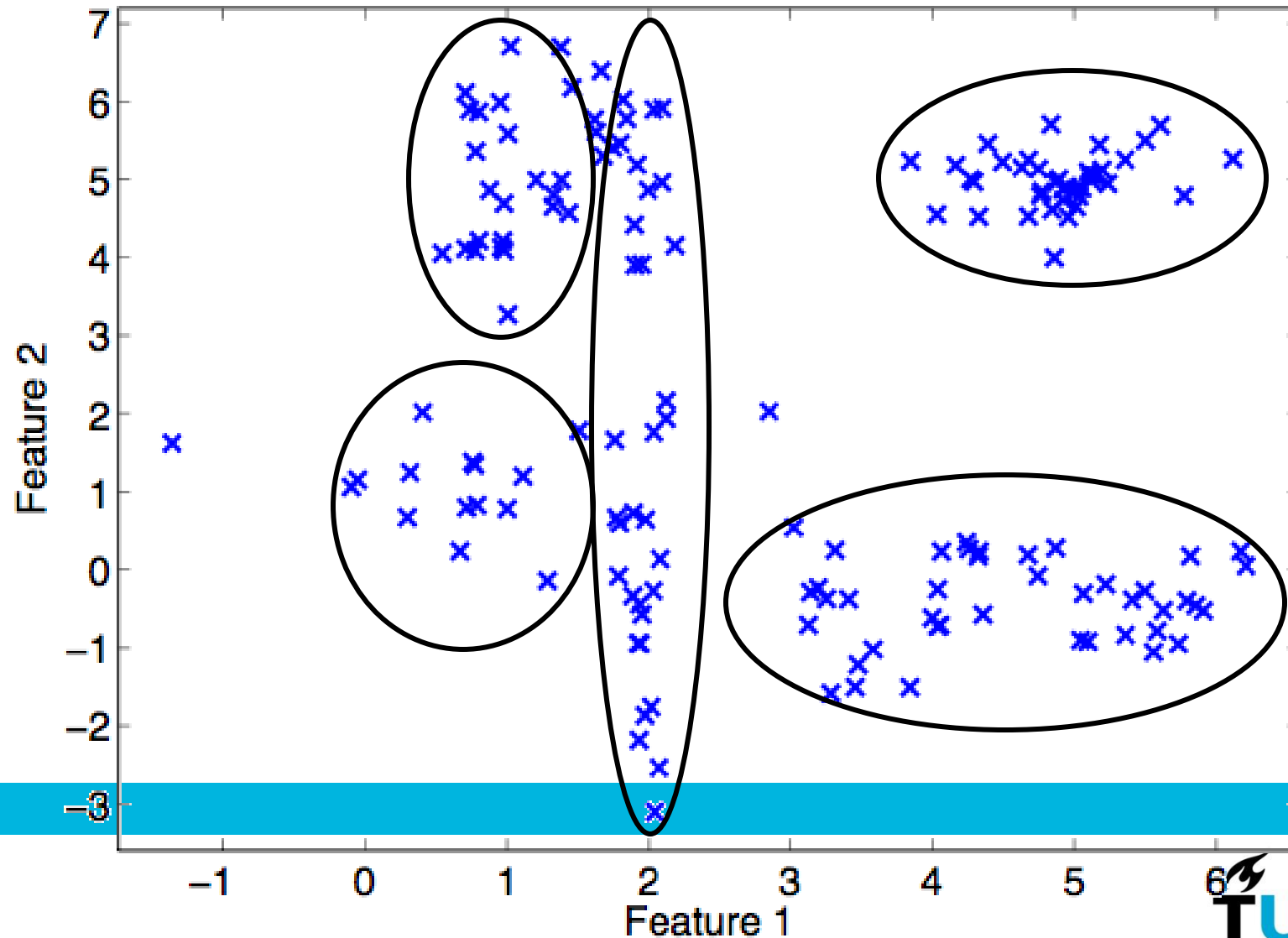
How Many Groups in Data?



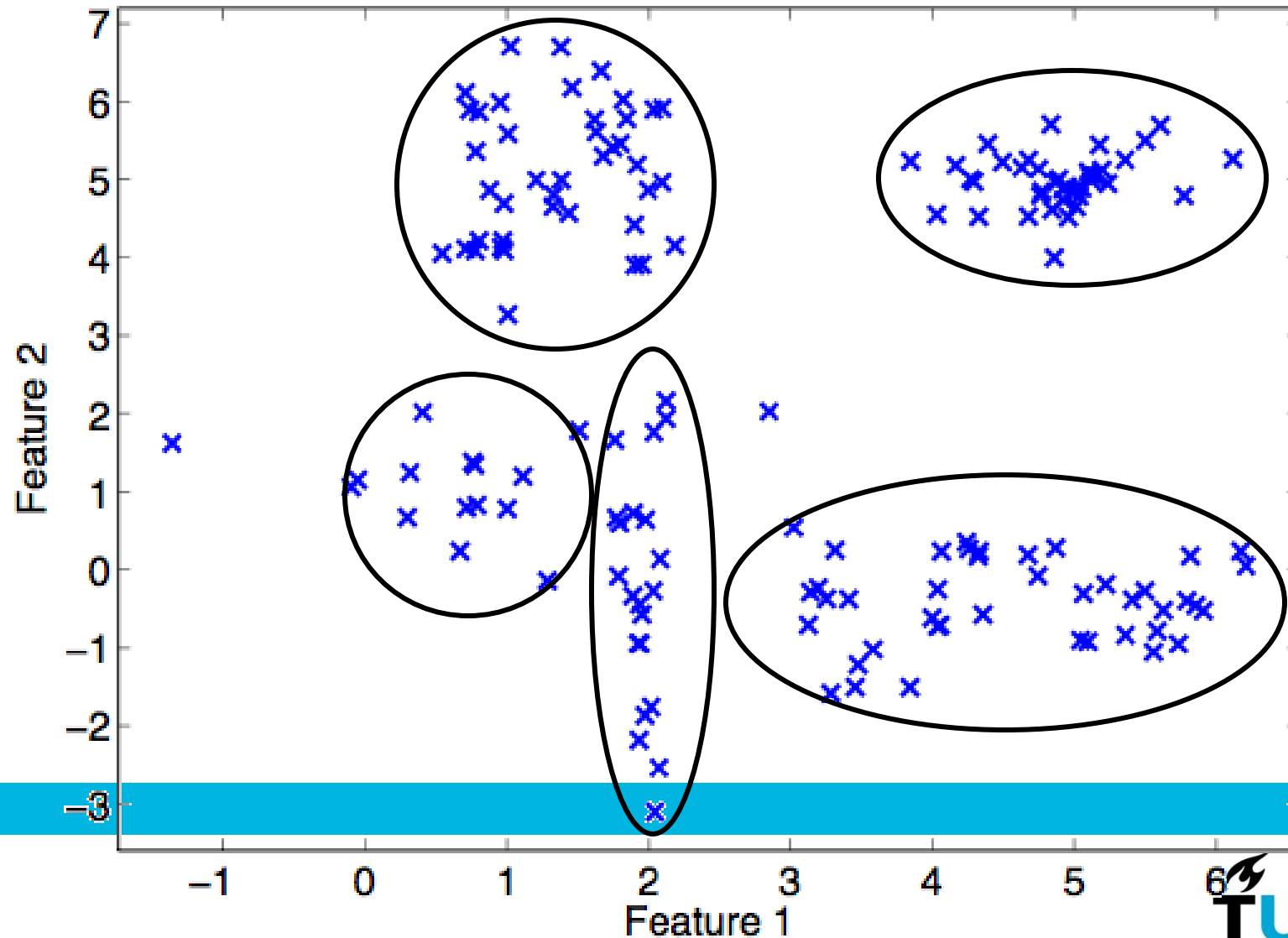
How Many Groups in Data?



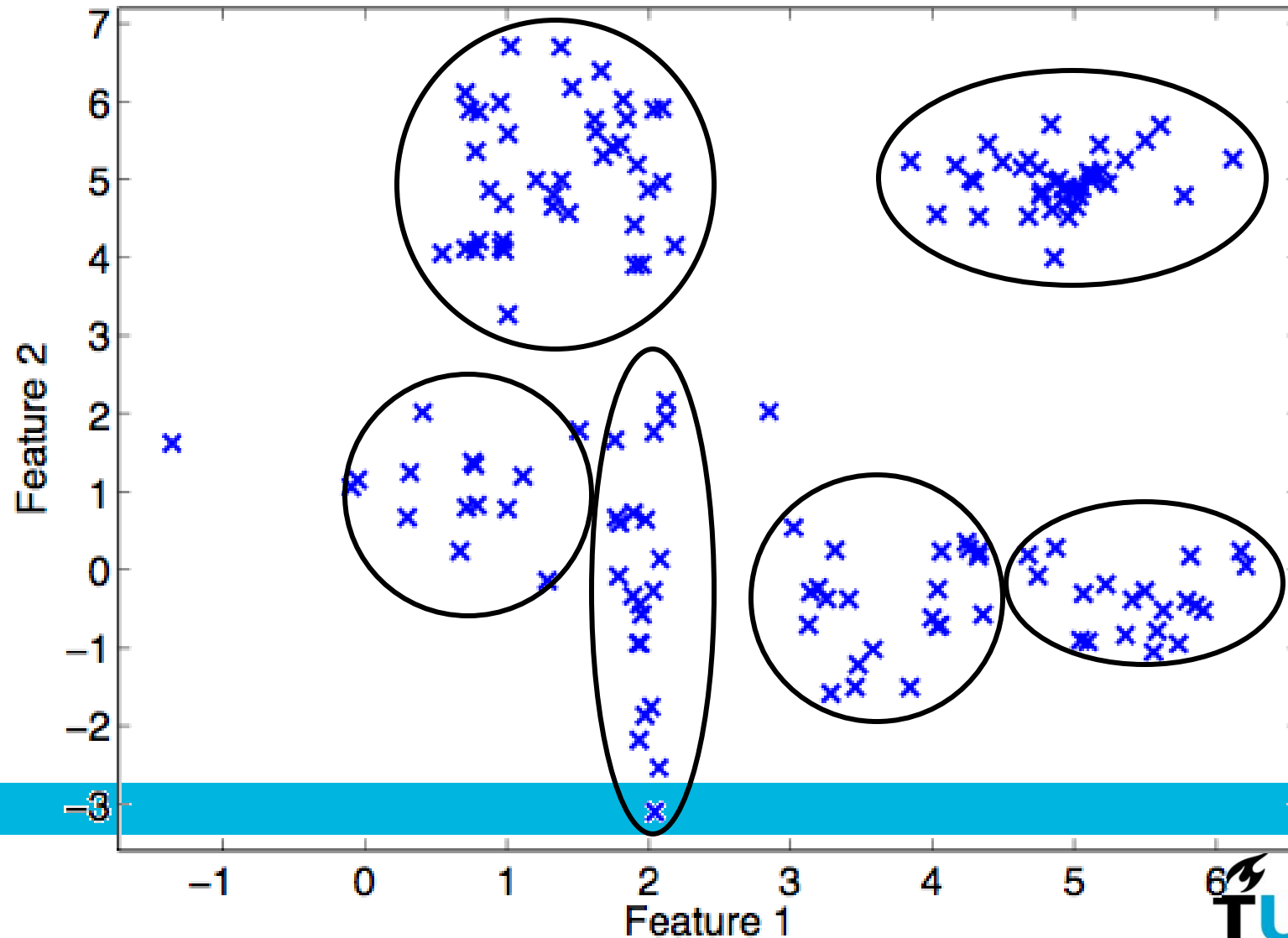
How Many Groups in Data?



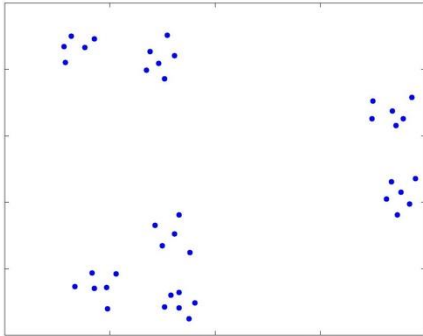
How Many Groups in Data?



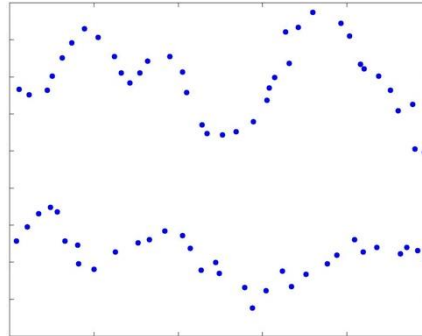
How Many Groups in Data?



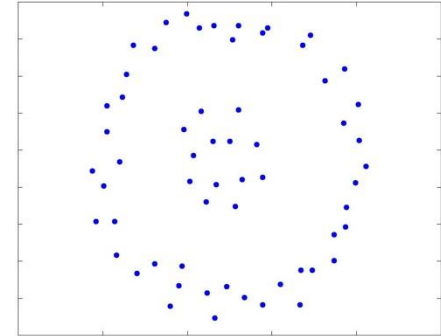
What is a cluster?



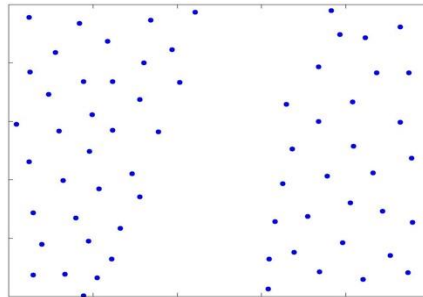
Shape: compact, convex
Separation: large



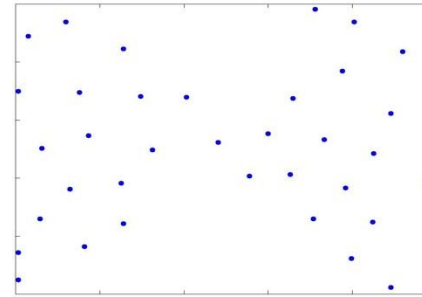
Shape: strings
Separation: large?



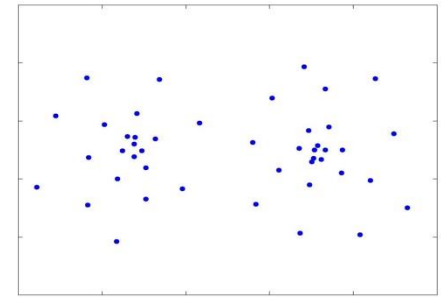
Shape: convex and circular
Separation: large?



Shape: ?
Separation: large?



Shape: loose, convex
Separation: small



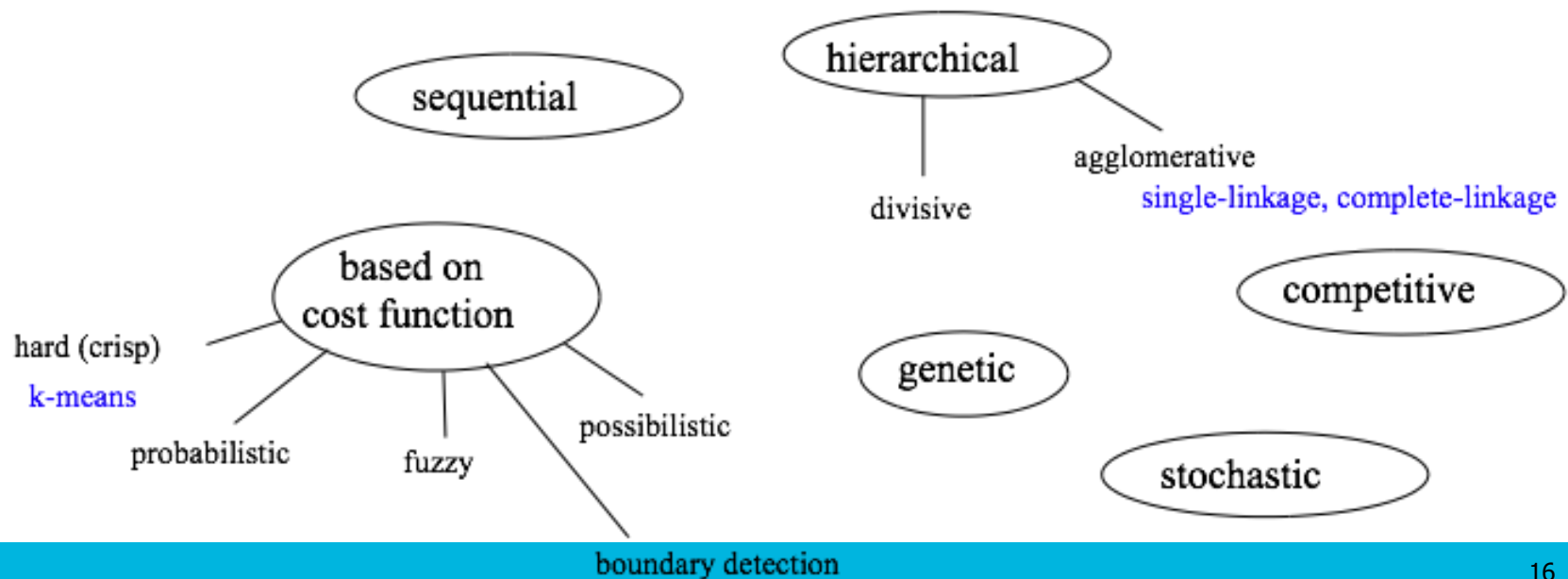
Shape: loose, convex
Separation: small

What is a cluster?

- Clustering: finding natural groups in data...
 - which themselves are far apart
 - in which objects are close together
- Define what is “far apart” and “close together”:
 - Need a distance measure or dissimilarity measure
 - This measure should capture what we think is important for the grouping
 - The choice for a certain distance measure is crucial!
- **There is no such thing as *the objective clustering***

Clustering Clustering Methods

- Very large field, huge number of methods
 - See for example **Theodoridis and Koutroumbas, Pattern Recognition, 2003**
 - More than 240 pages overview of cluster analysis



Chapters 11-16 from the book...

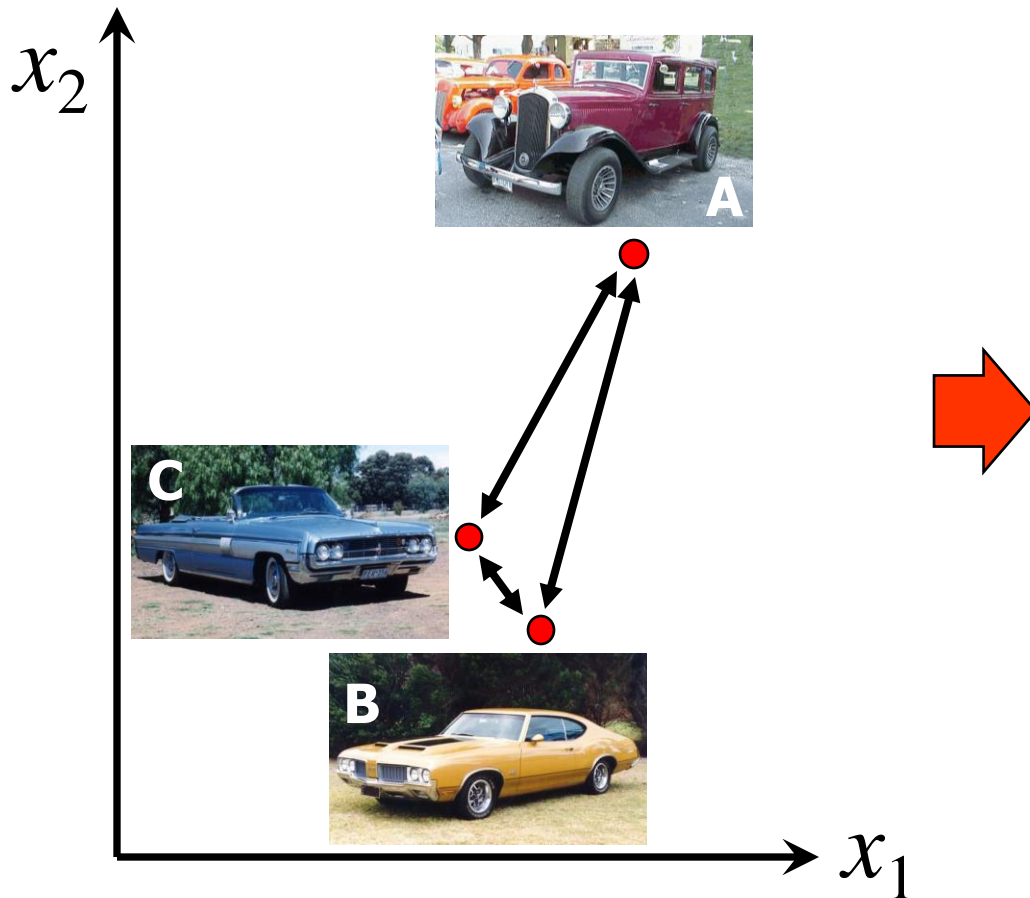
- In literature, an almost infinite number of methods is proposed
- The book tries to cover many of them
- We will discuss the most intuitive, and most used, clustering methods
- Ignore sections 12.3, 12.4, 12.5, 12.6, 12.7
- Ignore pages 661-692
- Ignore 14.3, 14.4, 14.6
- Ignore 15.3 till 15.12 (expect maybe 15.8)

Agenda

- Clustering measures
- Clustering methods (hard assignments)
 - Hierarchical clustering
 - k -means clustering
- Cluster validation
 - Fusion graphs
 - The Davies-Bouldin Index

Yes, you also can do soft assignments...

Dissimilarity measures



$$D = \begin{bmatrix} 0 & d(\mathbf{A}, \mathbf{B}) & d(\mathbf{A}, \mathbf{C}) \\ & 0 & d(\mathbf{B}, \mathbf{C}) \\ & & 0 \end{bmatrix}$$
$$= \begin{bmatrix} 0 & 10 & 11 \\ & 0 & 2 \\ & & 0 \end{bmatrix}$$

Dissimilarity measures

- Let $d(r, s)$ be the dissimilarity between objects r and s
- Formally, dissimilarity measures should satisfy

$$d(r, s) \geq 0, \forall r, s$$

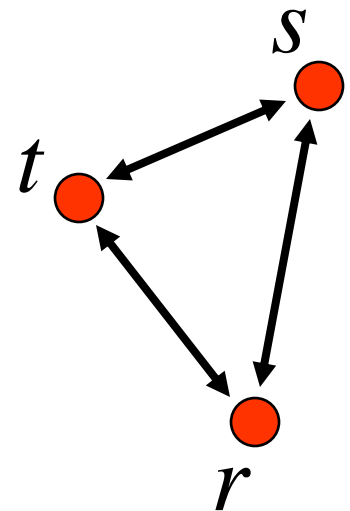
$$d(r, r) = 0, \forall r$$

$$d(r, s) = d(s, r), \forall r, s$$

- If in addition, the triangle inequality holds, the measure is a metric

$$d(r, t) + d(t, s) \geq d(r, s), \forall r, s, t$$

- **Most often used: Euclidean distance (metric)**



Distance measure

- Define a distance between objects:

- Euclidean:
$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^l (x_i - y_i)^2}$$

- City-block
$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^l |x_i - y_i|$$

- ℓ_p -metric
$$d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^l |x_i - y_i|^p \right)^{1/p}$$

More similarity measures

- Cosine similarity

$$s_{cos}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

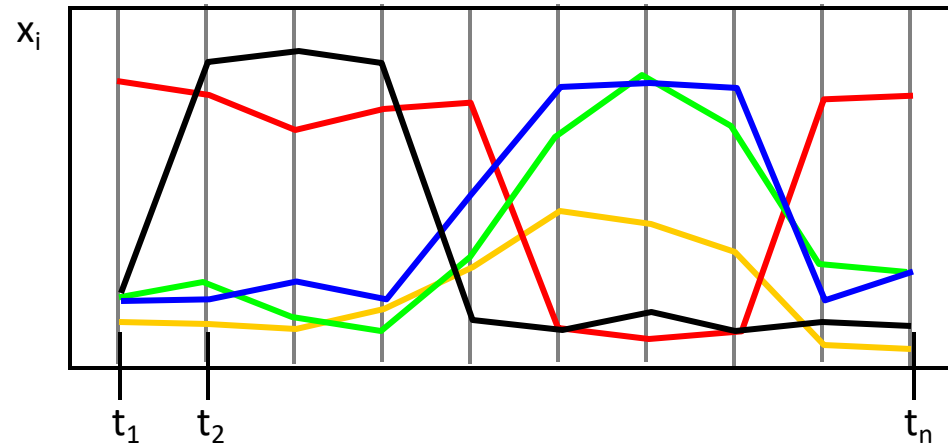
- Pearson's correlation coefficient

$$r_{Pearson}(\mathbf{x}, \mathbf{y}) = \frac{(\mathbf{x} - \mu_x)^T (\mathbf{y} - \mu_y)}{\|\mathbf{x} - \mu_x\| \|\mathbf{y} - \mu_y\|}$$

- and more... (for discrete features, mixed features, categorical features, ...)

Using different measures

- Example:
time series data
(gene expression)



Euclidean distance
match exact shape

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_{t=1}^n (x_{i,t} - x_{j,t})^2$$

$$\begin{aligned} d(\text{blue}, \text{green}) &< d(\text{blue}, \text{yellow}) \\ d(\text{blue}, \text{green}) &\ll d(\text{blue}, \text{red}) \\ d(\text{blue}, \text{green}) &\ll d(\text{blue}, \text{black}) \end{aligned}$$

Pearson correlation
ignore amplitude

$$1 - \rho_{ij}$$

$$\begin{aligned} d(\text{blue}, \text{green}) &\approx d(\text{blue}, \text{yellow}) \\ d(\text{blue}, \text{green}) &\ll d(\text{blue}, \text{red}) \\ d(\text{blue}, \text{green}) &\ll d(\text{blue}, \text{black}) \end{aligned}$$

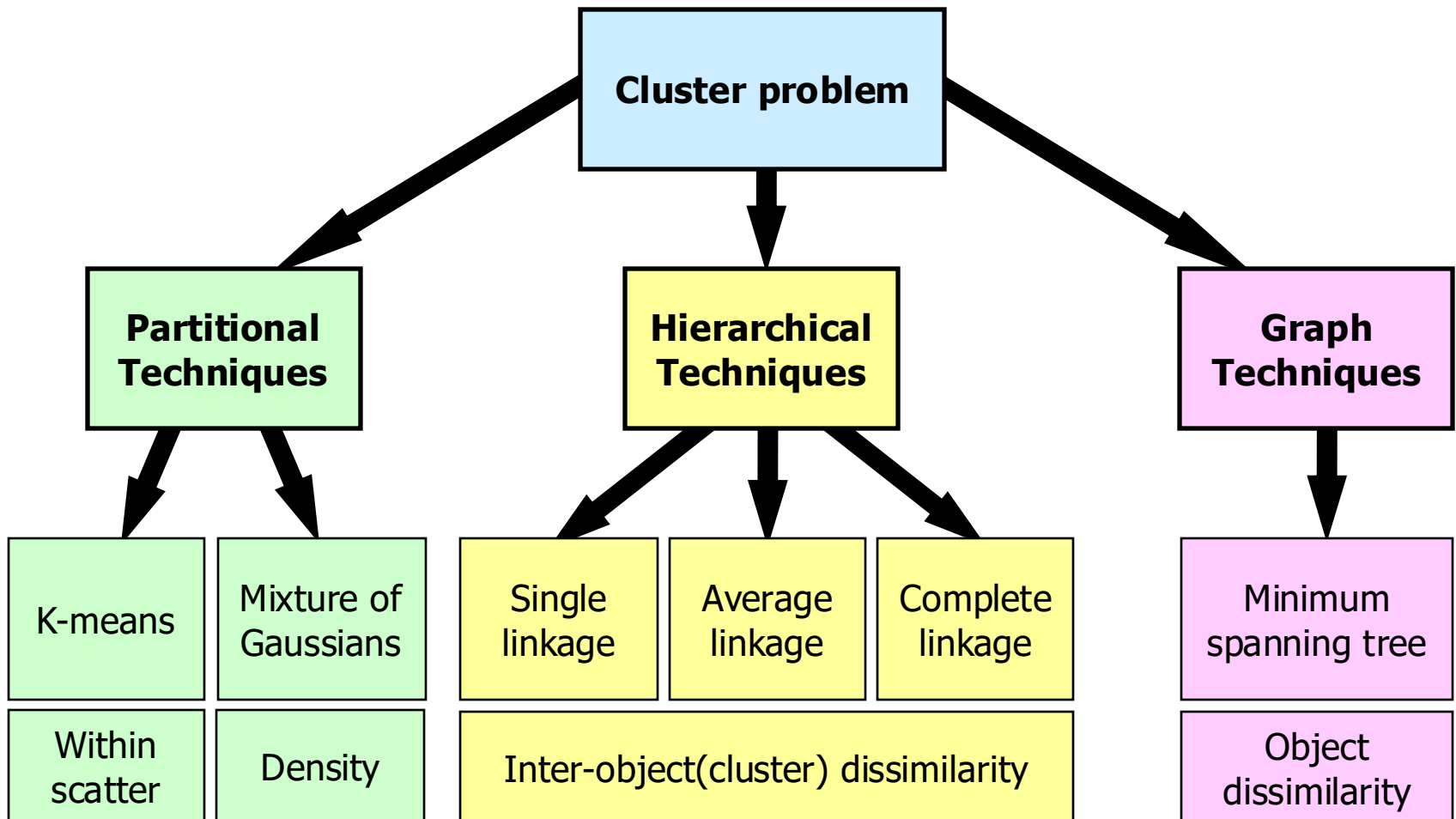
Absolute correlation
ignore amplitude & sign

$$1 - |\rho_{ij}|$$

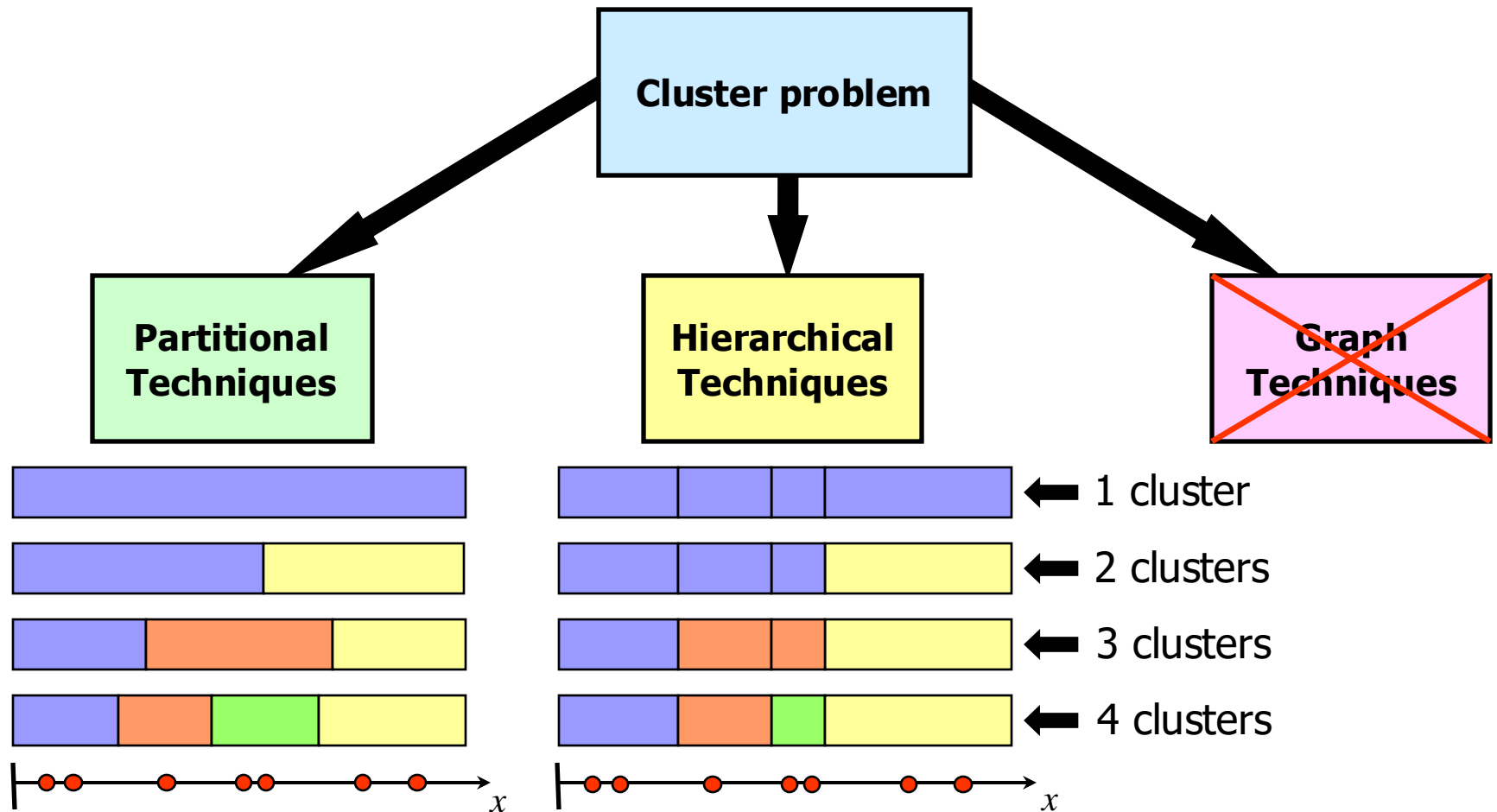
$$\begin{aligned} d(\text{blue}, \text{green}) &\approx d(\text{blue}, \text{yellow}) \\ d(\text{blue}, \text{green}) &\approx d(\text{blue}, \text{red}) \\ d(\text{blue}, \text{green}) &\ll d(\text{blue}, \text{black}) \end{aligned}$$

$$\rho_{ij} = \frac{\sum_{t=1}^n (x_{i,t} - \mu_i)(x_{j,t} - \mu_j)}{\sigma_i \sigma_j}$$

Clustering techniques

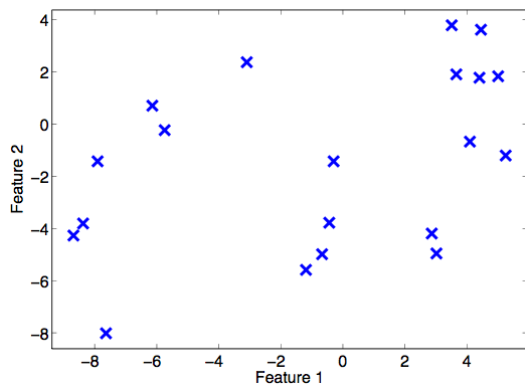


Clustering techniques

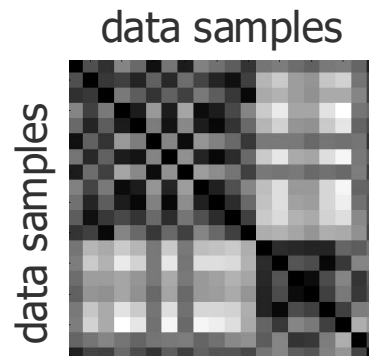


Agglomerative Hierarchical Clustering

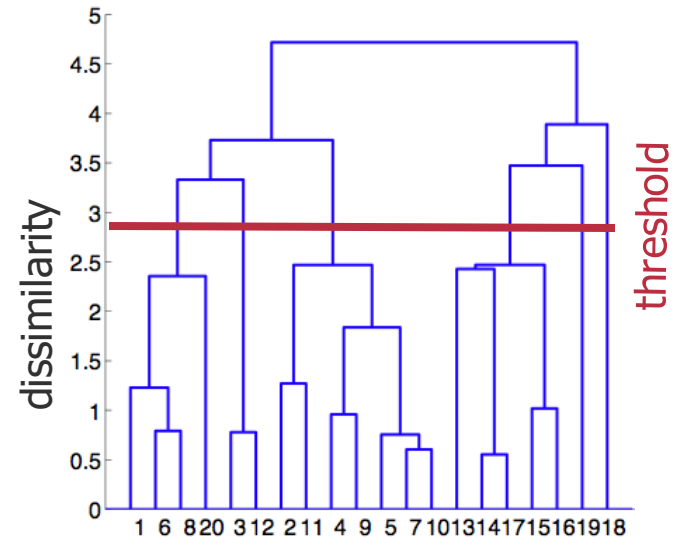
- **Starting from individual observations**, produce sequence of clusterings of increasing size
- At **each level**, **two clusters** chosen by criterion are merged



2D scatter plot of data



dissimilarity matrix



dendrogram

Agglomerative Hierarchical Clustering

1. Determine distances between all clusters
 2. Merge clusters that are **closest**
 3. IF #clusters > 1 THEN GOTO 1
- Which clusters to start with?
 - What is the distance between clusters?
 - Final number of clusters?

Different Merging Rules

- Two nearest objects in the clusters : single linkage

$$g(R, S) = \min_{ij} \{d(\mathbf{x}_i, \mathbf{x}_j) : \mathbf{x}_i \in R, \mathbf{x}_j \in S\}$$

- Two most remote objects in the clusters : complete linkage

$$g(R, S) = \max_{ij} \{d(\mathbf{x}_i, \mathbf{x}_j) : \mathbf{x}_i \in R, \mathbf{x}_j \in S\}$$

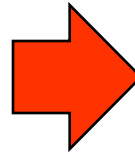
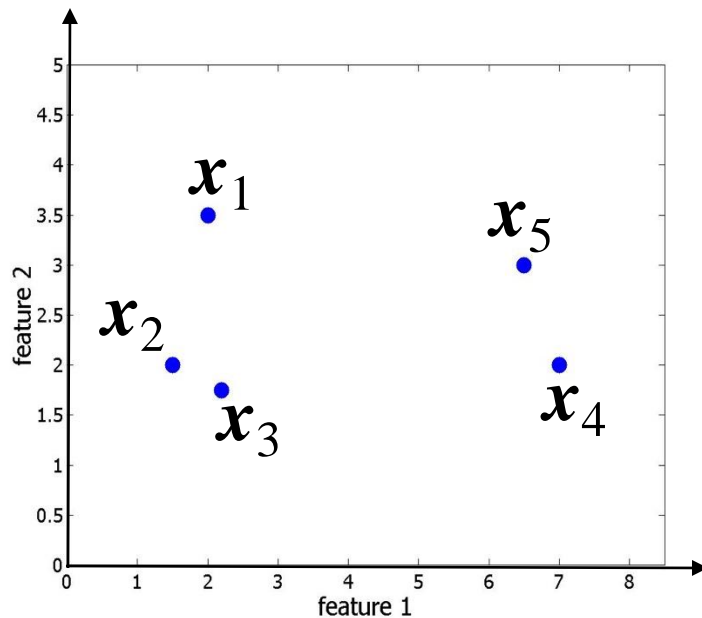
- Cluster centers : average linkage

$$g(R, S) = \frac{1}{|R||S|} \sum_{ij} \{d(\mathbf{x}_i, \mathbf{x}_j) : \mathbf{x}_i \in R, \mathbf{x}_j \in S\}$$

Hierarchical clustering

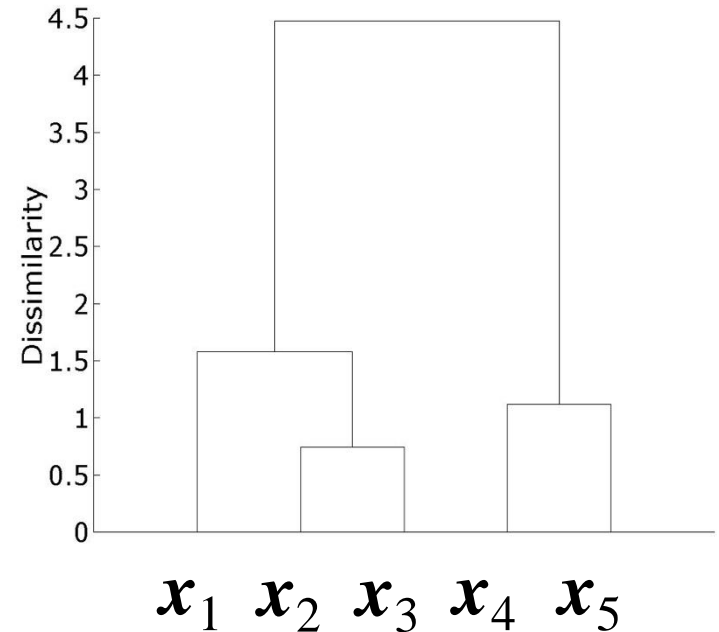
Input:

- dataset, $X: [n \times p]$, or directly:
- dissimilarity matrix, $D: [n \times n]$
- linkage type



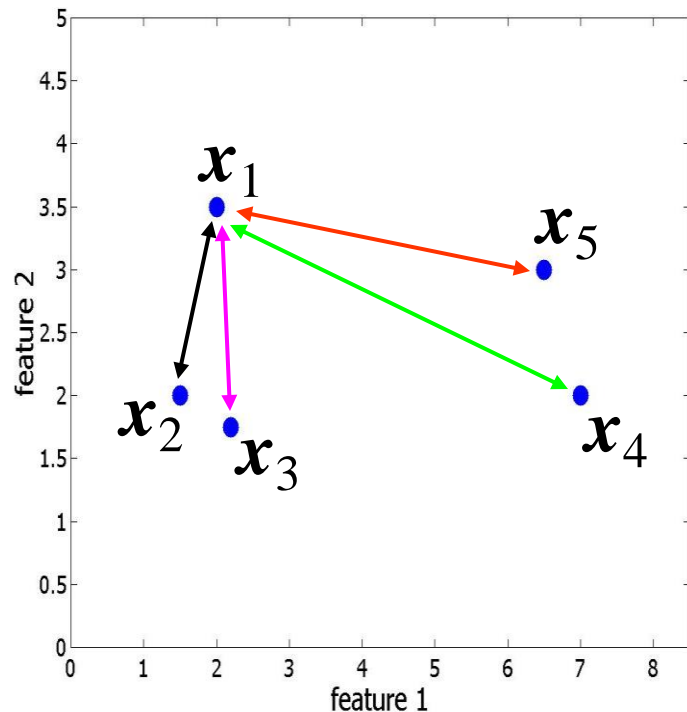
Output:

- dendrogram



Hierarchical clustering (2)

Dataset



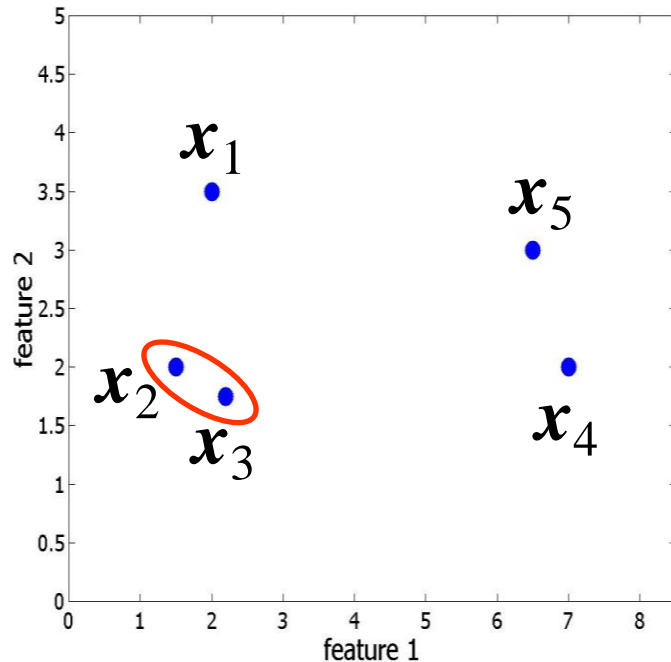
Euclidean distance matrix, D

	x_1	x_2	x_3	x_4	x_5
x_1	0.00	1.58	1.76	5.22	4.53
x_2		0.00	0.74	5.50	5.10
x_3			0.00	4.81	4.48
x_4				0.00	1.12
x_5					0.00

Hierarchical clustering (3)

- **Step 1:**

Find the most similar pair of objects: $\min_{(i,j)} \{d(i,j)\} = d(2,3)$

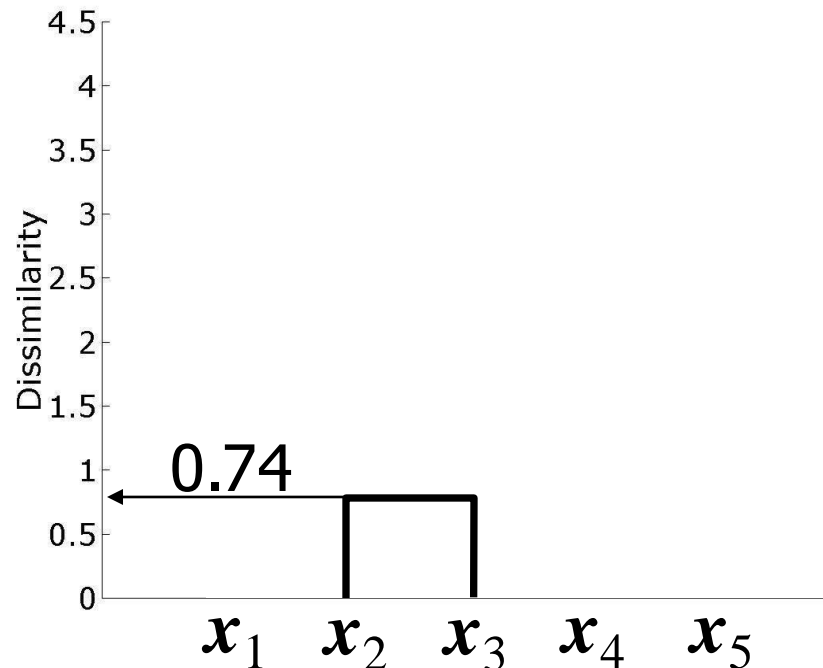


	x_1	x_2	x_3	x_4	x_5
x_1	0.00	1.58	1.76	5.22	4.53
x_2		0.00	0.74	5.50	5.10
x_3			0.00	4.81	4.48
x_4				0.00	1.12
x_5					0.00

Hierarchical clustering (4)

- **Step 2:**

Merge x_2 and x_3 into a single object, $[x_2, x_3]$;

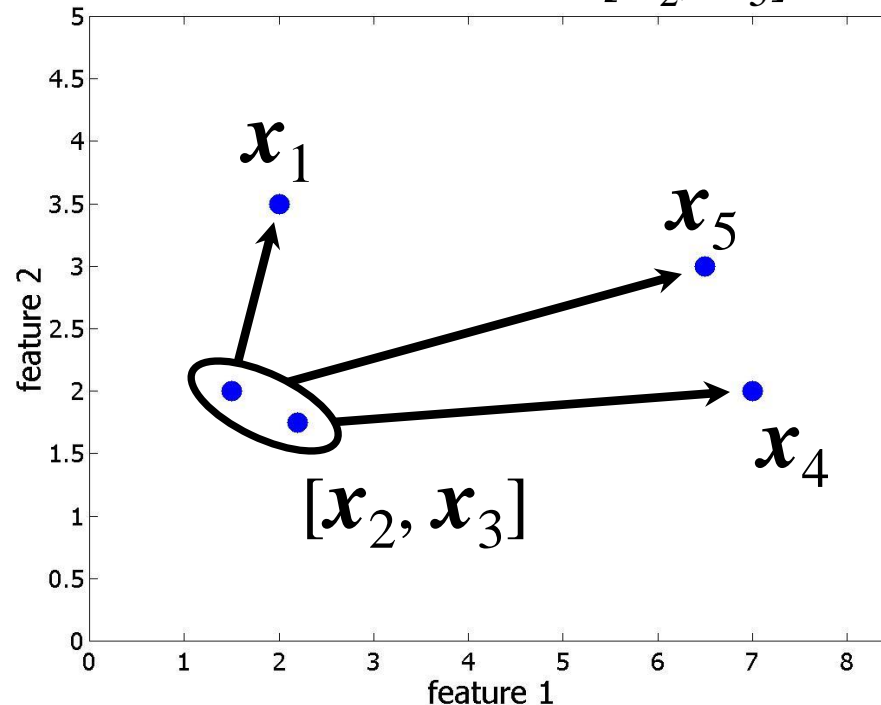


Hierarchical clustering (5)

- **Step 3:**

Recompute D –

what is the distance between $[x_2, x_3]$ and the rest?

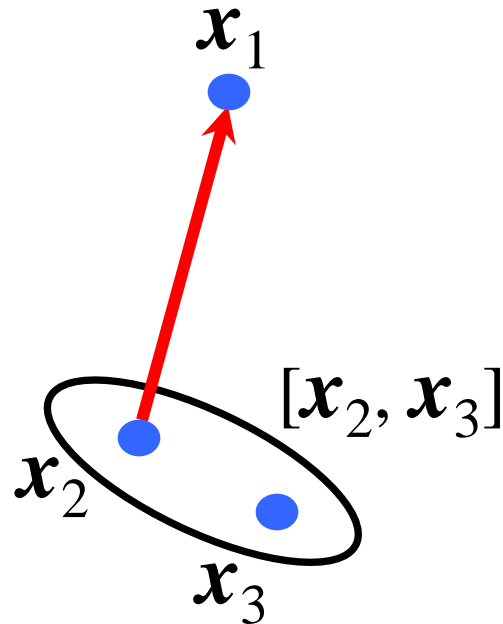


Hierarchical clustering (6)

- **Step 3:**

Recompute D –

single linkage: $d([x_2, x_3], x_1) = \min(d(x_1, x_2), d(x_1, x_3))$

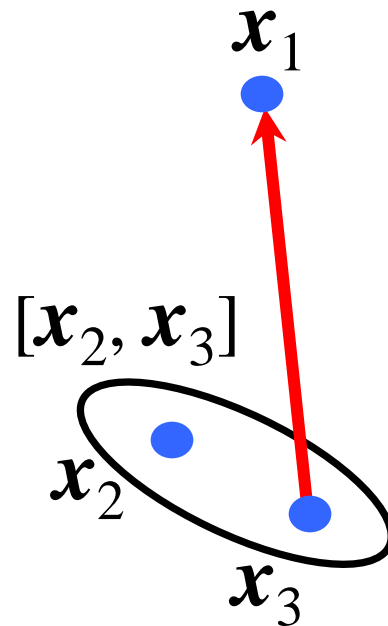


Hierarchical clustering (7)

- **Step 3:**

Recompute D –

complete linkage: $d([x_2, x_3], x_1) = \max(d(x_1, x_2), d(x_1, x_3))$

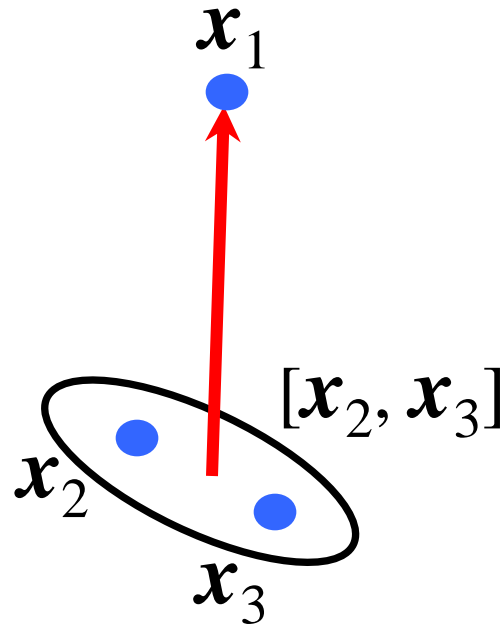


Hierarchical clustering (8)

- **Step 3:**

Recompute D –

average linkage: $d([x_2, x_3], x_1) = \text{mean}(d(x_1, x_2), d(x_1, x_3))$



Hierarchical clustering (9)

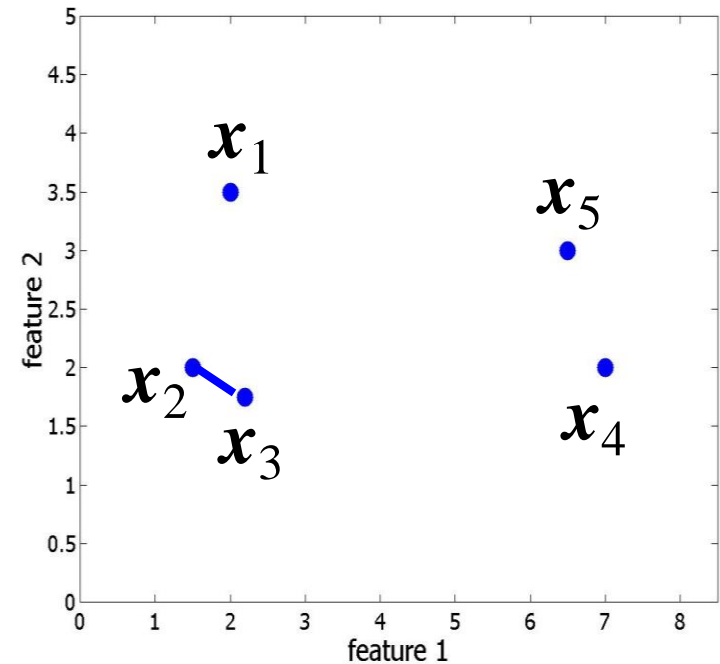
- **Step 3:**
Recompute D – **single linkage:**

	x_1	$[x_2, x_3]$	x_4	x_5
x_1	0.00	1.58	5.22	4.53
$[x_2, x_3]$		0.00	4.81	4.48
x_4			0.00	1.12
x_5				0.00

Hierarchical clustering (10)

- **Repeat, step 1:**

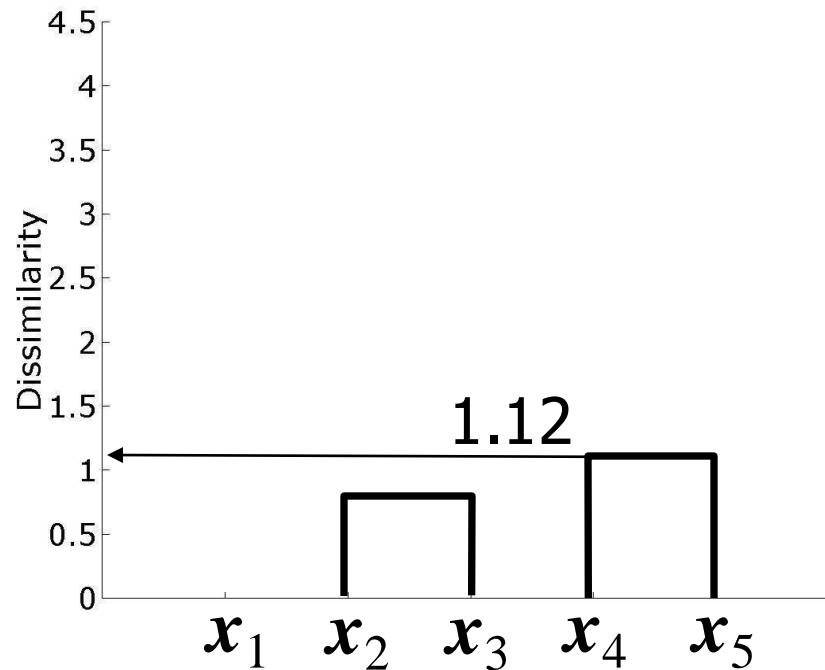
Find the most similar pair of objects: $\min_{(i,j)} \{d(i,j)\} = d(4,5)$



	x_1	$[x_2, x_3]$	x_4	x_5
x_1	0.00	1.58	5.22	4.53
$[x_2, x_3]$	1.58	0.00	4.81	4.48
x_4	5.22	4.81	0.00	1.12
x_5	4.53	4.48	1.12	0.00

Hierarchical clustering (11)

- **Repeat, step 2:**
Merge x_4 and x_5 into a single object, $[x_4, x_5]$;



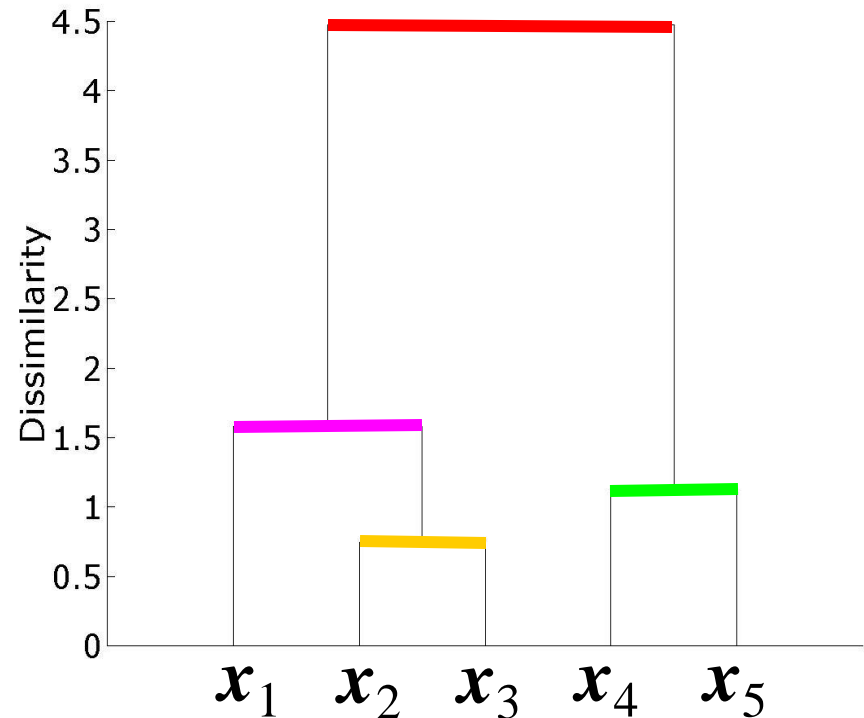
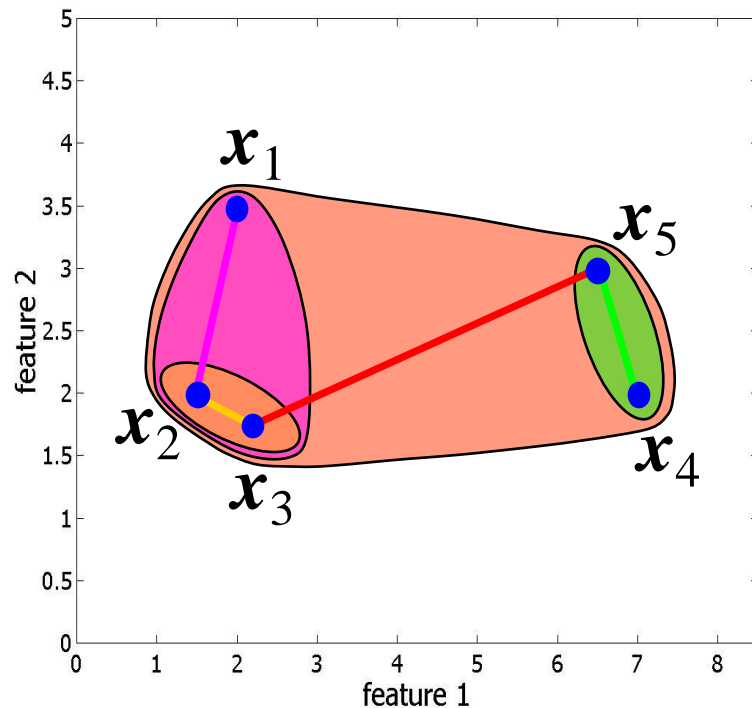
Hierarchical clustering (12)

- **Repeat, step 3:**
Recompute D (single linkage):

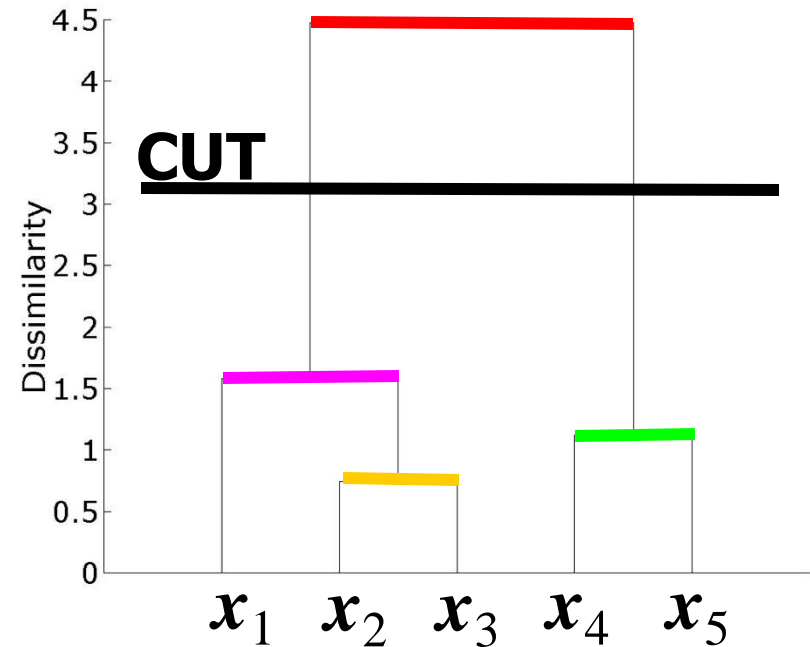
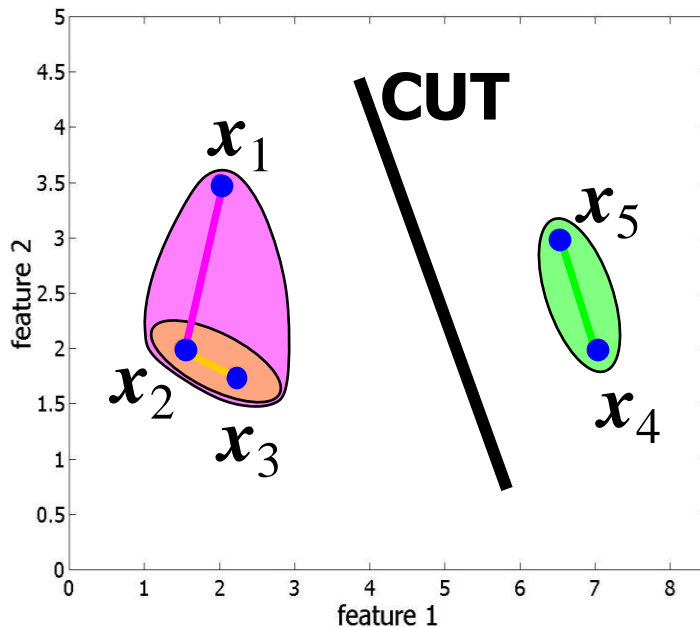
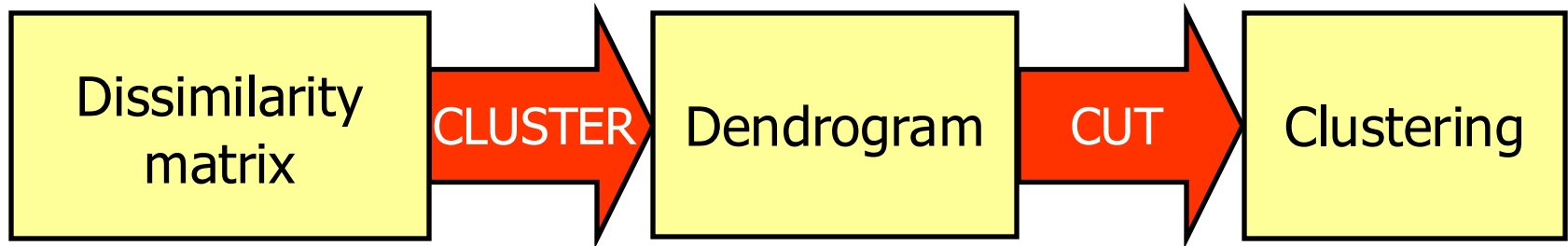
	x_1	$[x_2, x_3]$	$[x_4, x_5]$
x_1	0.00	1.58	4.53
$[x_2, x_3]$		0.00	4.48
$[x_4, x_5]$			0.00

Hierarchical clustering (13)

- Repeat steps 1-3 until a single cluster remains...

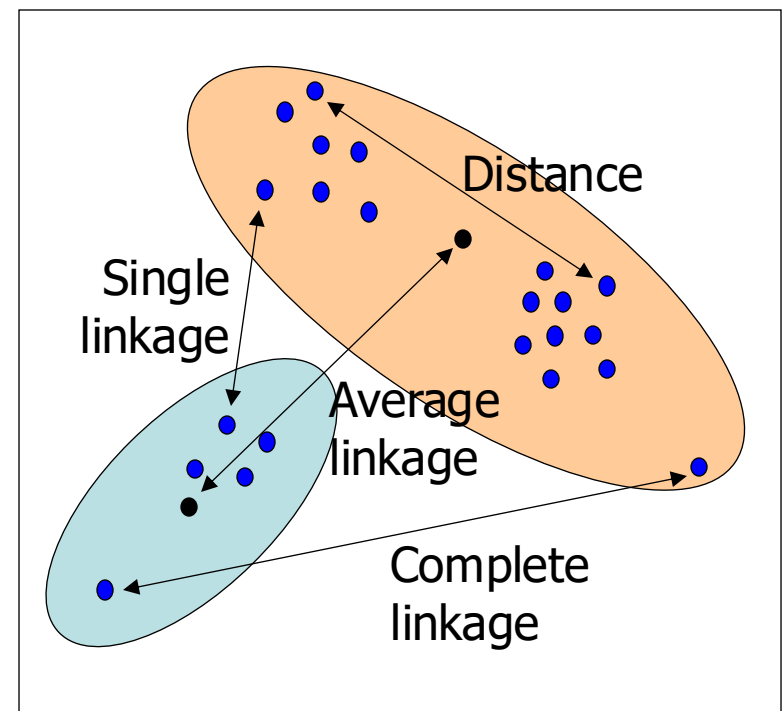


Hierarchical clustering (14)

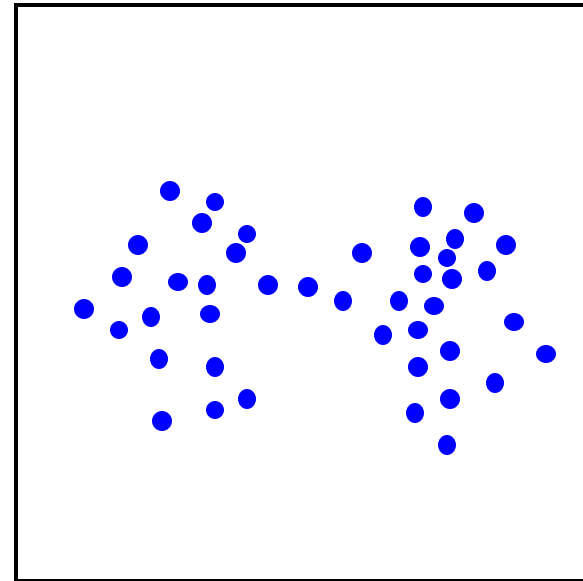
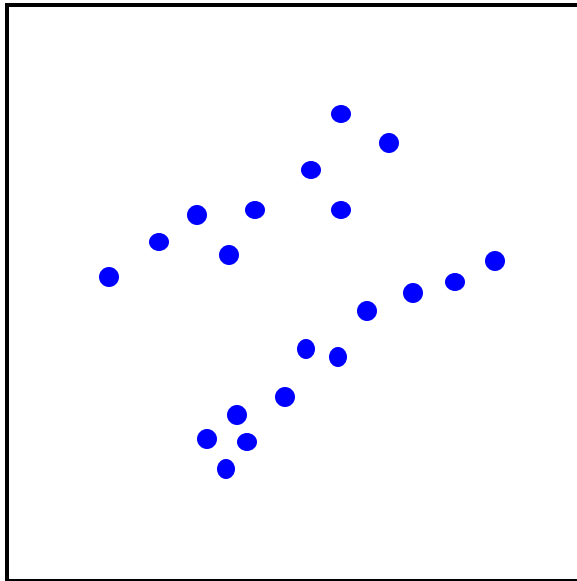


Hierarchical clustering (15)

- Hierarchical clustering: repeatedly group closest clusters
- Important choices:
 - *Distance measure* between objects: Euclidean, correlation, ...
 - *Linkage* between clusters: single, average, complete



Linkage and cluster shape

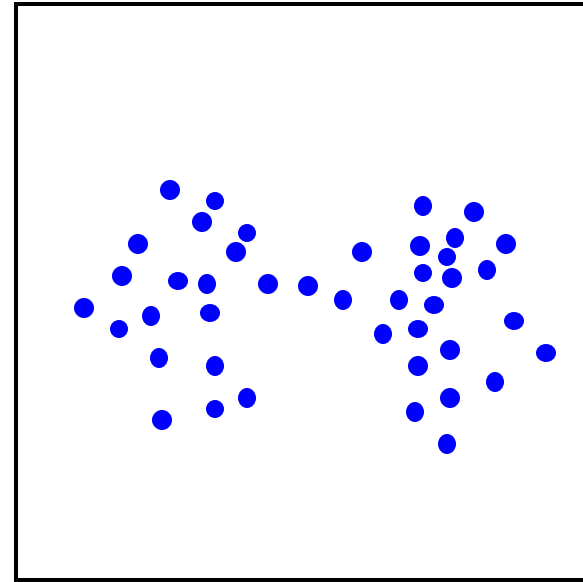
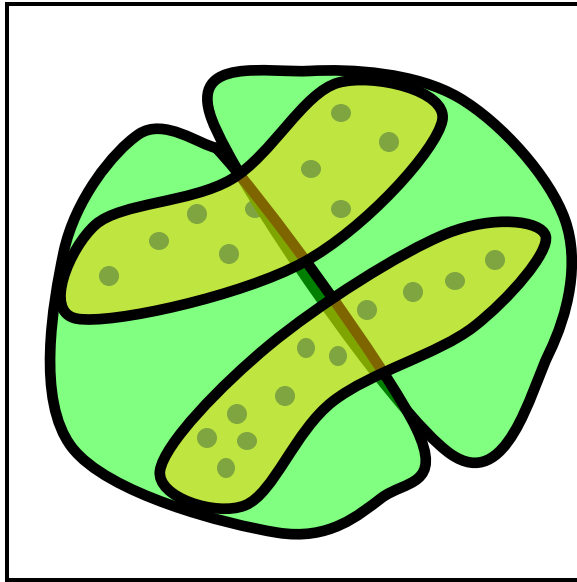


complete linkage



single linkage

Linkage and cluster shape (2)

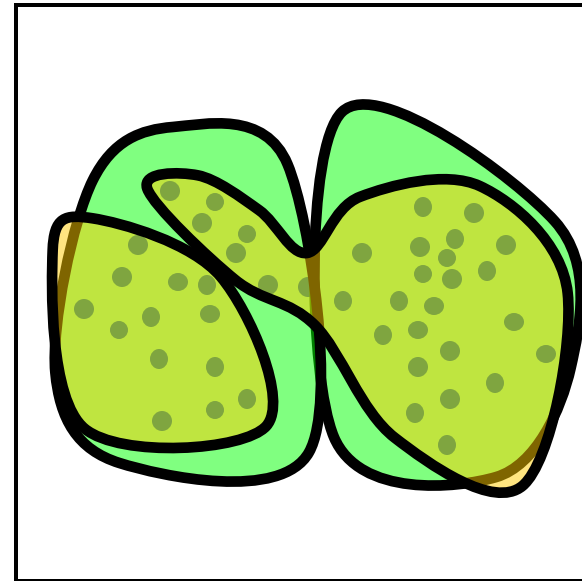
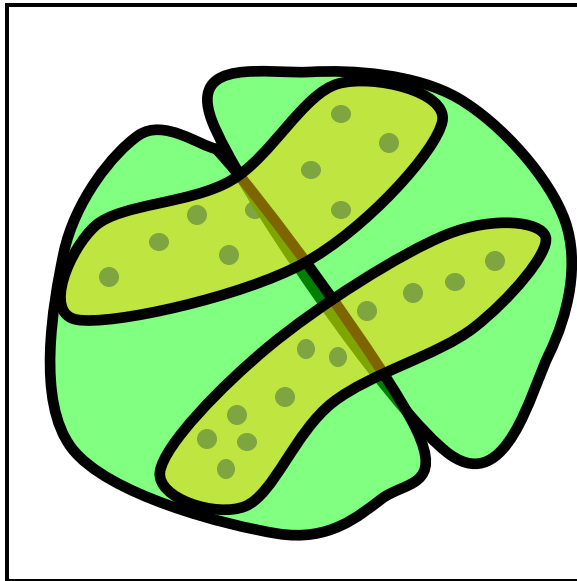


Complete linkage



Single linkage

Linkage and cluster shape (3)



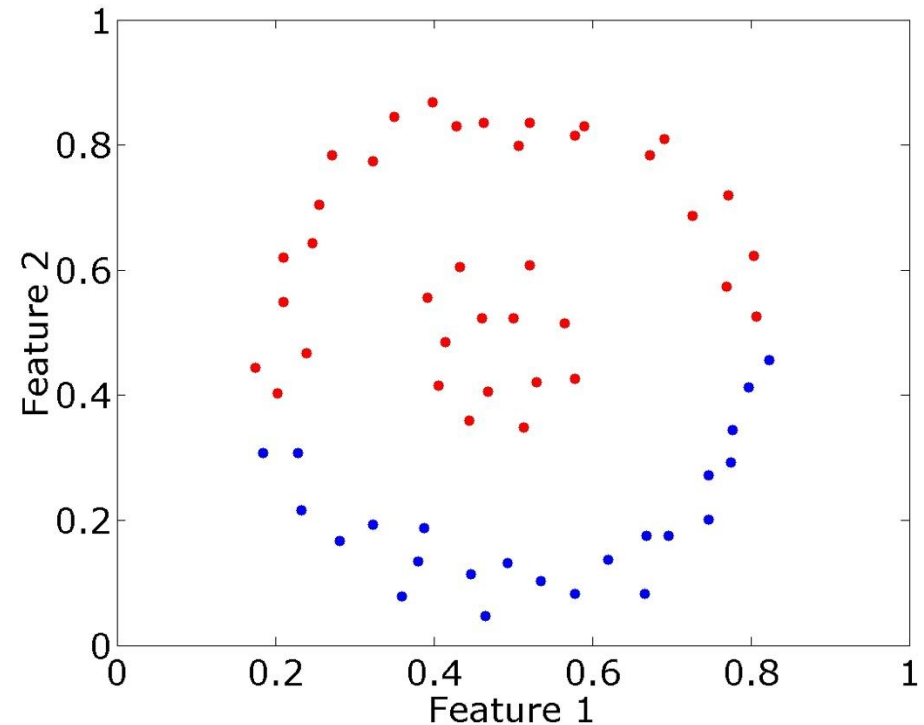
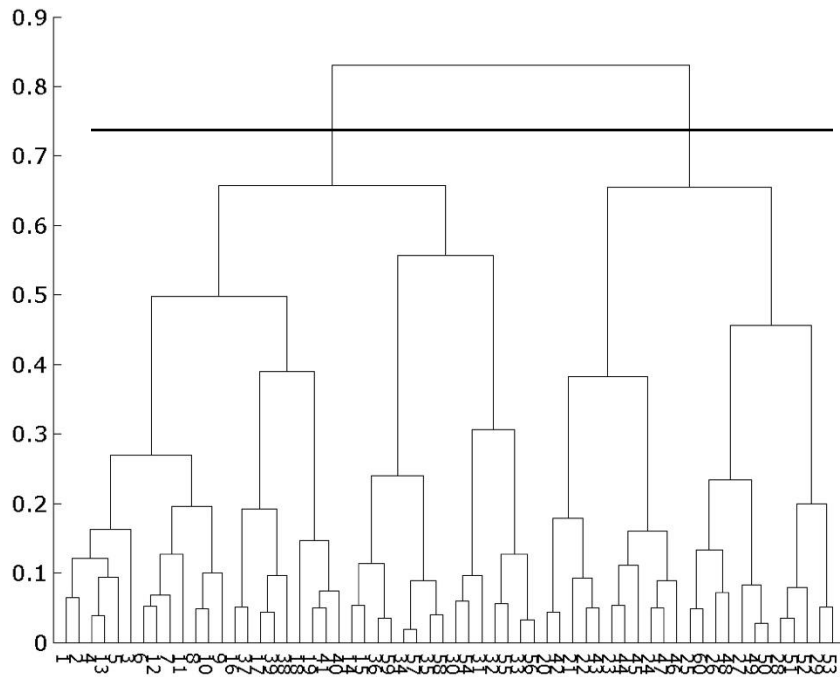
Complete linkage



Single linkage

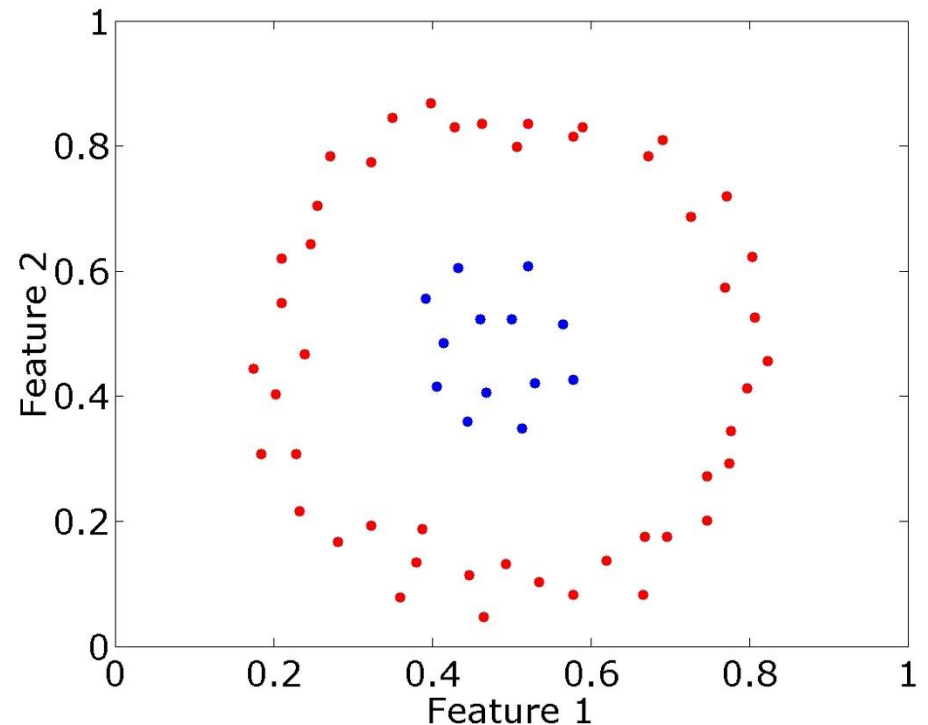
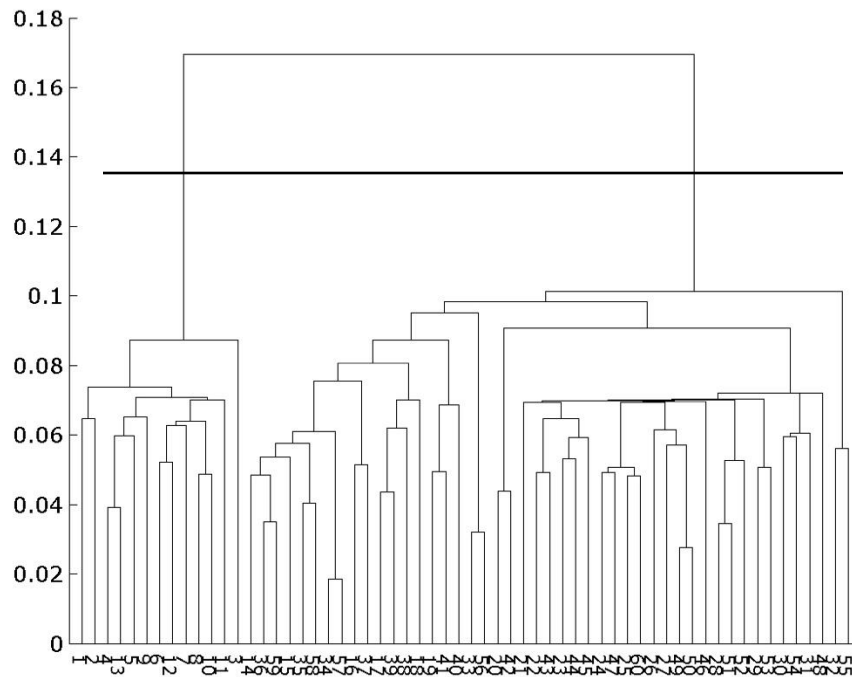
Hierarchical clustering examples

Euclidean, complete linkage



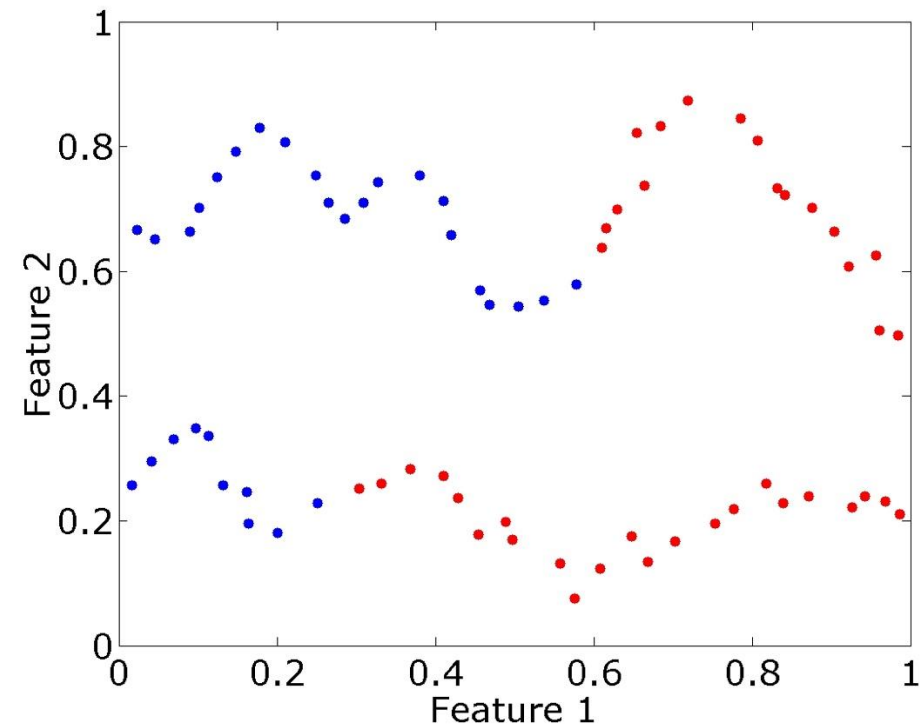
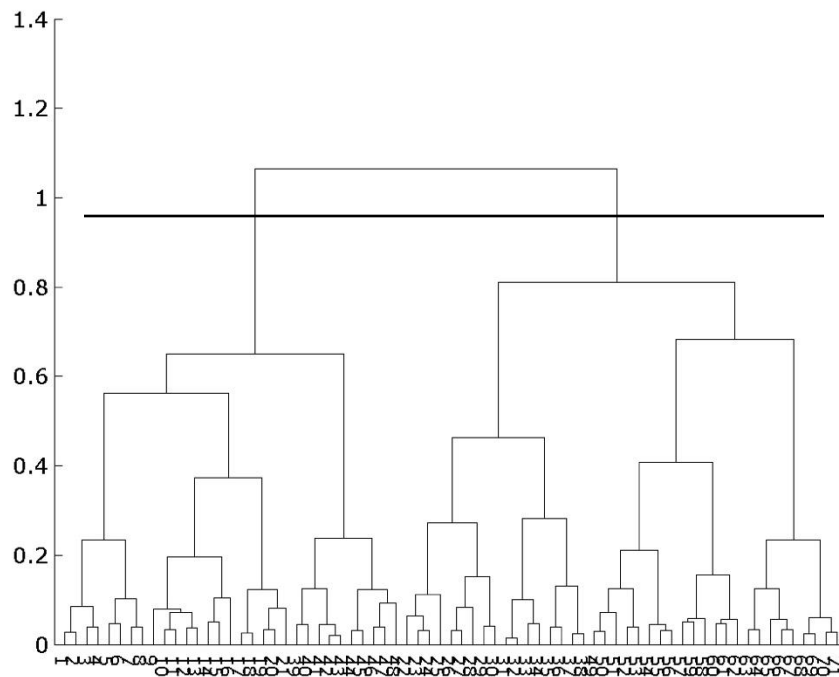
Hierarchical clustering examples

Euclidean, single linkage



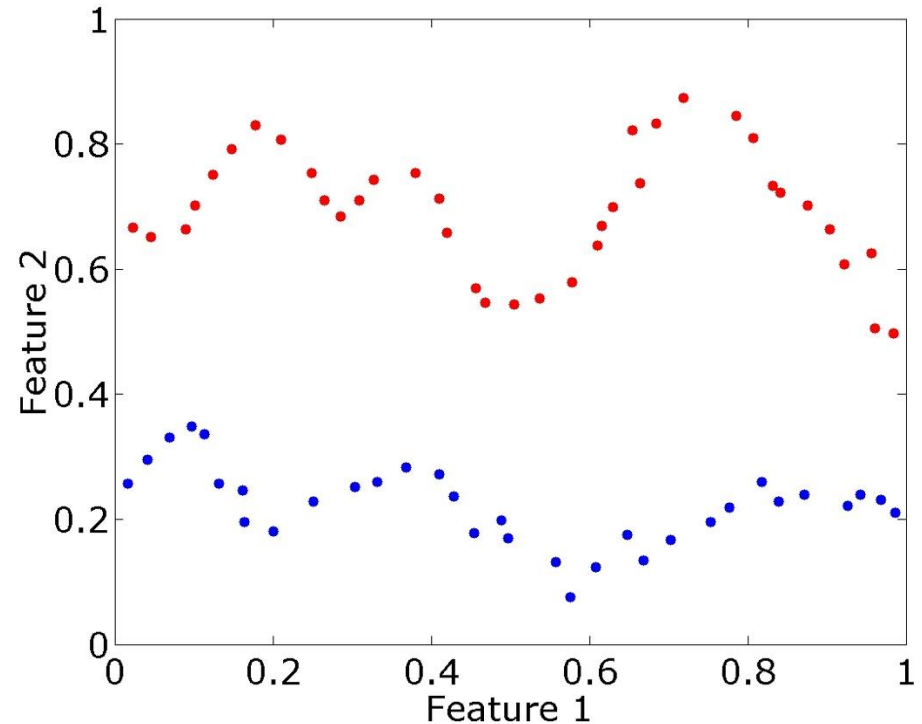
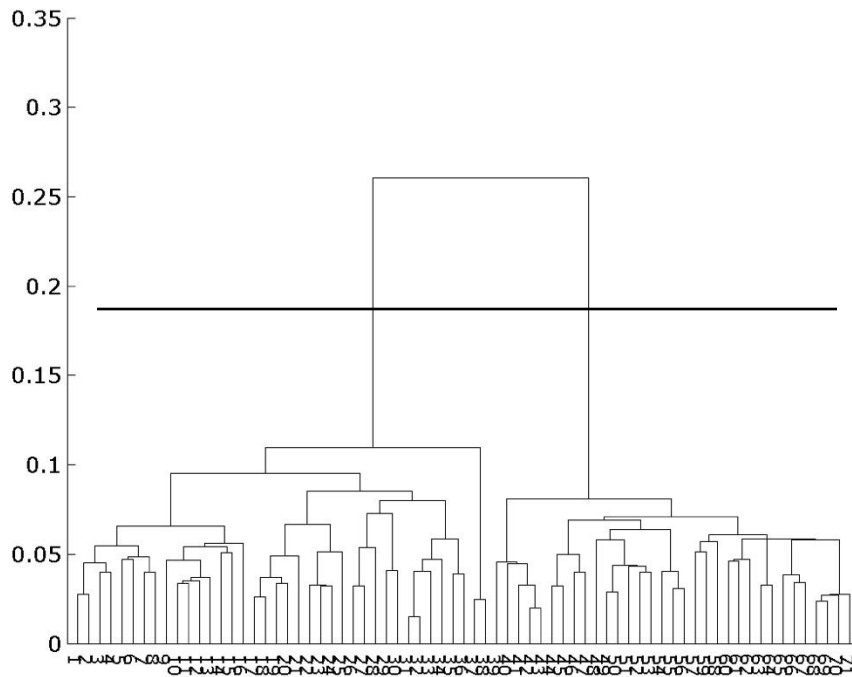
Hierarchical clustering examples

Euclidean, complete linkage



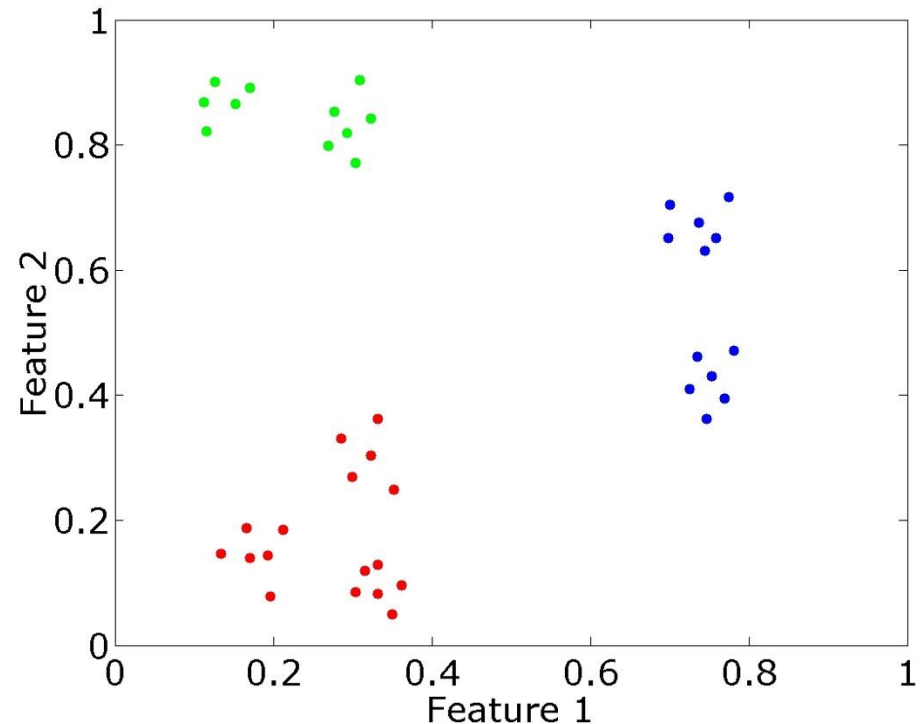
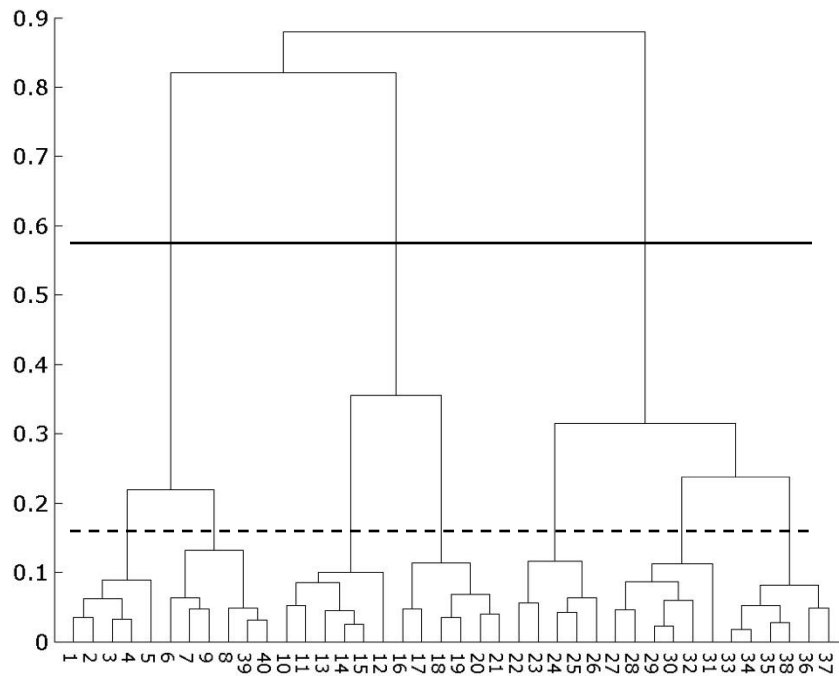
Hierarchical clustering examples

Euclidean, single linkage



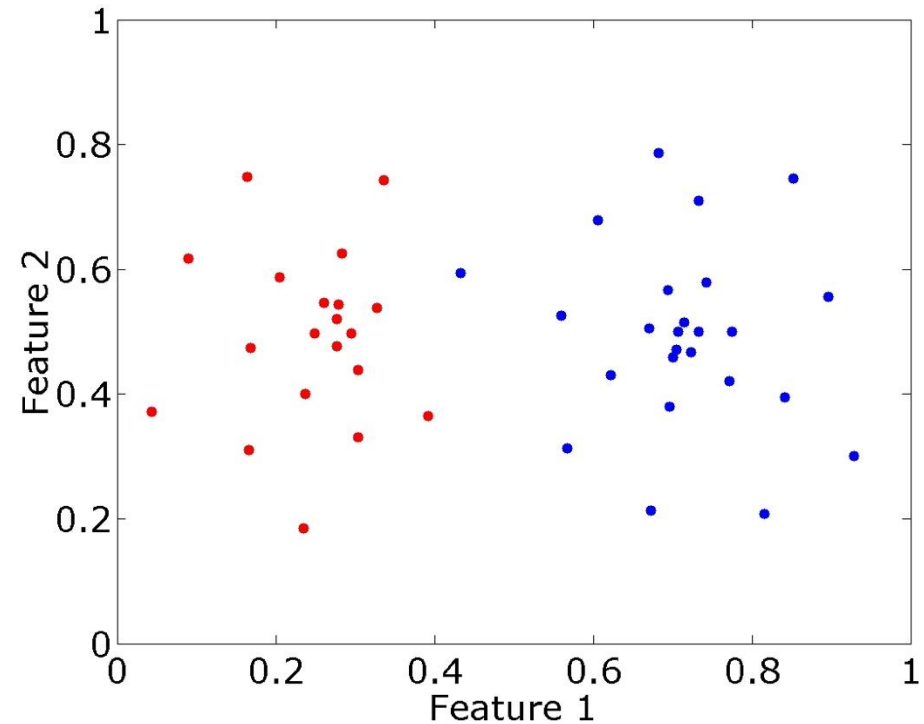
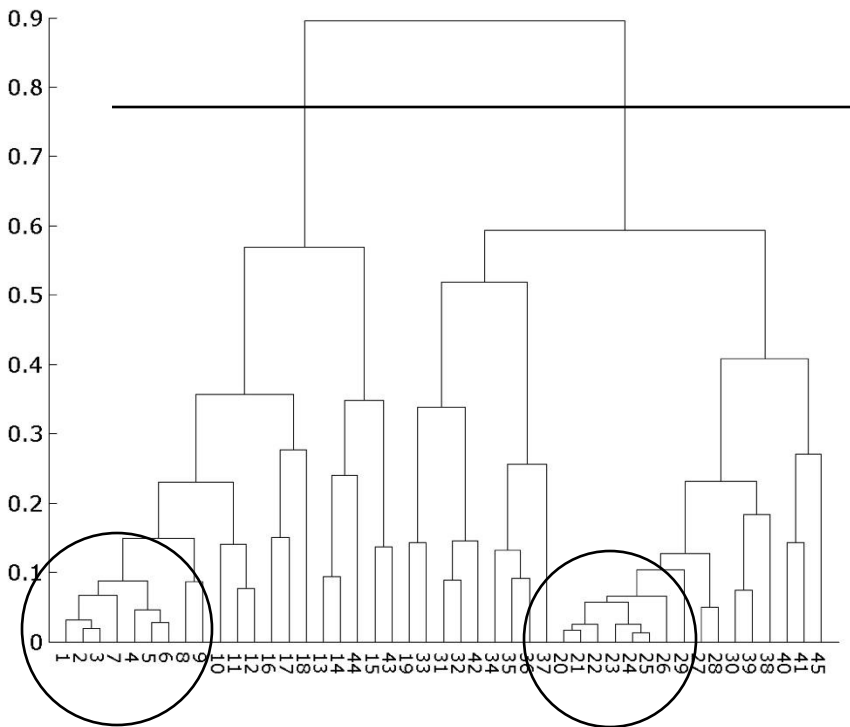
Hierarchical clustering examples

Euclidean, complete linkage

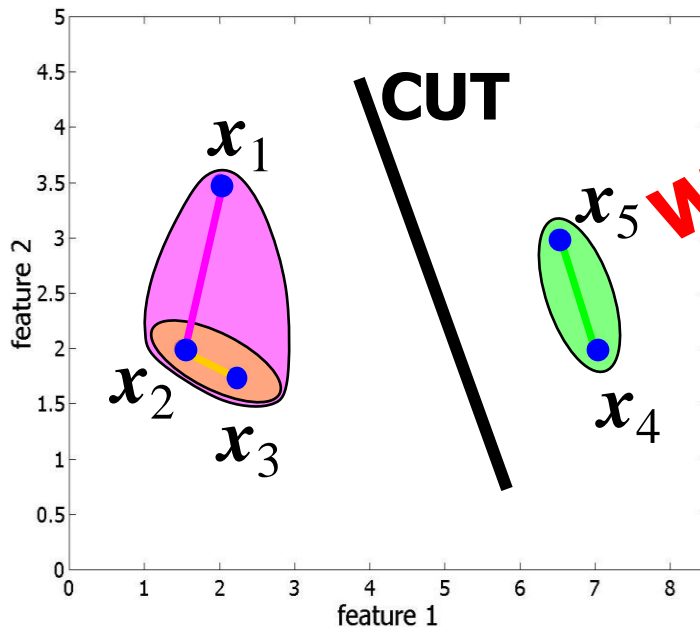
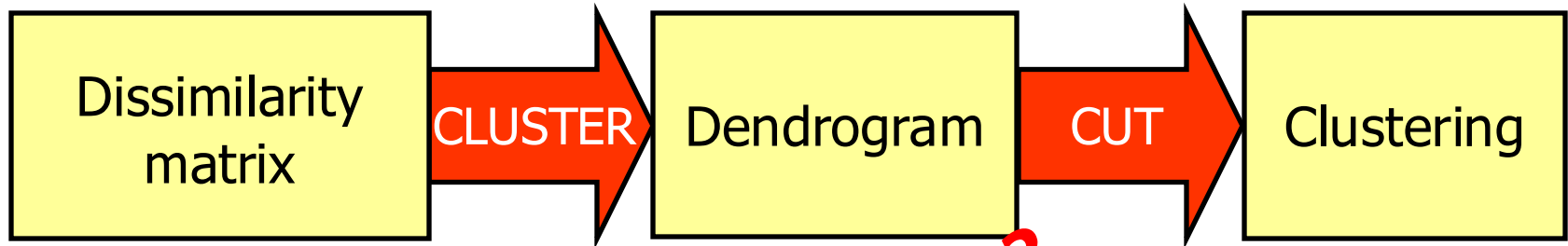


Hierarchical clustering examples

Euclidean, complete linkage

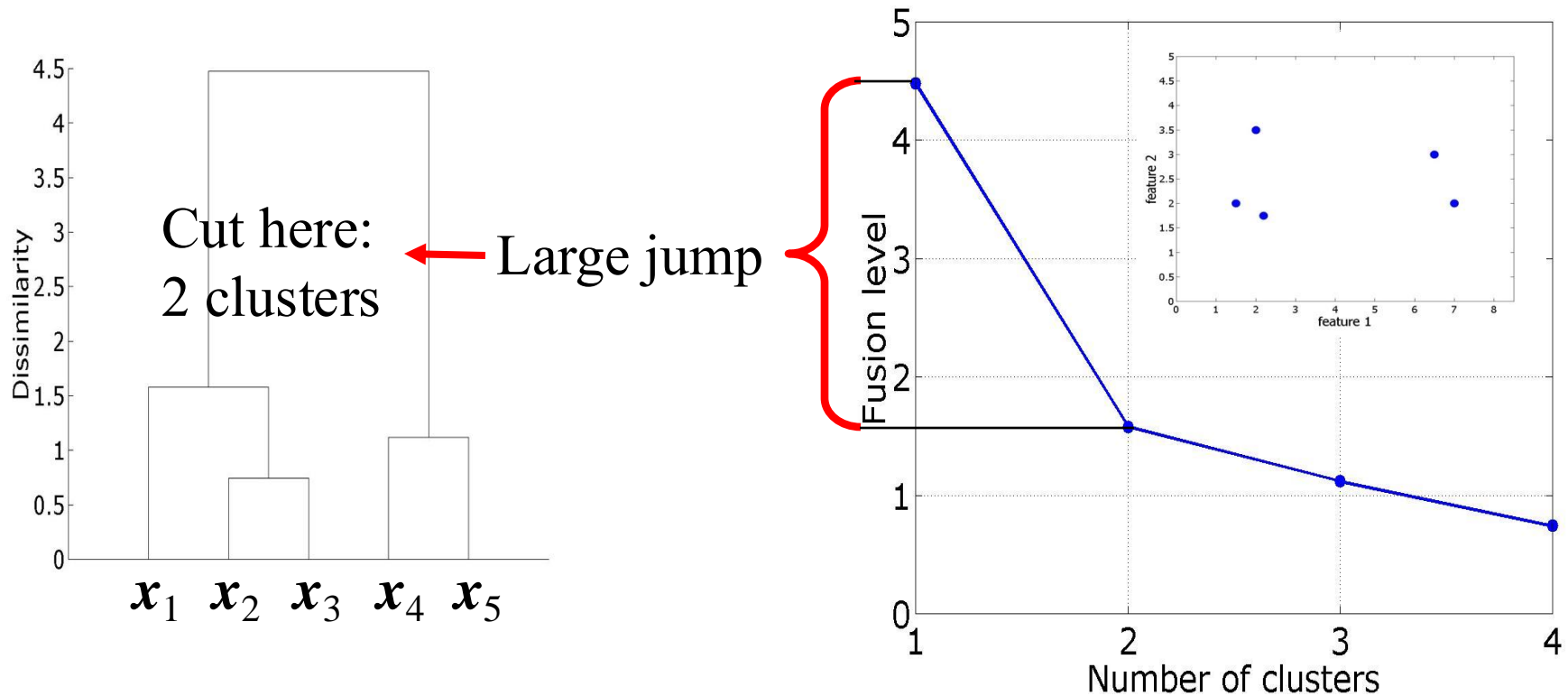


Remember our example...?



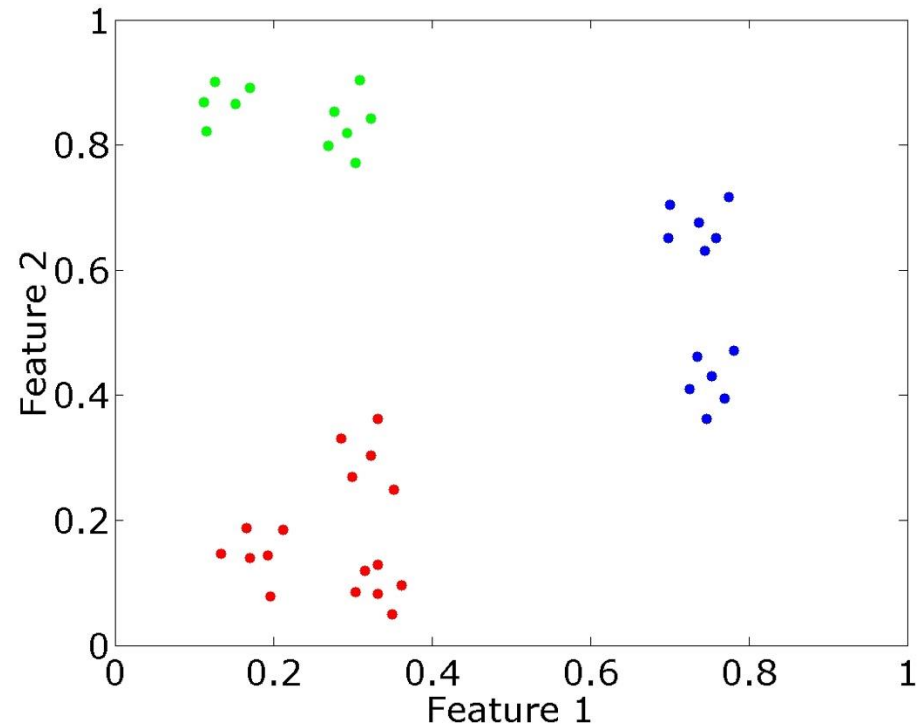
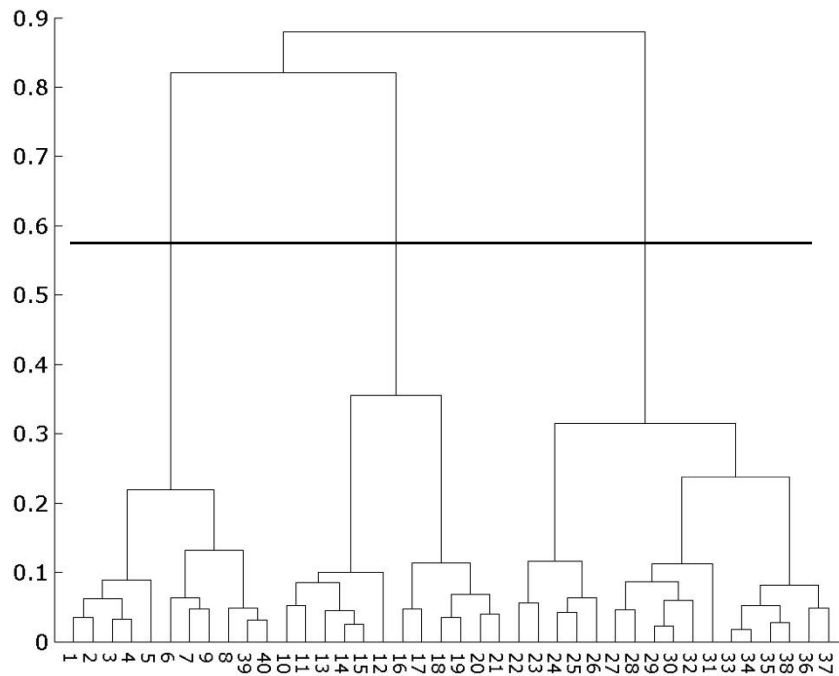
Fusion graph

- fusion level



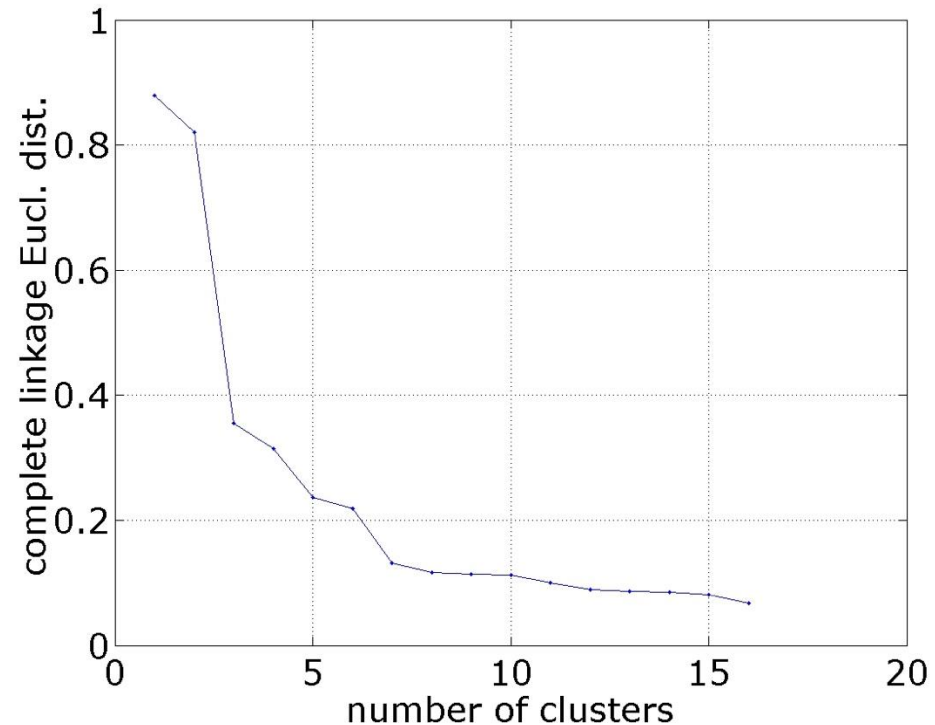
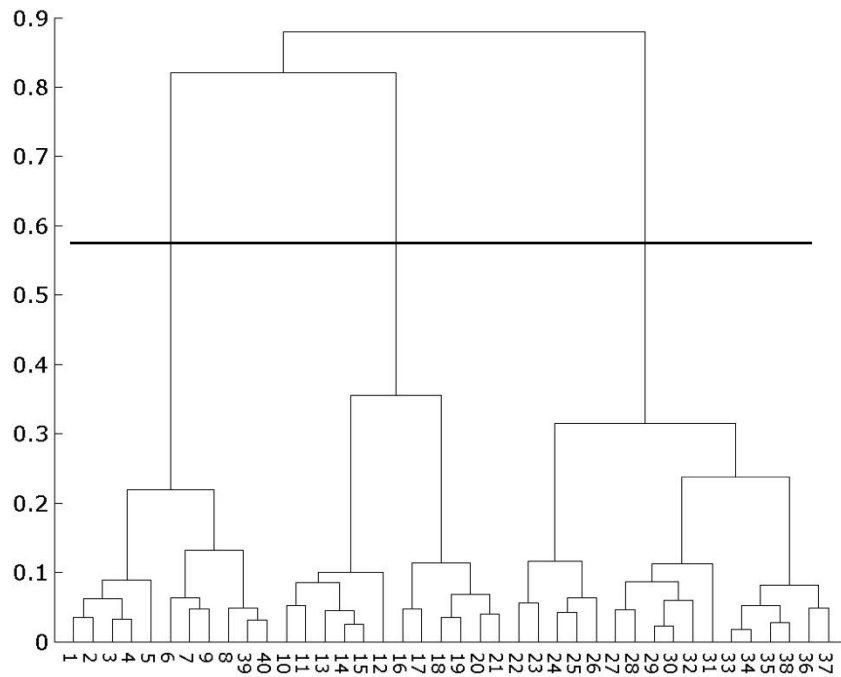
Fusion graph (2)

(Euclidean; complete linkage)



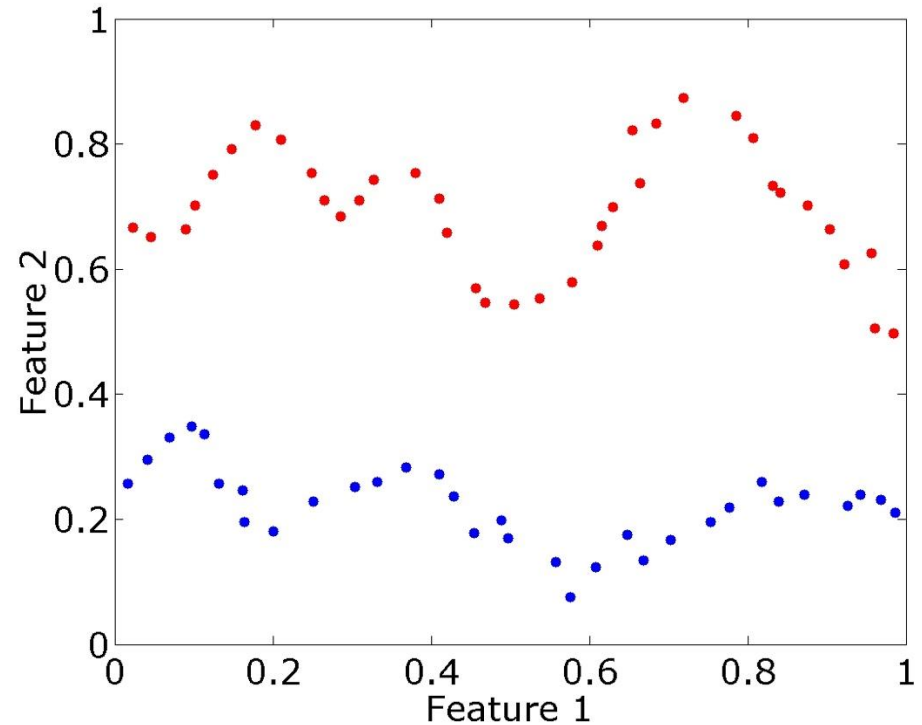
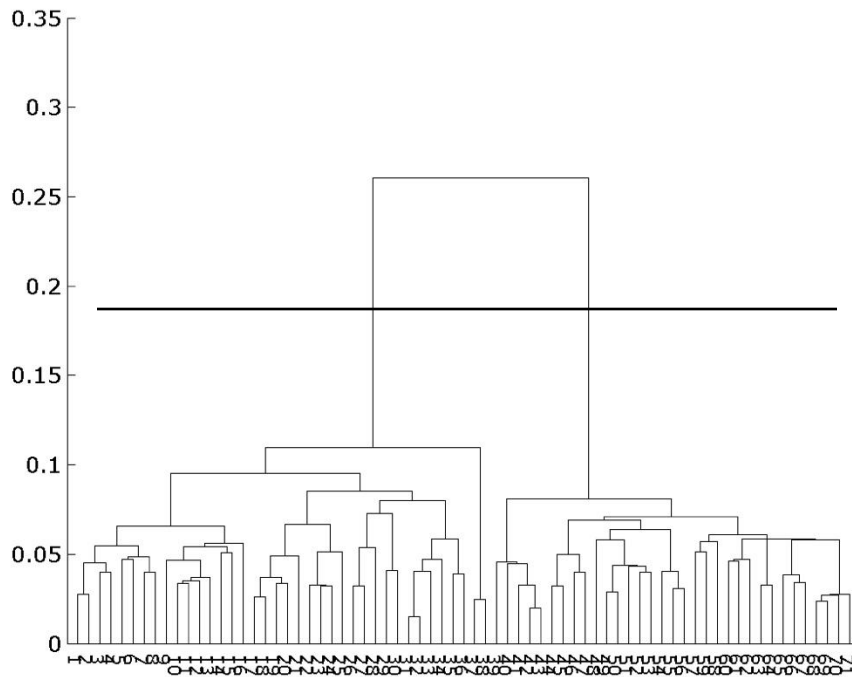
Fusion graph (3)

(Euclidean; complete linkage)



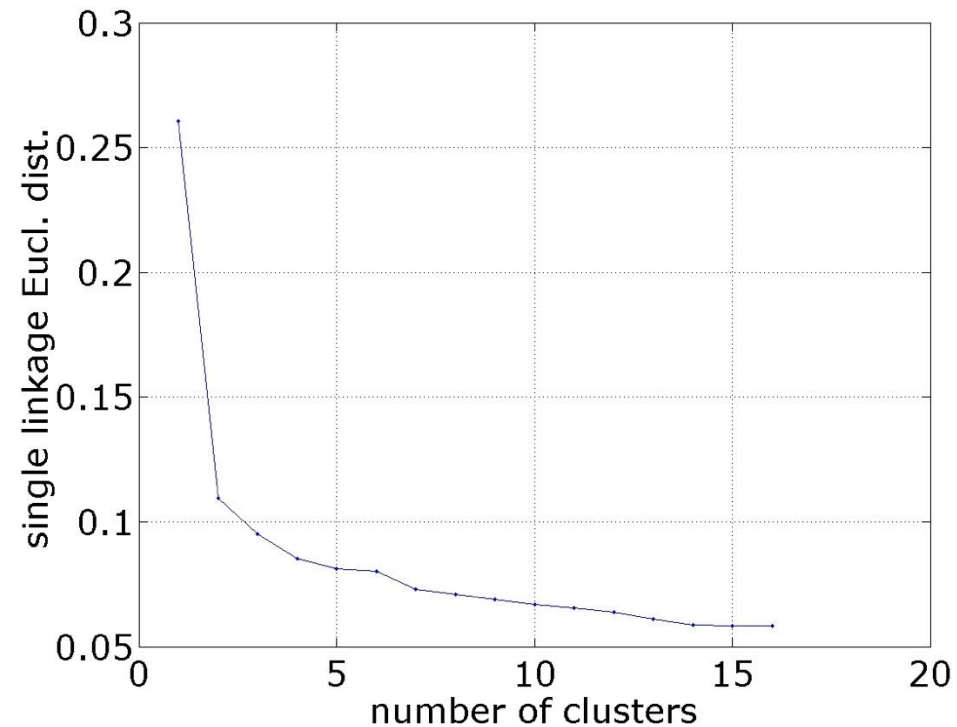
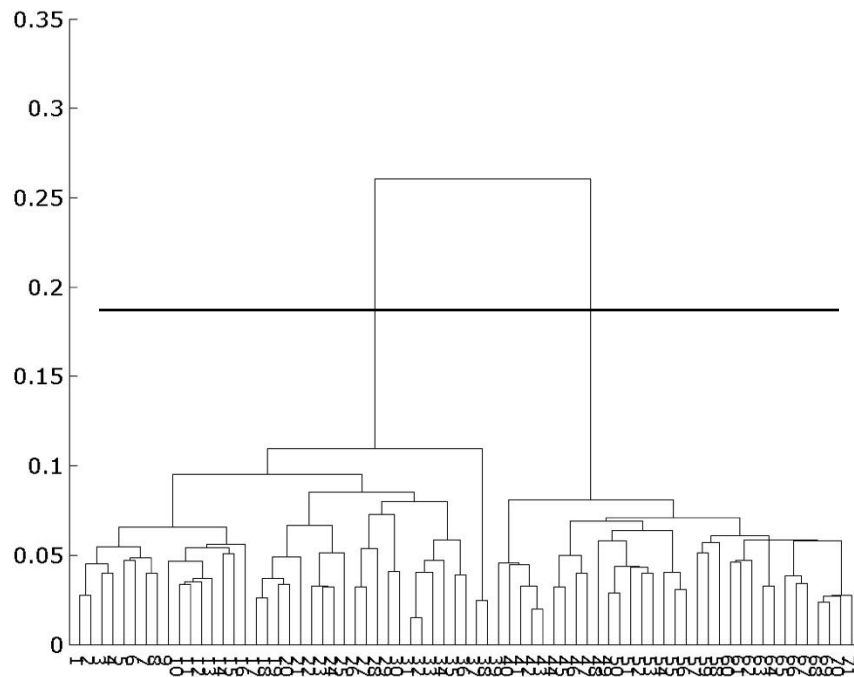
Fusion graph (4)

(Euclidean; single linkage)



Fusion graph (5)

(Euclidean; single linkage)



Hierarchical clustering (16)

- Advantages:
 - dendrogram gives overview of all possible clusterings
 - linkage type allows to find clusters of varying shapes (convex and non-convex)
 - different dissimilarity measures can be used
- Disadvantages:
 - computationally intensive:
 $O(n^2)$ in complexity and memory
 - clusterings limited to “hierarchical nestings”

Any other way of clustering?

- Hierarchical:



- ?????? E.g., *K*-means (sum-of-squares)

Any other way of clustering?

- Let's start with what we expect to see at the end after conducting a clustering!
-
- Basically, our goal is to **partition the data into some number K of clusters.**
- $K = ???$

Intuitively, what should a cluster look like?

- Well...a group of data points...
- whose **inter-point distances** are **small** compared with the **distances to points outside of the cluster**.



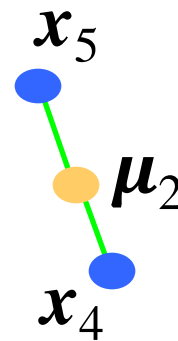
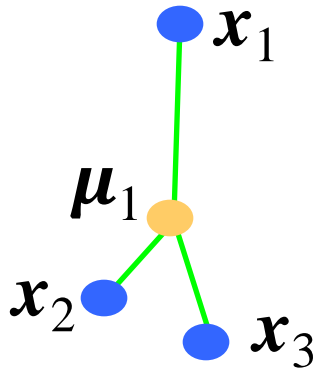
A notion to represent the k^{th} cluster

- $\mu_k, k=1, \dots, K$, a prototype associated with the k^{th} cluster.
- But what is μ_k ? A simple but sensible choice...



A notion to represent the k^{th} cluster

- μ_k , $k=1, \dots, K$, a prototype associated with the k^{th} cluster.
- But what is μ_k ? A simple but sensible choice...
- the centre of the cluster



Our goal is then...

- To find an **assignment of data points** to clusters, as well as a set of vectors μ_k , such that
- the **sum of the squares of the distances** of each data point to its closest μ_k , is a **minimum**.

K-means clustering,
a.k.a., sum-of-squares clustering

Our goal is then...

- To find an **assignment of data points** to clusters, as well as a set of vectors μ_k , such that
- the **sum of the squares of the distances** of each data point to its closest μ_k , is a **minimum**.
- But **where** to put μ_k ...

K-means clustering,
a.k.a., sum-of-squares clustering

K-means clustering

- 1. choose number of clusters (K)
- 2. position prototypes ($\mu_k, k=1, \dots, K$) **randomly**
- 3. assign samples to closest prototype
- 4. compute mean of samples assigned to same prototype → **new prototype position**

Repeat steps 3 and 4 as long as prototypes move

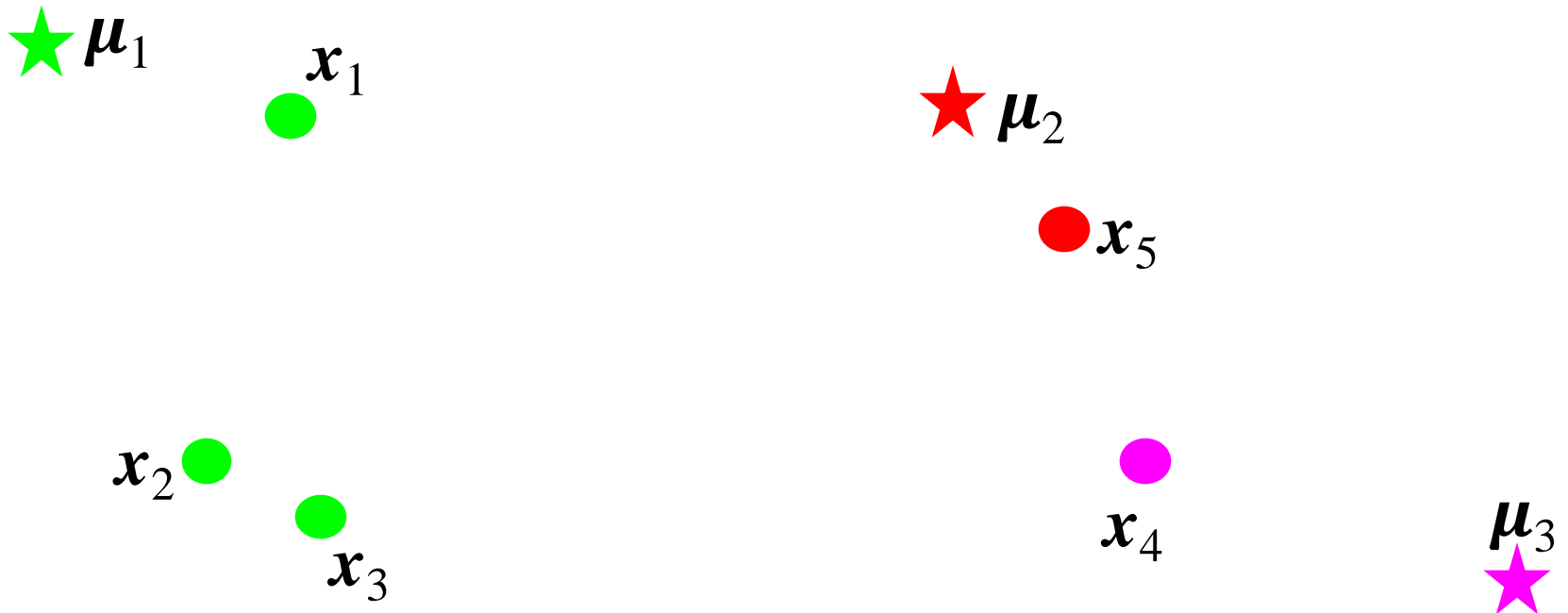
K-means clustering

- **Step 1:** Choose number of clusters/prototypes
- **Step 2:** Position prototypes randomly



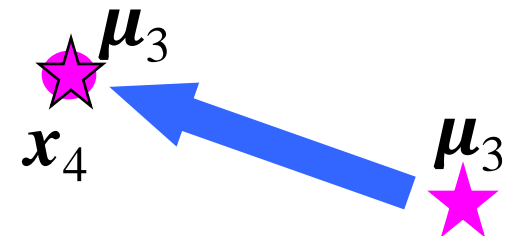
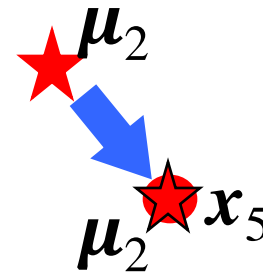
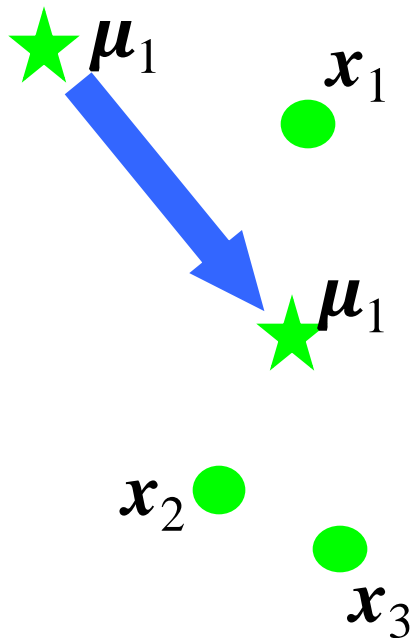
K-means clustering

- **Step 3:** Assign samples to closest prototype



K-means clustering

- **Step 4:** Compute mean of samples assigned to same prototype → **new prototype positions**



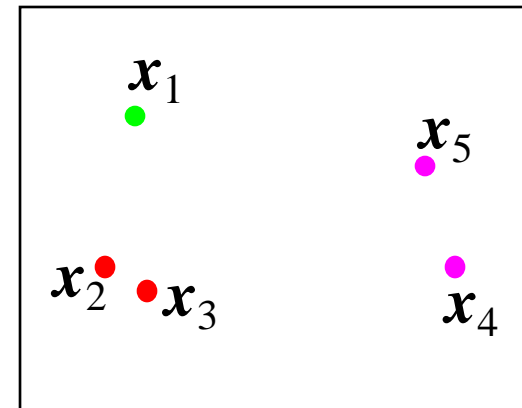
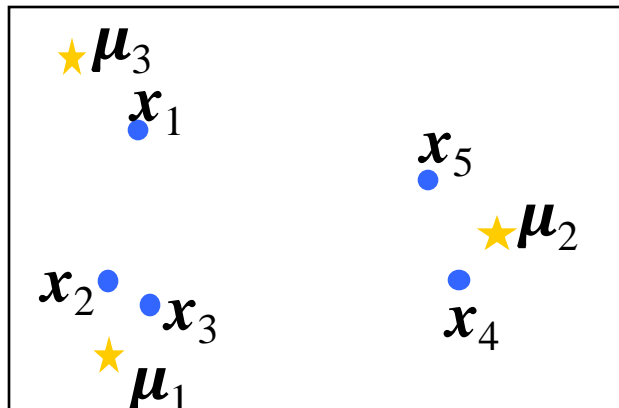
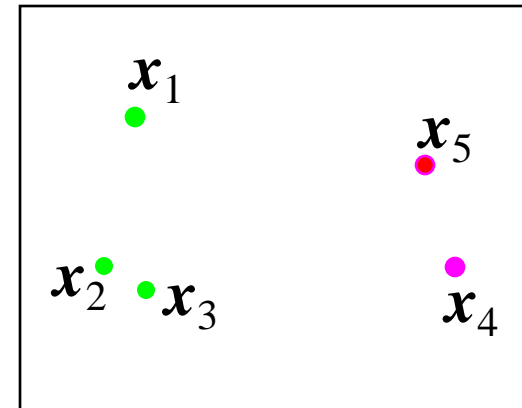
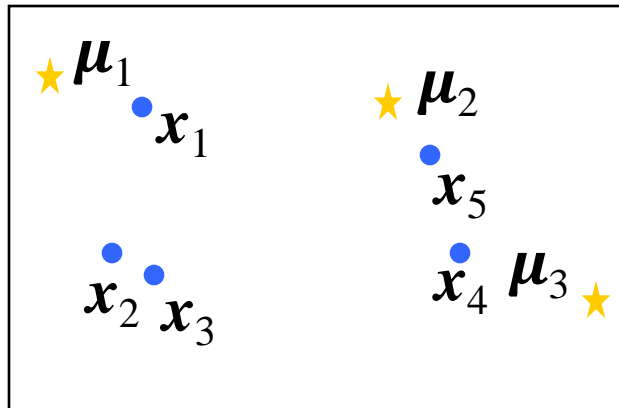
K-means clustering

- **Repeat** as long as prototype positions change:
 - **Step 3:** Assign samples
 - **Step 4:** Recompute prototype positions



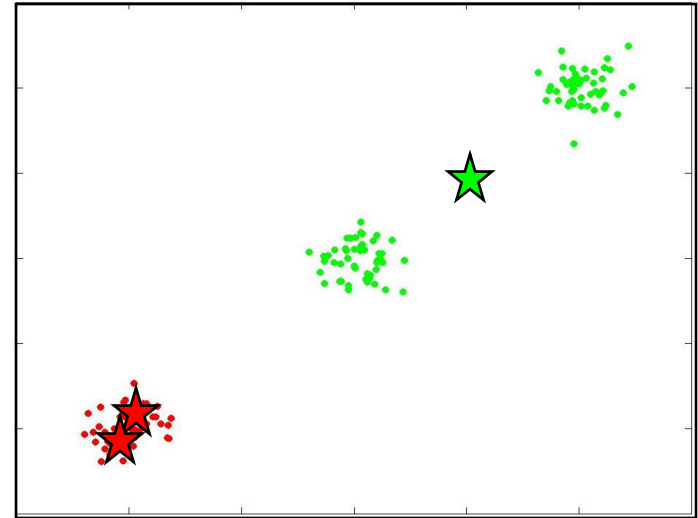
K-means problems

- Clustering depends on initialization



K-means problems

- Algorithm can get stuck in local minima
- Solution:
 - start from ***I* different random initialisations**;
 - **keep the best clustering** (lowest $\text{Tr}(S_W)$);
 - For high-dimensional data, many restarts are necessary (e.g. $I = 10000$)!

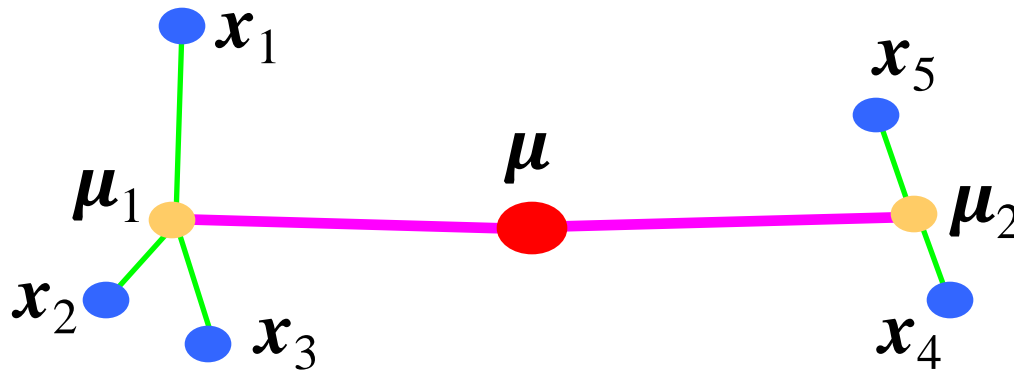


K-means clustering

- Recall from Week 4 (within and between group scatter):

$$S_W = \sum_{j=1}^K \frac{n_k}{N} \Sigma_k ,$$

$$S_B = \sum_{k=1}^K \frac{n_k}{N} (\mu_k - \mu)(\mu_k - \mu)^T, \quad \mu = \sum_{k=1}^K \frac{n_k}{N} \mu_k$$



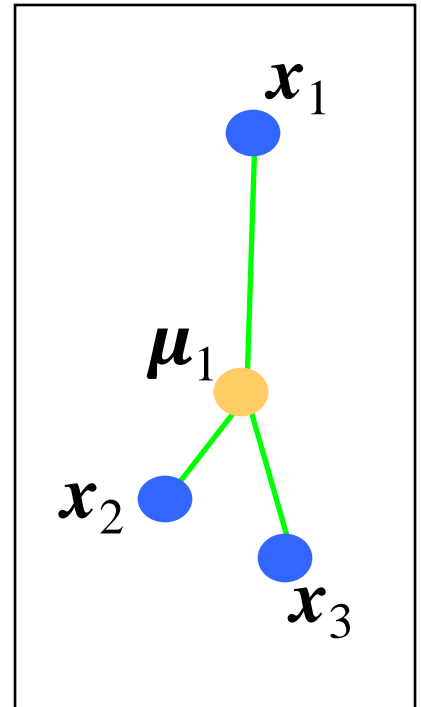
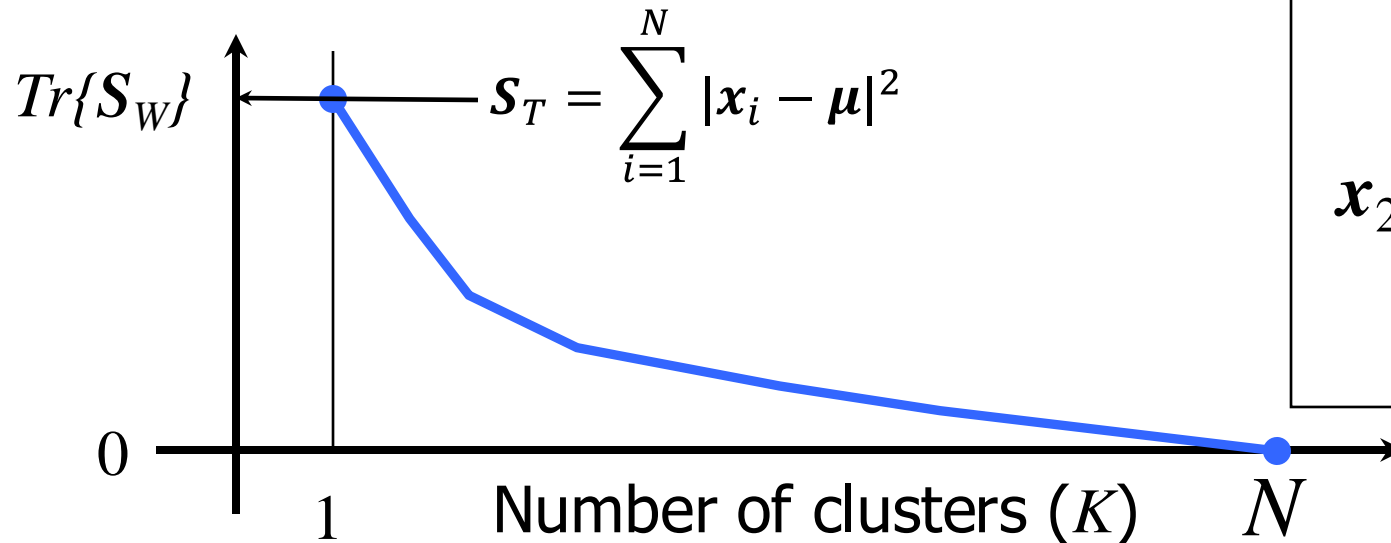
$$K = 2, N = 5, \\ n_1 = 3, n_2 = 2$$

K-means clustering

- Minimize:

$$Tr\{\mathbf{S}_W\} = \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^{n_k} |\mathbf{x}_i - \boldsymbol{\mu}_k|^2$$

(sum of cluster within-scatters)



Our goal is then... (Recall)

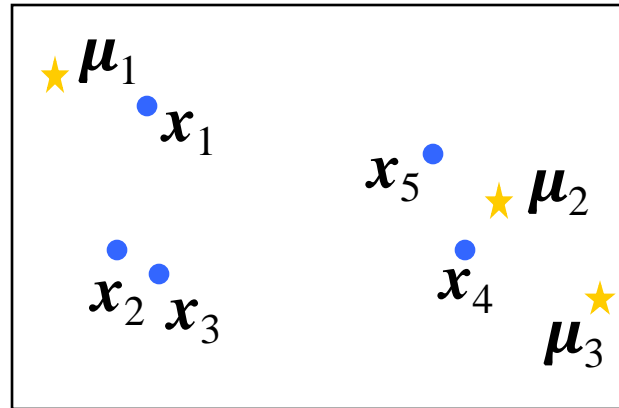
- To find an **assignment of data points** to clusters, as well as a set of vectors μ_k , such that
- the **sum of the squares of the distances** of each data point to its closest μ_k , is a **minimum**.

K-means clustering,
a.k.a., sum-of-squares clustering

K-means problems

10.5

- Some prototypes do not capture any data points (μ_3)



- Possible solution:
 - remove cluster and continue with $K - 1$ means
 - alternatively, split largest cluster into two or add a random cluster to continue with K means

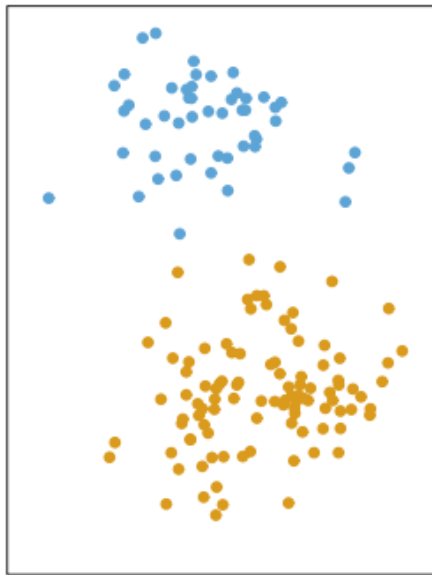
K-means clustering

- Choose a different K :

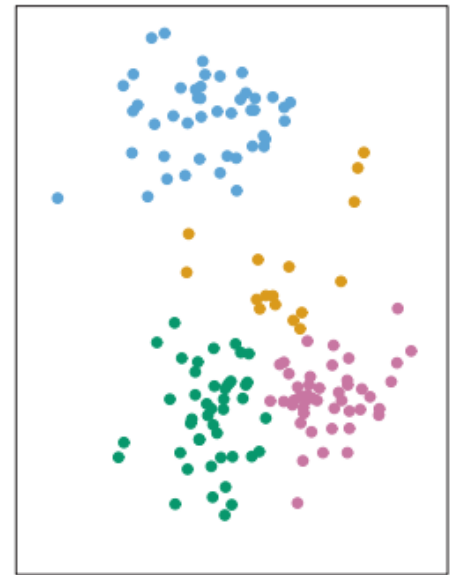
$K = 2$



$K = 3$



$K = 4$



K-means summary

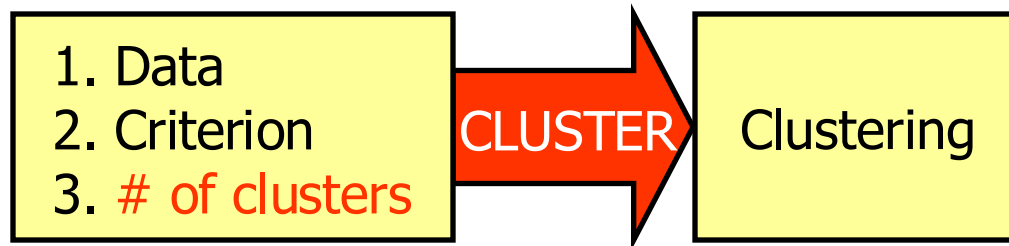
- Disadvantages:
 - Finds only convex clusters (“round shapes”)
 - Sensitive to initialization
 - Can get stuck in local minima
- Advantages:
 - Simple
 - Fast

Clustering

- Hierarchical:

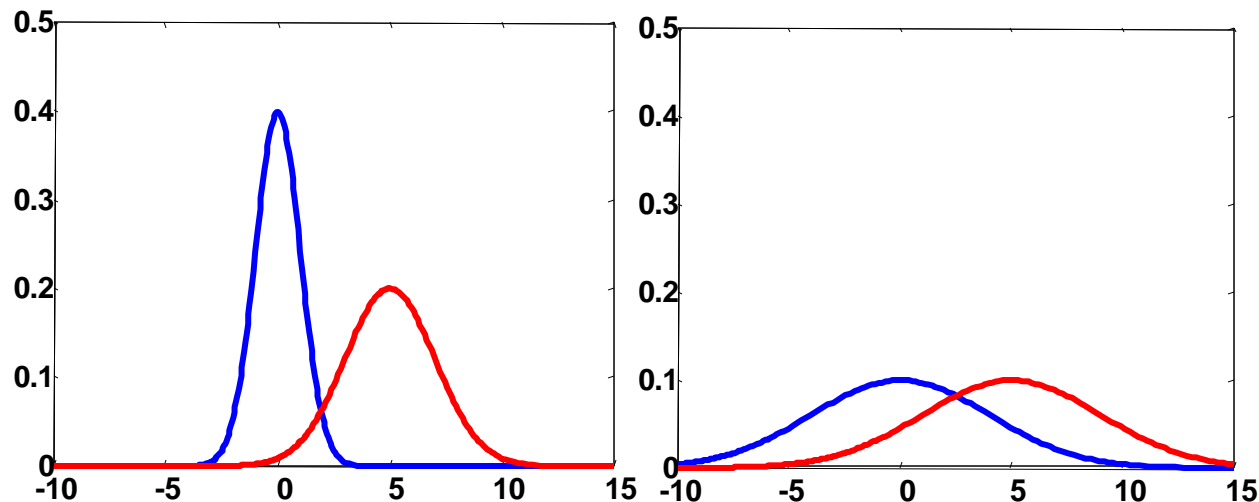


- Sum-of-squares:



Validity measures

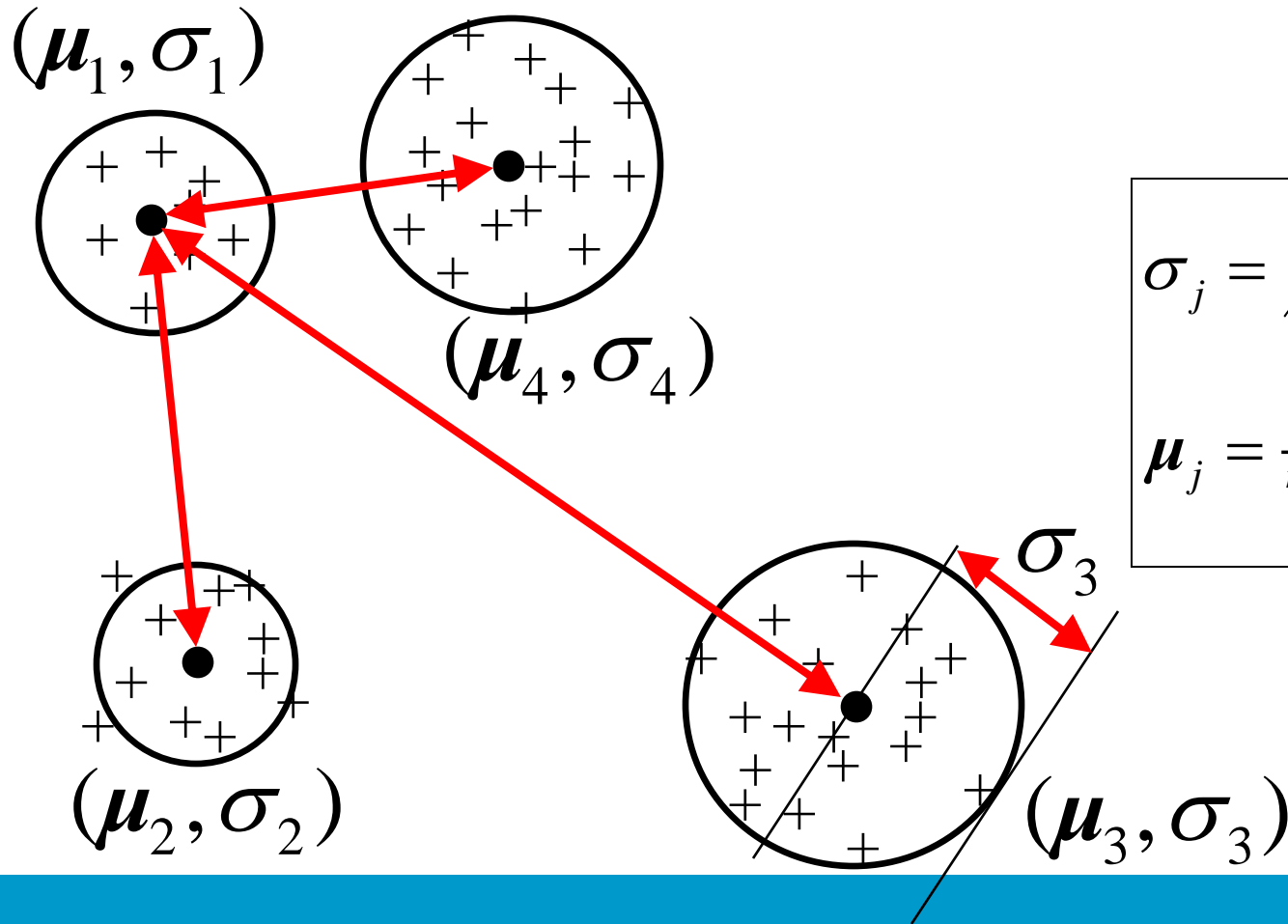
- Many are based on within and between group scatter
- The **larger the between** group scatter and **the smaller the within** group scatter, **the better**
- Example: Davies-Bouldin index



Davies-Bouldin index

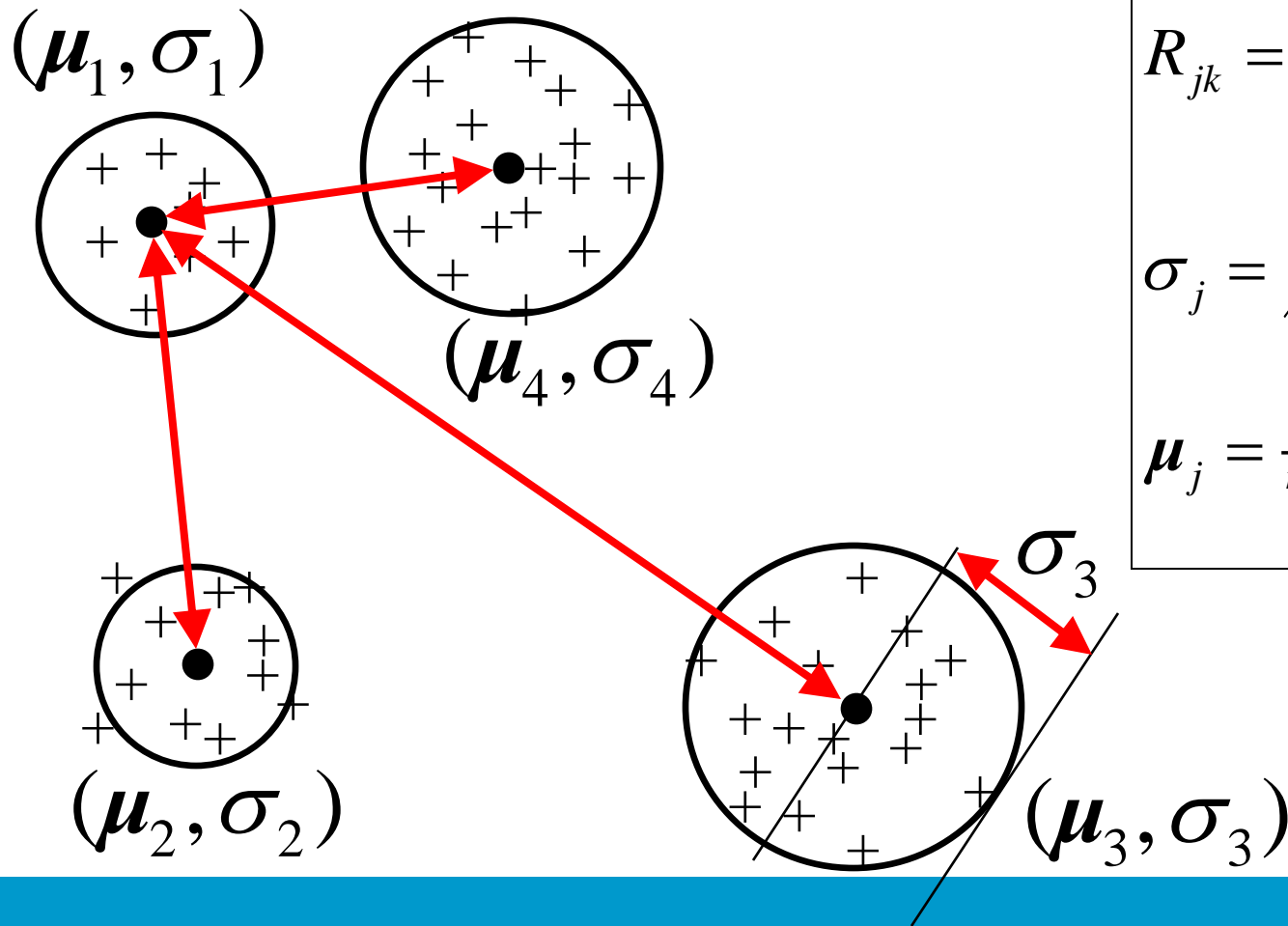
- **Assumption: clusters are spherical**
- For a good clustering, it should hold that:
 - objects are compactly organized within a cluster
 - clusters are far apart
- D.L. Davies and D.W. Bouldin, IEEE Transactions on Pattern Analysis and Machine Intelligence 1, pp. 224-227, 1979

Davies-Bouldin index (2)



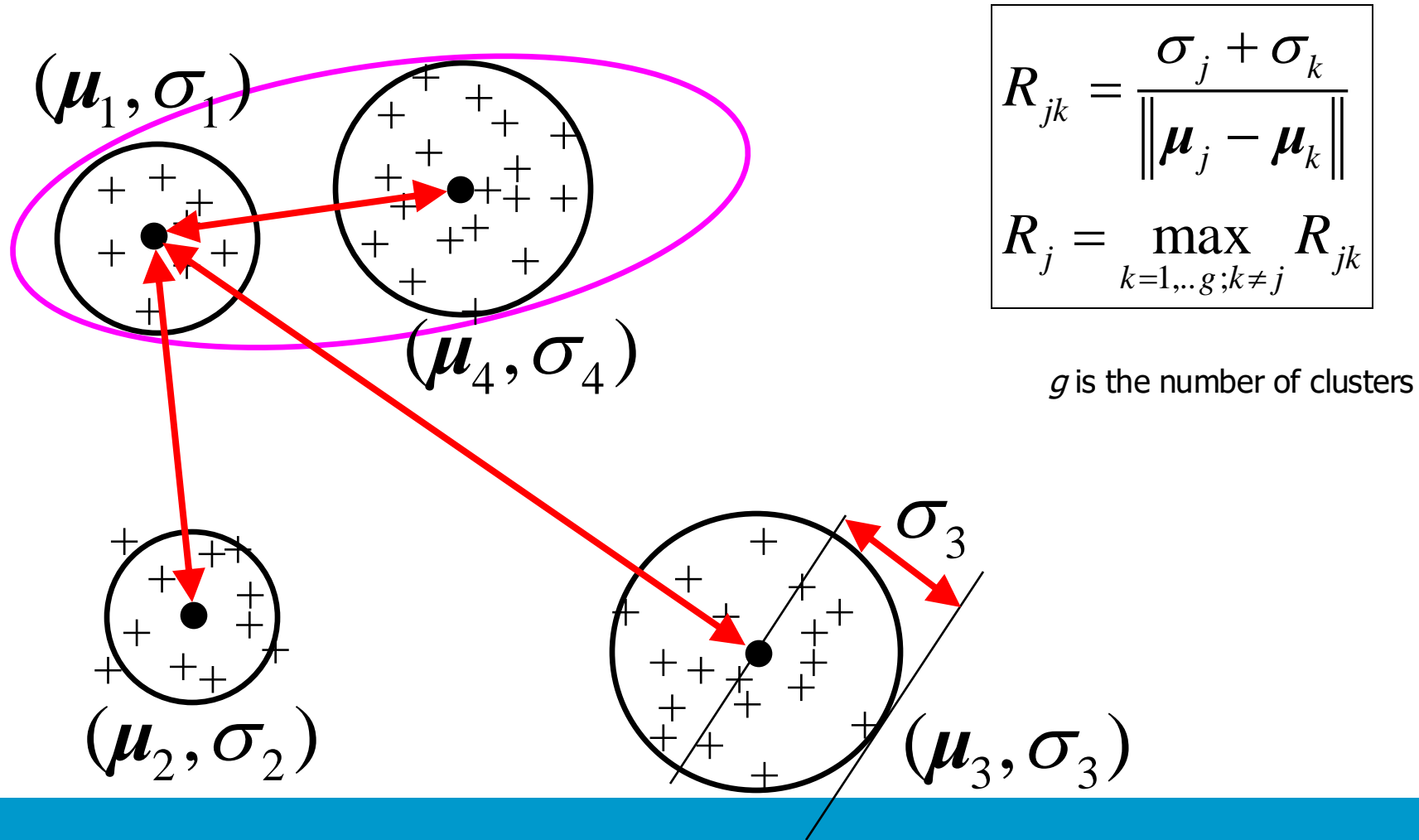
$$\sigma_j = \sqrt{\frac{1}{n_j} \sum_{x_i \in C_j} \|x_i - \mu_j\|^2}$$
$$\mu_j = \frac{1}{n_j} \sum_{x_i \in C_j} x_i$$

Davies-Bouldin index (3)



$$R_{jk} = \frac{\sigma_j + \sigma_k}{\|\mu_j - \mu_k\|}$$
$$\sigma_j = \sqrt{\frac{1}{n_j} \sum_{x_i \in C_j} \|x_i - \mu_j\|^2}$$
$$\mu_j = \frac{1}{n_j} \sum_{x_i \in C_j} x_i$$

Davies-Bouldin index (3)



Davies-Bouldin index (4)

$$R_{jk} = \frac{\sigma_j + \sigma_k}{\|\mu_j - \mu_k\|}$$

$$R_j = \max_{k=1, \dots, g; k \neq j} R_{jk}$$

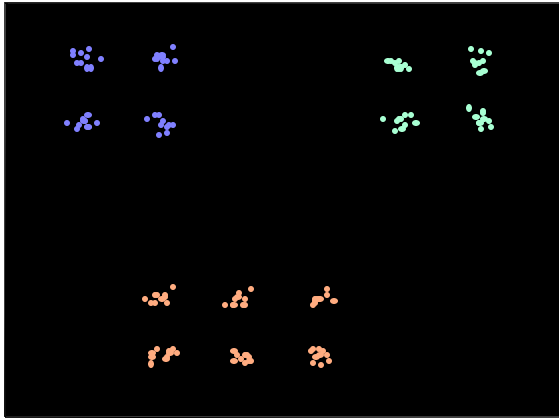
$$I_{DB} = \frac{1}{g} \sum_{j=1}^g R_j$$

Paired cluster criterion

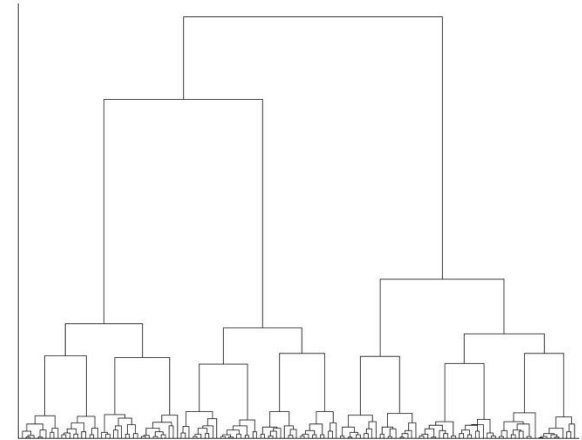
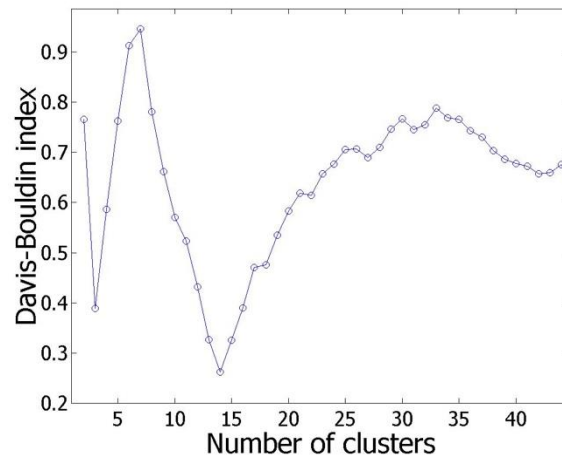
Worst-case value per cluster

Average worst-case

Davies-Bouldin index (5)



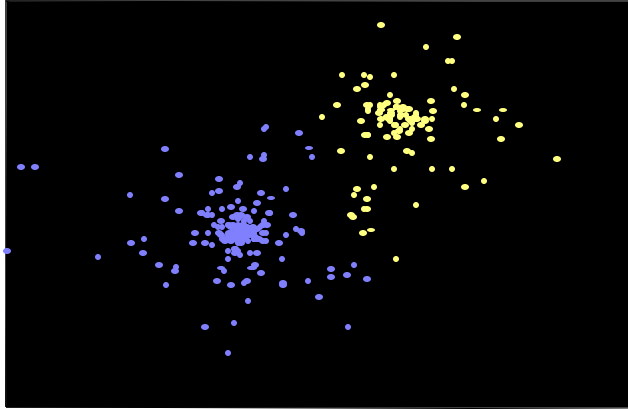
Dataset



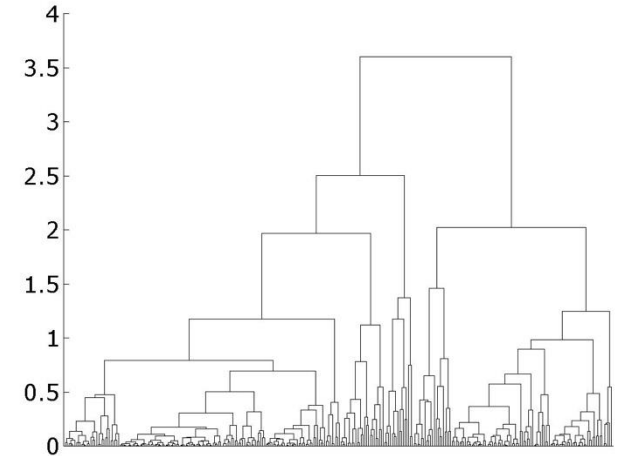
Complete link

Davis Bouldin:
3 or 14 clusters

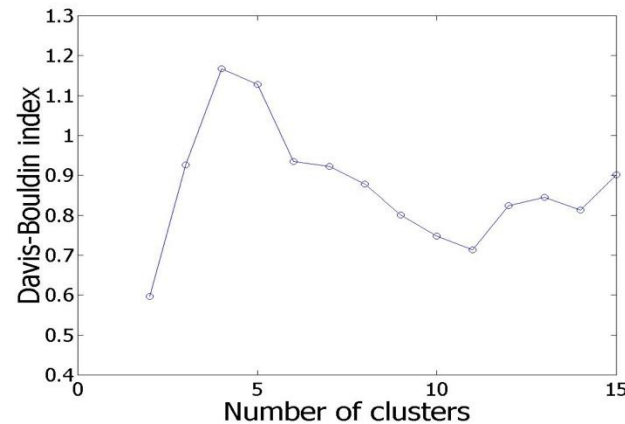
Davies-Bouldin index (5)



Dataset

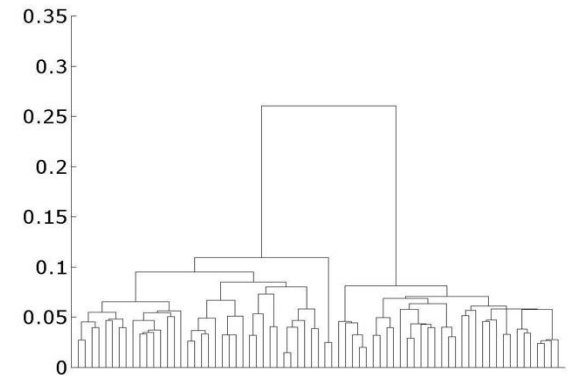
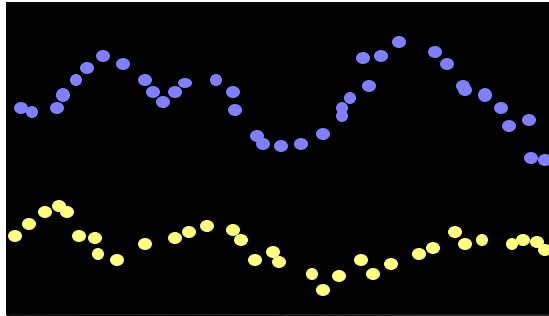


Complete link

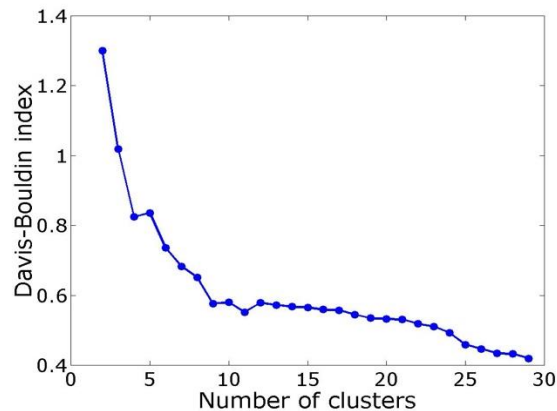


Davis Bouldin:
2 or 11 clusters

Davies-Bouldin index (6)



Davis-Bouldin:



Single link

