

Learning Systems (DT8008)

# Unsupervised Learning: Clustering

Dr. Mohamed-Rafik Bouguelia  
[mohamed-rafik.bouguelia@hh.se](mailto:mohamed-rafik.bouguelia@hh.se)

Halmstad University

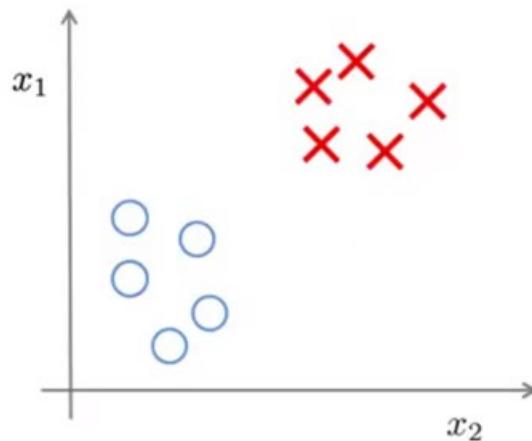
# **Classification ≠ Clustering**

Supervised

Unsupervised

# Supervised vs. Unsupervised Learning

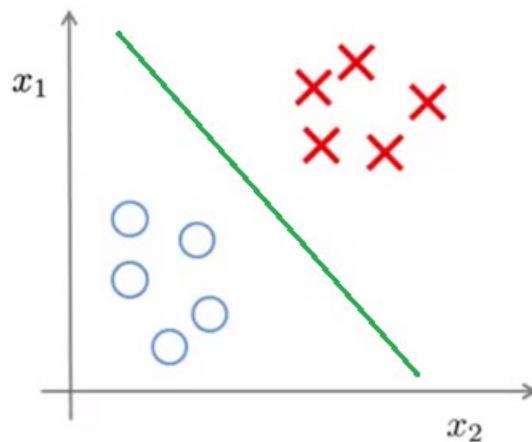
## Supervised learning



**Training set:**  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$

# Supervised vs. Unsupervised Learning

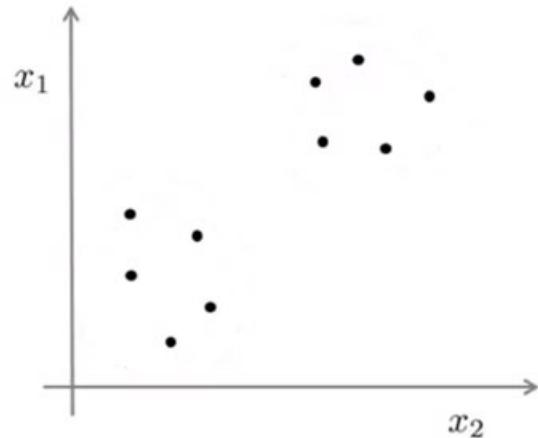
## Supervised learning



**Training set:**  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$

# Supervised vs. Unsupervised Learning

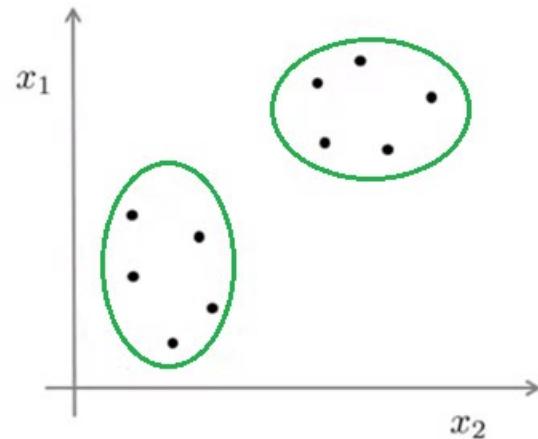
## Unsupervised learning



Training set:  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$

# Supervised vs. Unsupervised Learning

## Unsupervised learning



Training set:  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$

# Supervised vs. Unsupervised Learning

Of the following examples, which would you address using an unsupervised learning algorithm? (Check all that apply.)

- Given email labeled as spam/not spam, learn a spam filter.
- Given a set of news articles found on the web, group them into set of articles about the same story.
- Given a database of customer data, automatically discover market segments and group customers into different market segments.
- Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not.

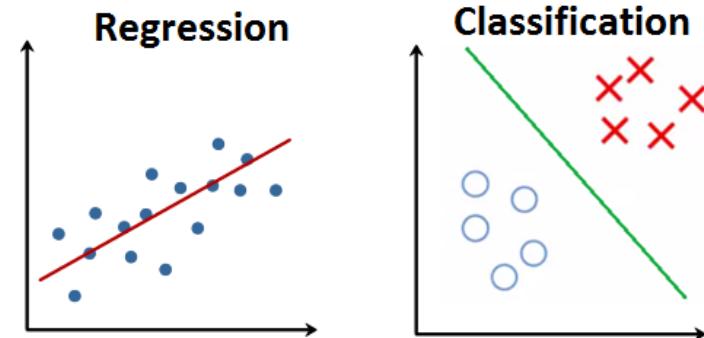
# Supervised vs. Unsupervised Learning

Of the following examples, which would you address using an unsupervised learning algorithm? (Check all that apply.)

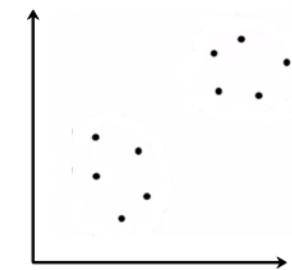
- Given email labeled as spam/not spam, learn a spam filter.
- Given a set of news articles found on the web, group them into set of articles about the same story.
- Given a database of customer data, automatically discover market segments and group customers into different market segments.
- Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not.

# Supervised vs. Unsupervised Learning

- Supervised learning uses labeled data pairs  $(x^{(i)}, y^{(i)})$  to learn some function  $h: X \rightarrow y$ 
  - e.g. Regression, Classification



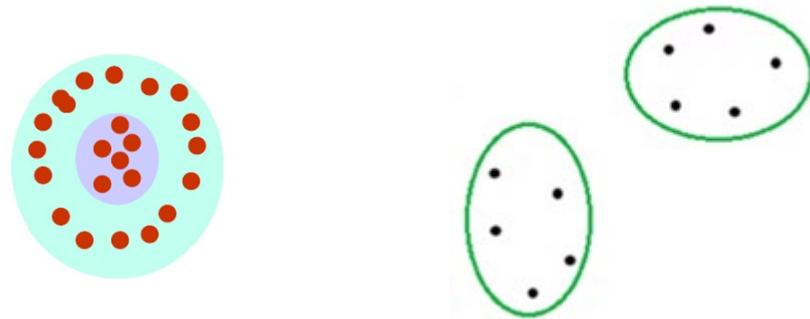
- What if we don't have labels?
- Unsupervised learning → without labels (without supervision from human)
- Semi-supervised learning → few labelled instances and lot of unlabeled instances.
- **Clustering** is the unsupervised grouping of data-points.
  - e.g. for exploring the data or for knowledge discovery ...



# Supervised vs. Unsupervised Learning

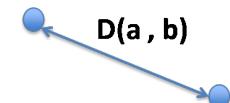
- Example of unsupervised learning tasks:
  - Clustering
  - Outlier and Anomaly detection
  - Dimensionality Reduction (e.g. with PCA)
  - Unsupervised feature learning (e.g. with Autoencoders).
- We want to explore the data to find some intrinsic structures in them.
  - **Clustering** is when the clusters are **not known**
  - If the clusters were **known** (i.e. they are known classes) and the problem was to place a new instance into the proper class, then this is **classification** (not clustering)
- It is often easier to obtain unlabeled data than labeled data (which requires human intervention).

# Clustering

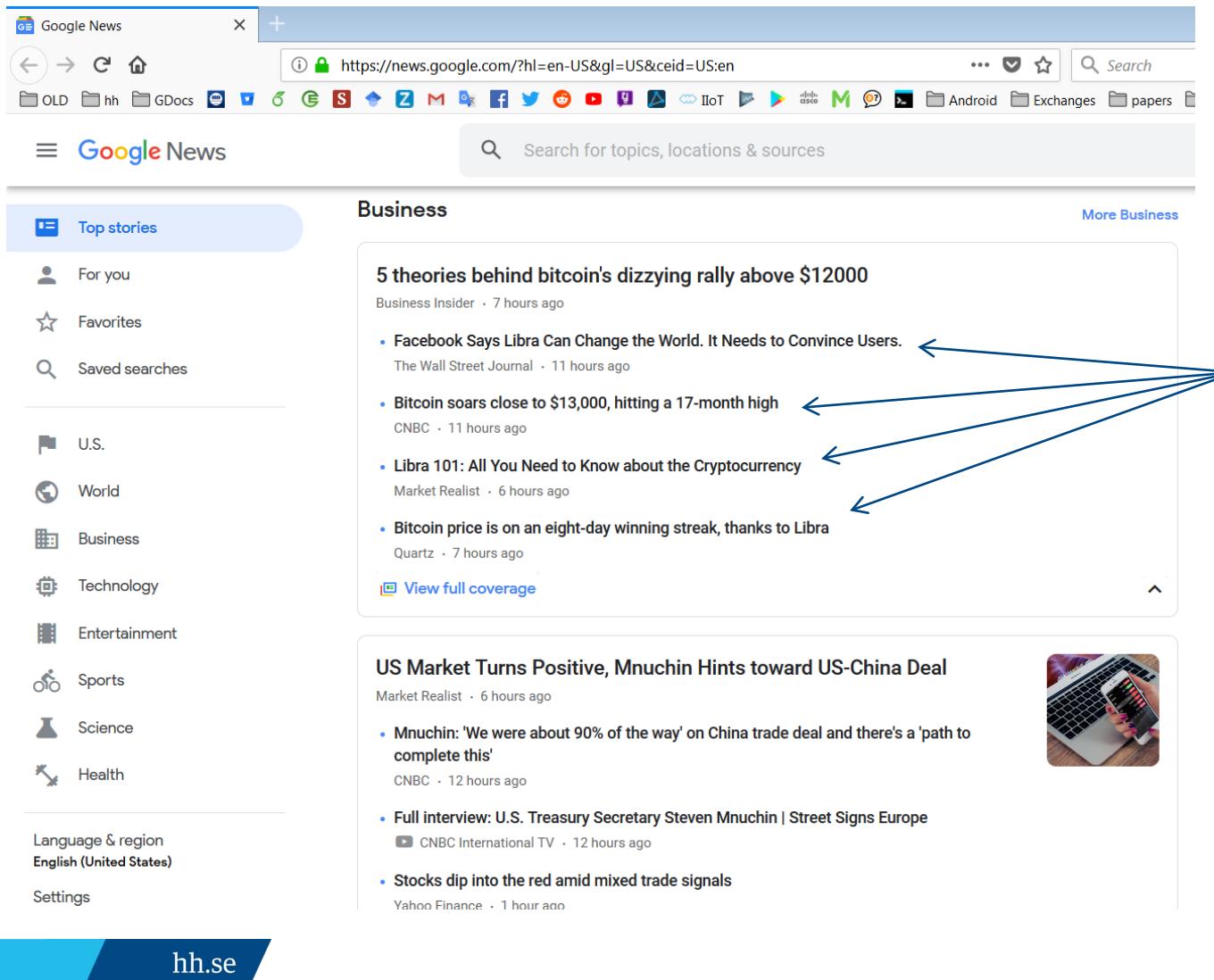


- Clustering is grouping data points such that
  - Instances within each cluster are similar to each other (**minimize intra-clusters distance**)
  - Instances from different clusters are dissimilar (**maximize inter-clusters distance**)
- Goal is to discover interesting structures (unknown subgroups in the data).
- The dissimilarity is defined using a distance measure
  - e.g. The Euclidean distance
- Examples:
  - Groups of breast cancer patients grouped by their gene expression measurements
  - Groups of shoppers characterized by their browsing and purchase histories
  - Movies grouped by the ratings assigned by movie viewers.
  - ... etc.

$$D(a, b) = \sqrt{\sum_{i=1}^n (b_i - a_i)^2}$$



# Example application of clustering



The screenshot shows a Google News interface on a computer. The main content area is titled 'Business' and displays a cluster of news stories about Bitcoin and Libra. The stories include:

- 5 theories behind bitcoin's dizzying rally above \$12000 (Business Insider, 7 hours ago)
- Facebook Says Libra Can Change the World. It Needs to Convince Users. (The Wall Street Journal, 11 hours ago)
- Bitcoin soars close to \$13,000, hitting a 17-month high (CNBC, 11 hours ago)
- Libra 101: All You Need to Know about the Cryptocurrency (Market Realist, 6 hours ago)
- Bitcoin price is on an eight-day winning streak, thanks to Libra (Quartz, 7 hours ago)

Blue arrows point from the text 'Automatically grouping together the stories (news articles on the Web) that talk about the same topic.' to the cluster of stories in the 'Business' section.

**Top stories**

- For you
- Favorites
- Saved searches

U.S. World Business Technology Entertainment Sports Science Health

Language & region English (United States) Settings

hh.se

Automatically grouping together the stories (news articles on the Web) that talk about the same topic.

# Clustering is “subjective”

- **Unsupervised** learning (e.g clustering) is more subjective than **supervised** learning.

What is a natural grouping among these objects?



Clustering is subjective

A cluster containing Marge, Homer, Bart, and Mole.	A cluster containing Moe, Marge, Marge, and Bart.	A cluster containing Marge, Marge, Marge, and Lisa.	A cluster containing Homer, Mole, Moe, and Marge.
Simpson's Family	School Employees	Females	Males

One possible clustering (with two clusters).

Another possible clustering (with two clusters)

# Clustering is “subjective”

- **Unsupervised** learning (e.g clustering) is more subjective than **supervised** learning.

People can be clustered by different criteria

User ID	Country	Age	Gender	Blood	Heartbeat	Weight	Height	Sports	Income	Profession
1										
2										
3										
4	China	young	male			healthy			poor	
5										
6		old								
7										
8			female							
9										
10	US	young	male			unhealthy			rich	
11										
12		old								
13			female							
14						healthy				
15									poor	

# Distance metrics

# Distance functions

Note: a “*distance*” (i.e. dissimilarity) is the inverse of “*similarity*”.

- Let  $x \in R^d$  and  $z \in R^d$  be two data-points (vectors of dimension  $d$ ).
  - Minkowski distance (a general distance function) :

$$\|x - z\|_p = \left[ \sum_{i=1}^d (x_i - z_i)^p \right]^{\frac{1}{p}}$$

- Euclidean distance (i.e. when  $p = 2$ ) :

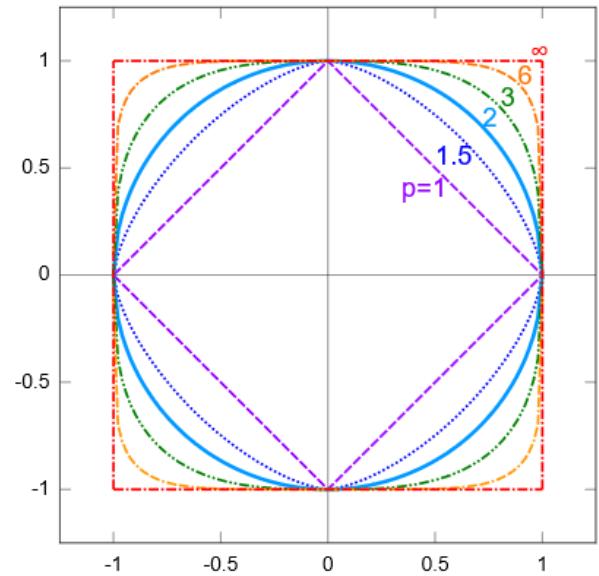
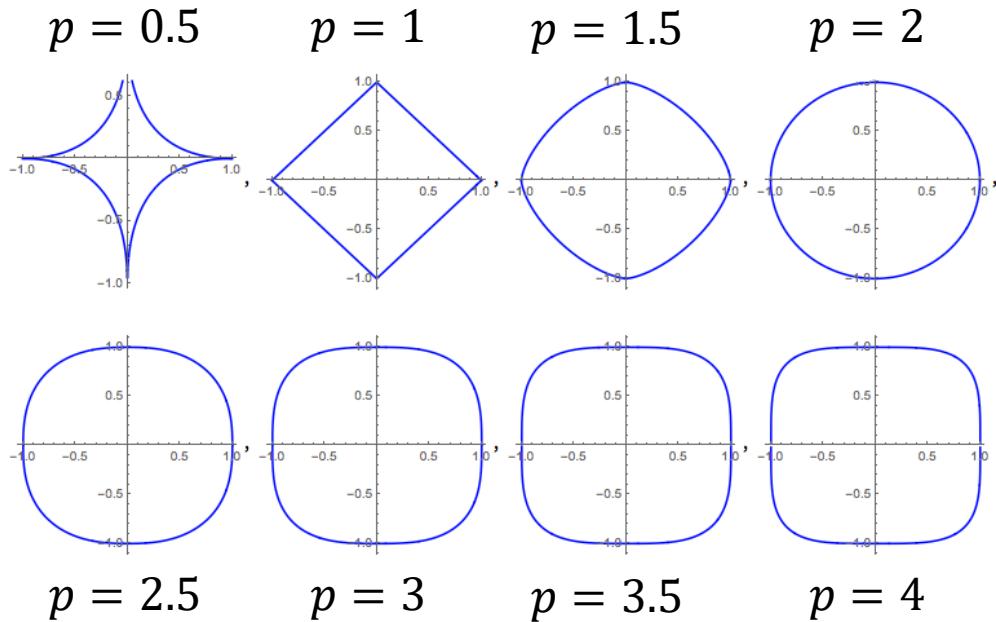
$$\|x - z\|_2 = \sqrt{\sum_{i=1}^d (x_i - z_i)^2}$$

- Manhattan distance (i.e. when  $p = 1$ ) :

$$\|x - z\|_1 = \sum_{i=1}^d |x_i - z_i|$$

# Norm of a vector

- Consider the vector  $v = x - z$ . The distance  $\|x - z\|_p$  is just the  $p$ -norm of the vector  $v$ .



All the points on the blue curve have the same  $p$ -norm  $\|v\|_p$  (Minkowski distance from  $v$  to the origin).

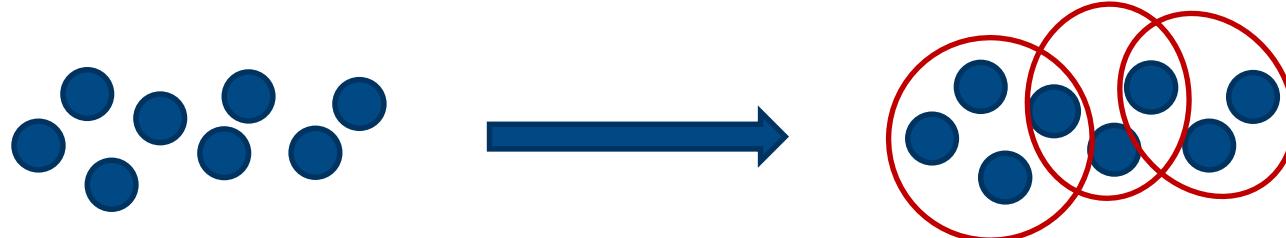
# **Some types of clustering techniques**

# Types of clustering techniques

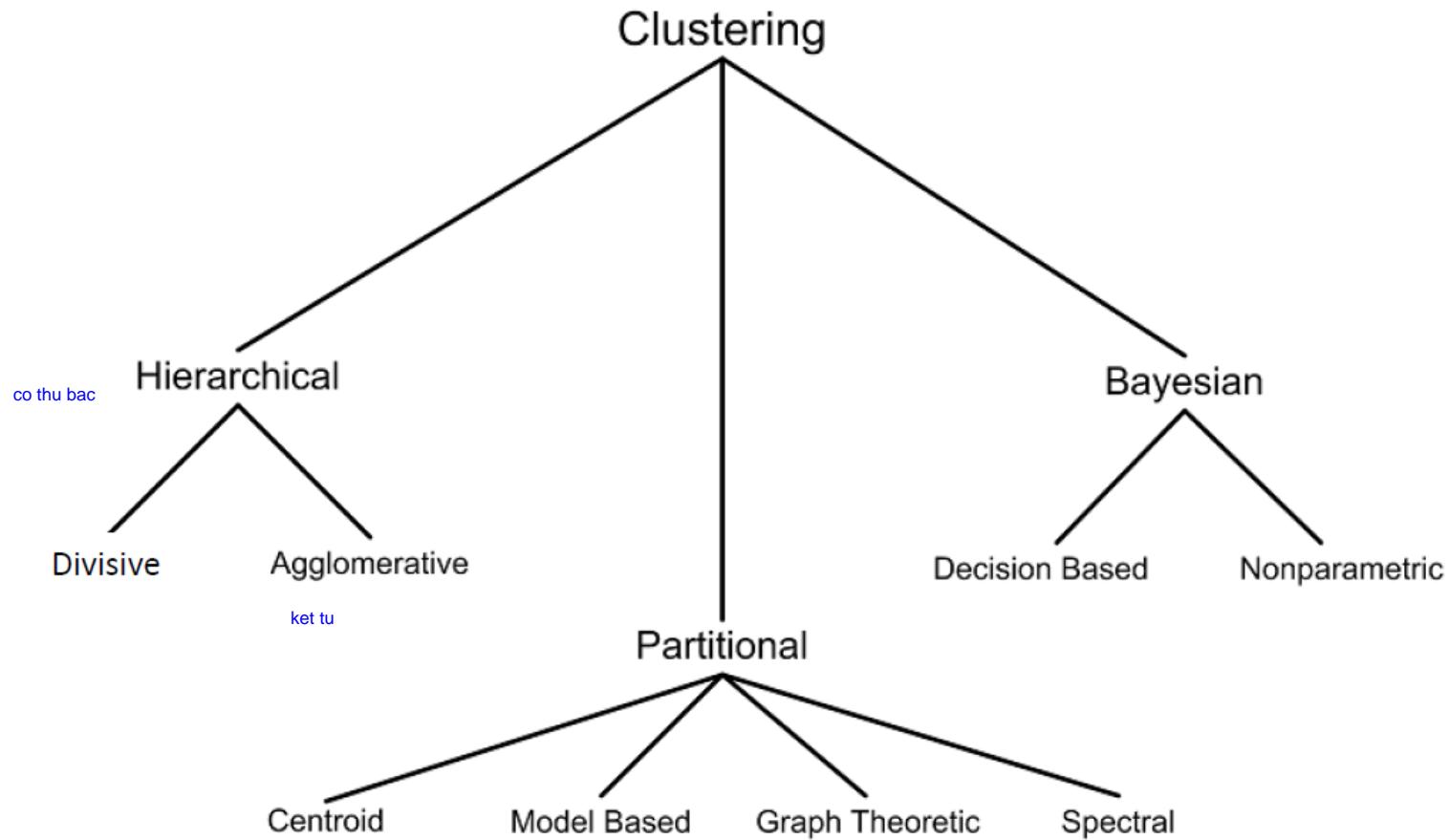
- **Hard clustering** (focus of this lecture): each data belongs to only one cluster.
  - More commonly used.



- **Soft clustering**: a data can belong to more than one cluster.
  - Example: some documents belongs to multiple topics
  - Membership to clusters is modeled as a probability distribution



# Types of clustering techniques



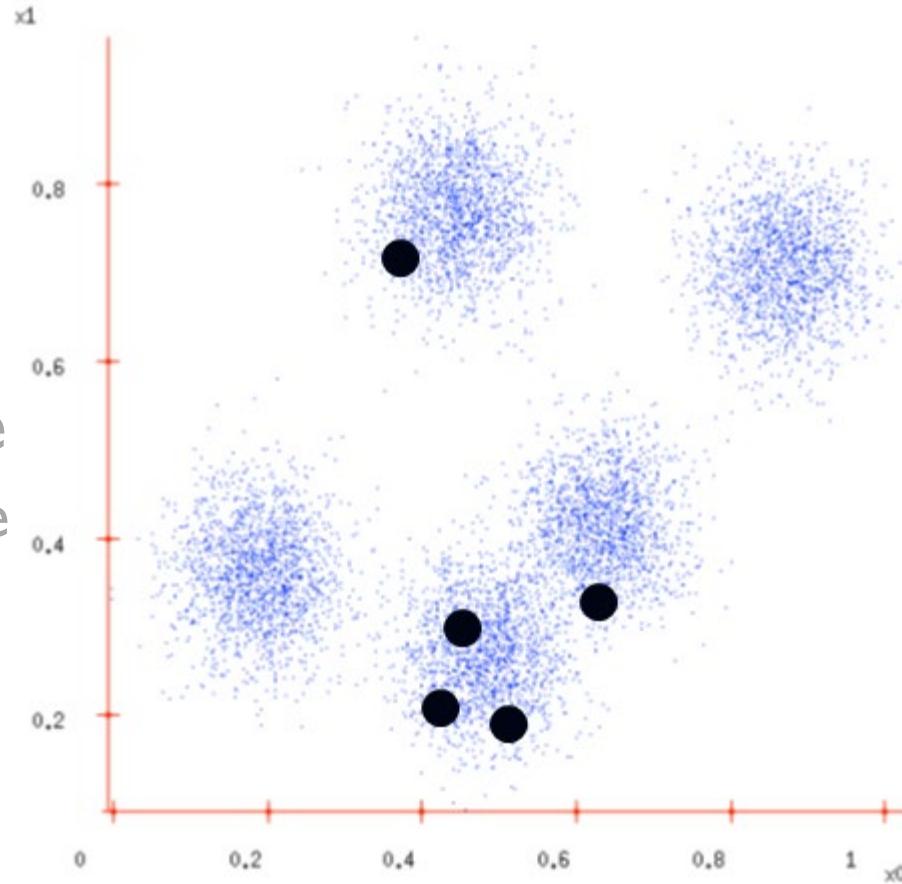
# Types of clustering techniques

- **Partitional** algorithms typically determine all clusters at once, but can also be used as divisive algorithms in the hierarchical clustering.
- **Hierarchical** algorithms find successive clusters using previously established clusters. These algorithms can be either **agglomerative** (“*bottom-up*”) or **divisive** (“*top-down*”):
  - 1 **Agglomerative algorithms** begin with each element as a separate cluster and merge them into successively larger clusters;
  - 2 **Divisive algorithms** begin with the whole set and proceed to divide it into successively smaller clusters.

# **The K-means clustering method**

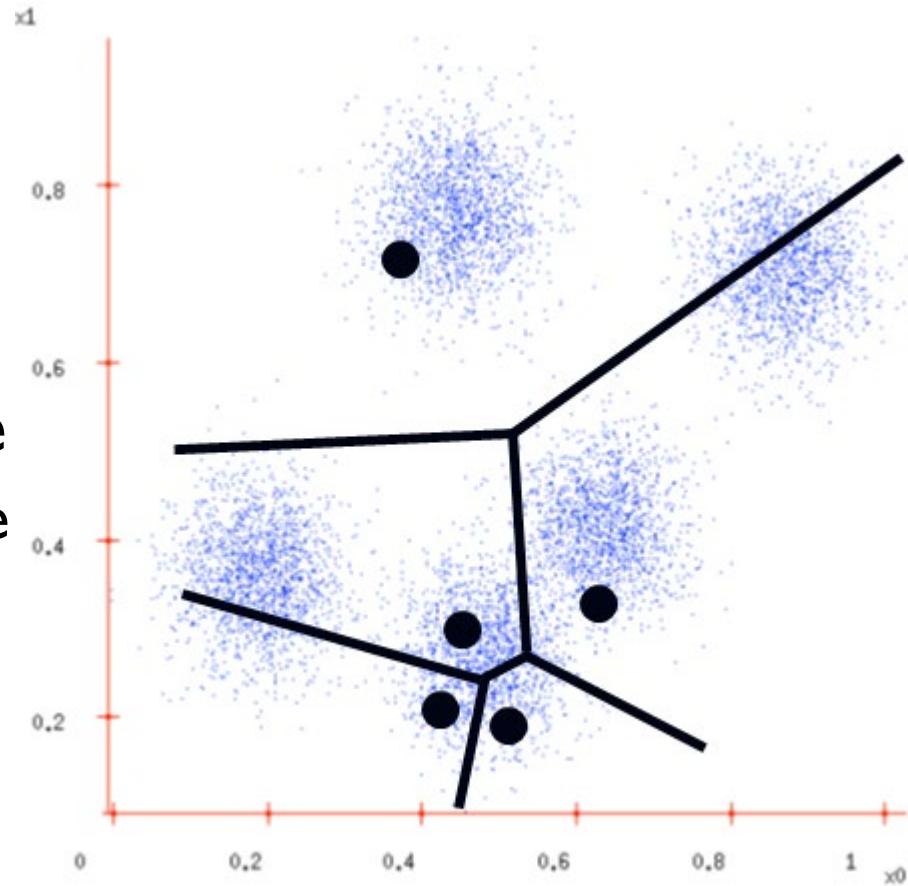
# K-Means (algorithm)

- Input: data, K
  - Randomly choose K cluster centers (centroids).
  - Loop until convergence
    - Assign each point to the cluster of the closest centroid.
    - Re-estimate the cluster centroids based on the data assigned to each.



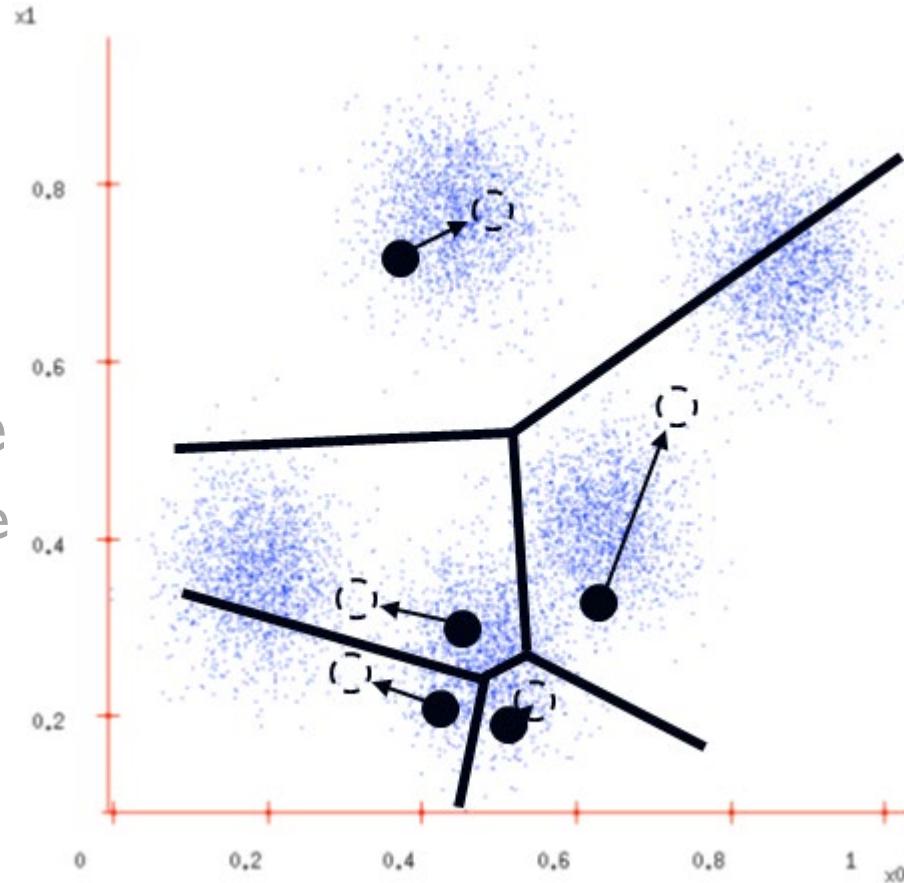
# K-Means (algorithm)

- Input: data, K
  - Randomly choose K cluster centers (centroids).
  - Loop until convergence
    - Assign each point to the cluster of the closest centroid.
    - Re-estimate the cluster centroids based on the data assigned to each.



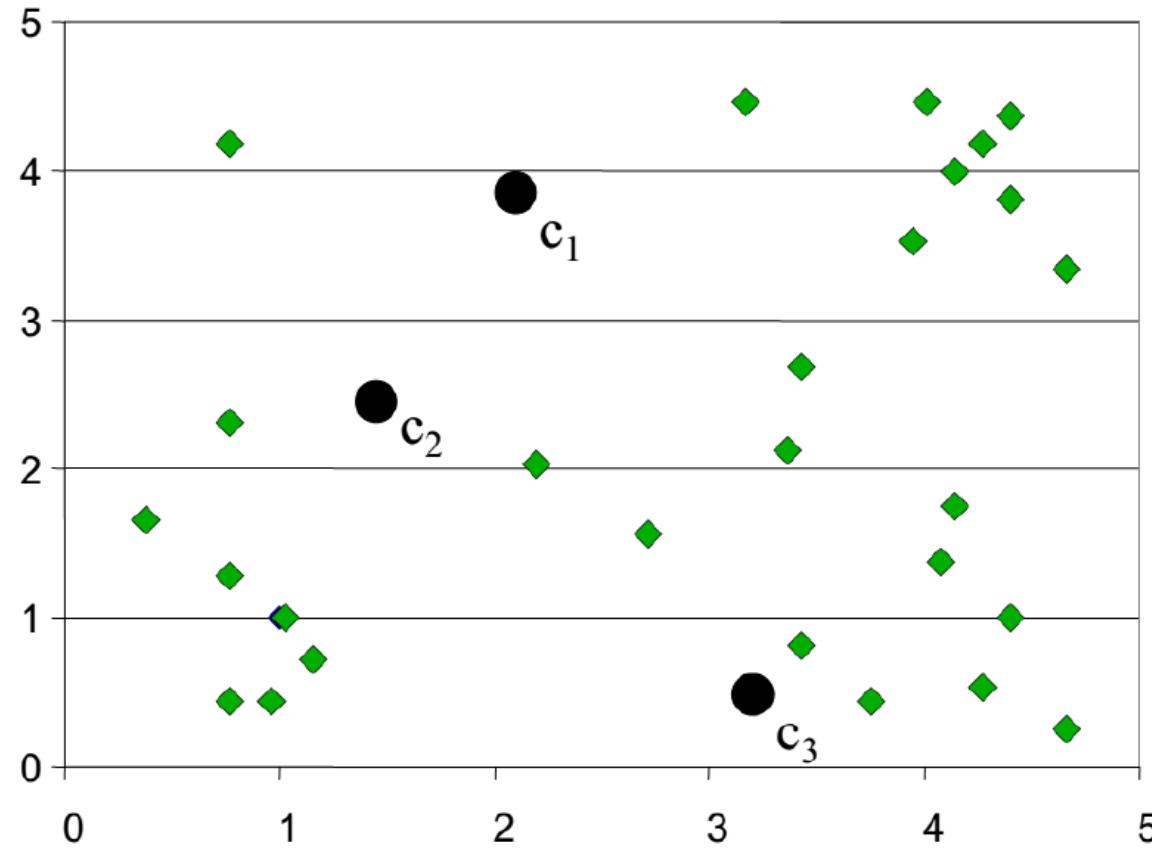
# K-Means (algorithm)

- Input: data, K
  - Randomly choose K cluster centers (centroids).
  - Loop until convergence
    - Assign each point to the cluster of the closest centroid.
    - Re-estimate the cluster centroids based on the data assigned to each.



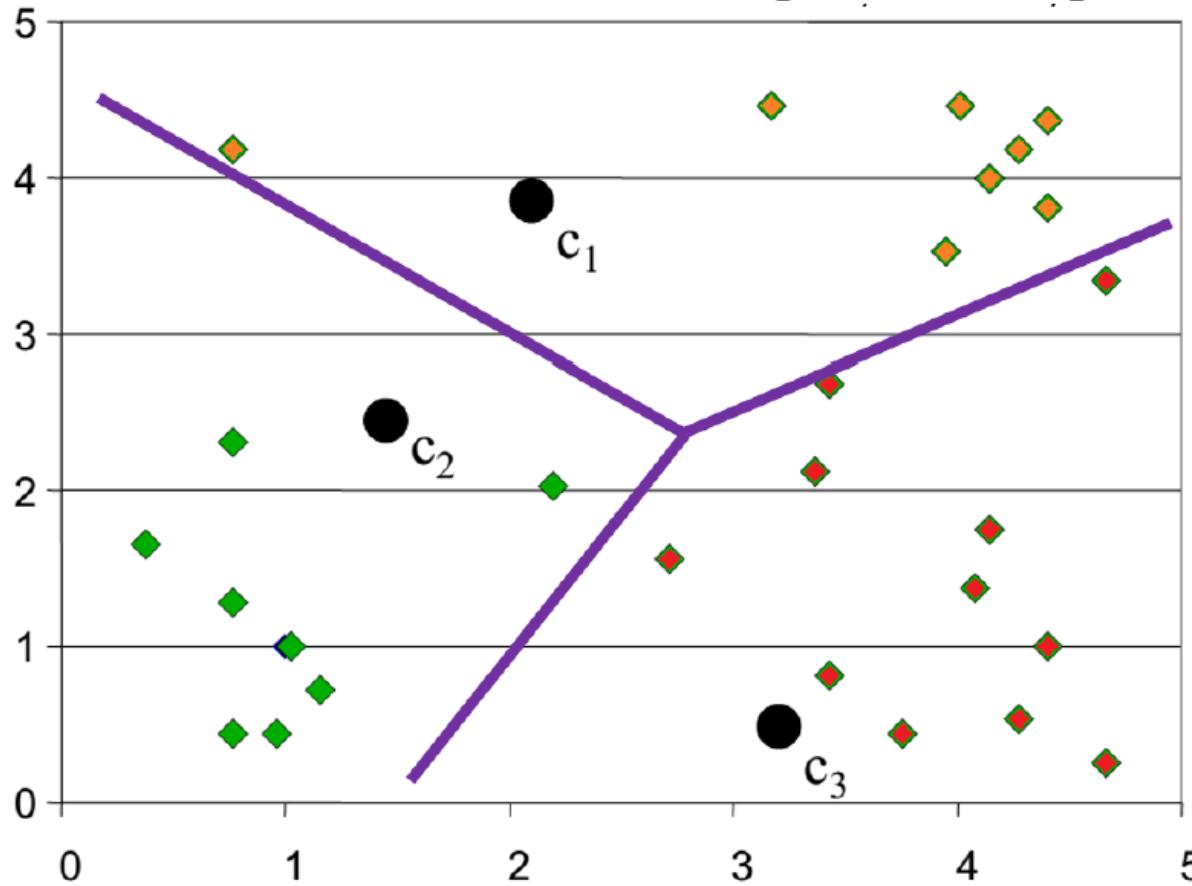
# K-Means (example)

Randomly initialize the cluster centers (synaptic weights)



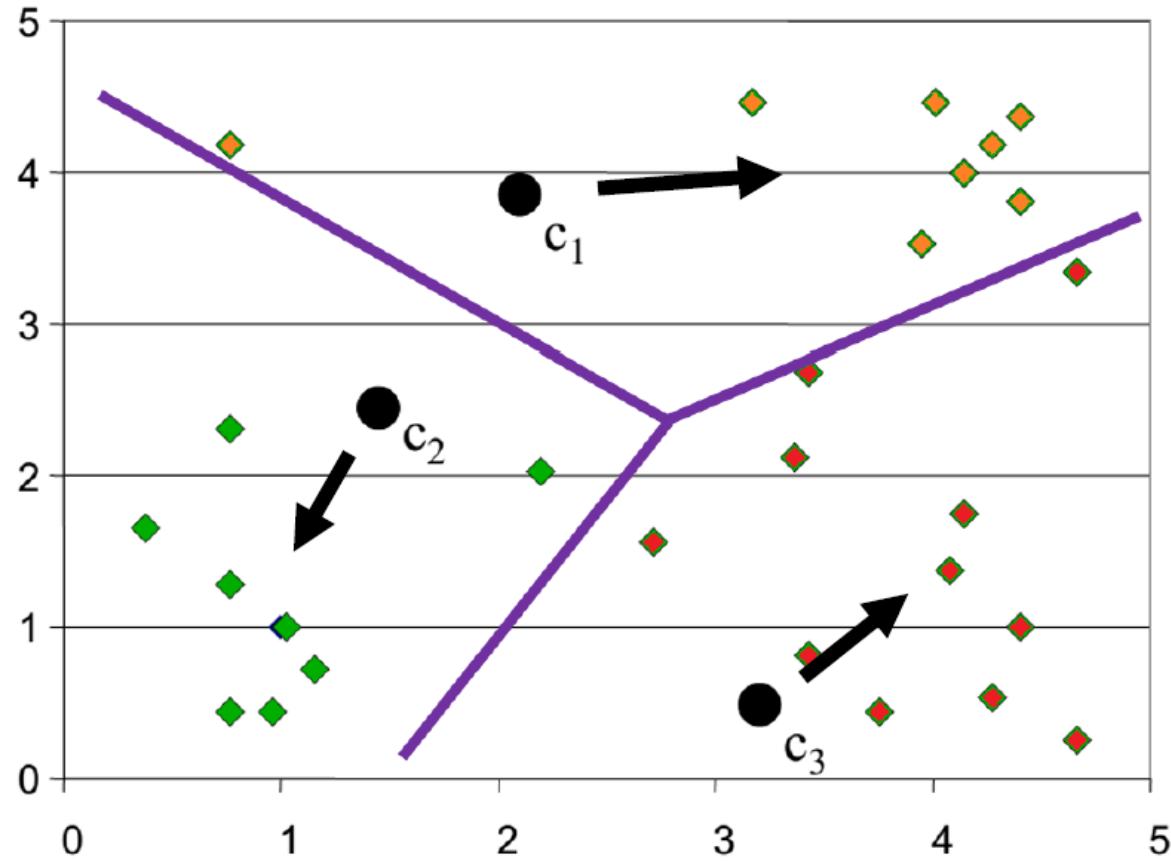
# K-Means (example)

Determine cluster membership for each input



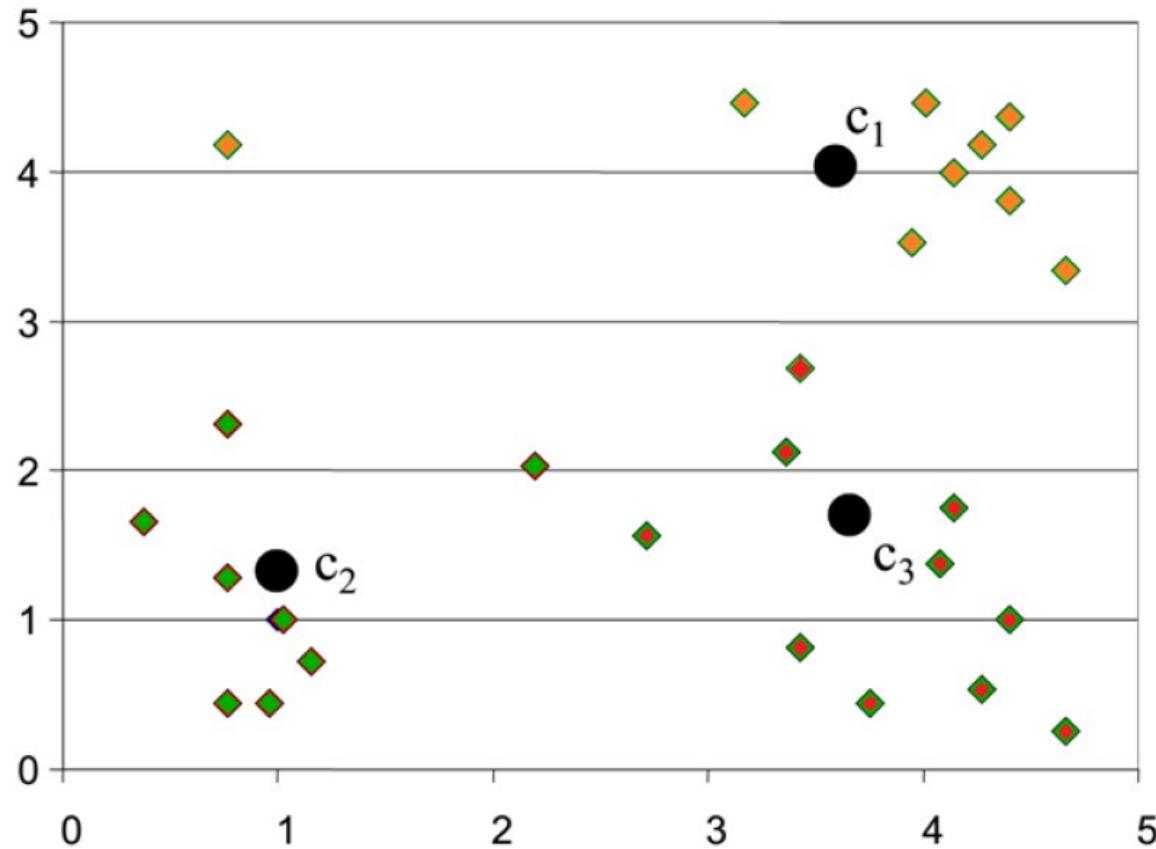
# K-Means (example)

Re-estimate cluster centers (adapt synaptic weights)



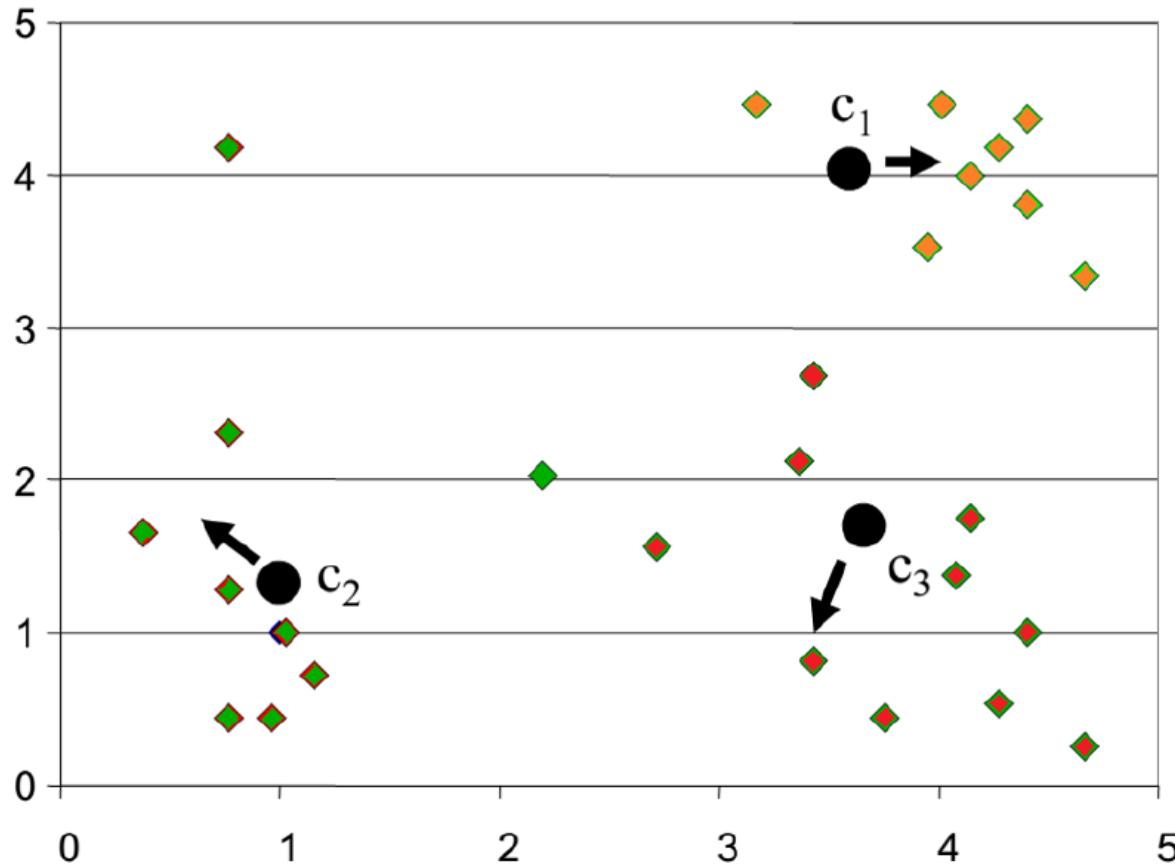
# K-Means (example)

Result of first iteration



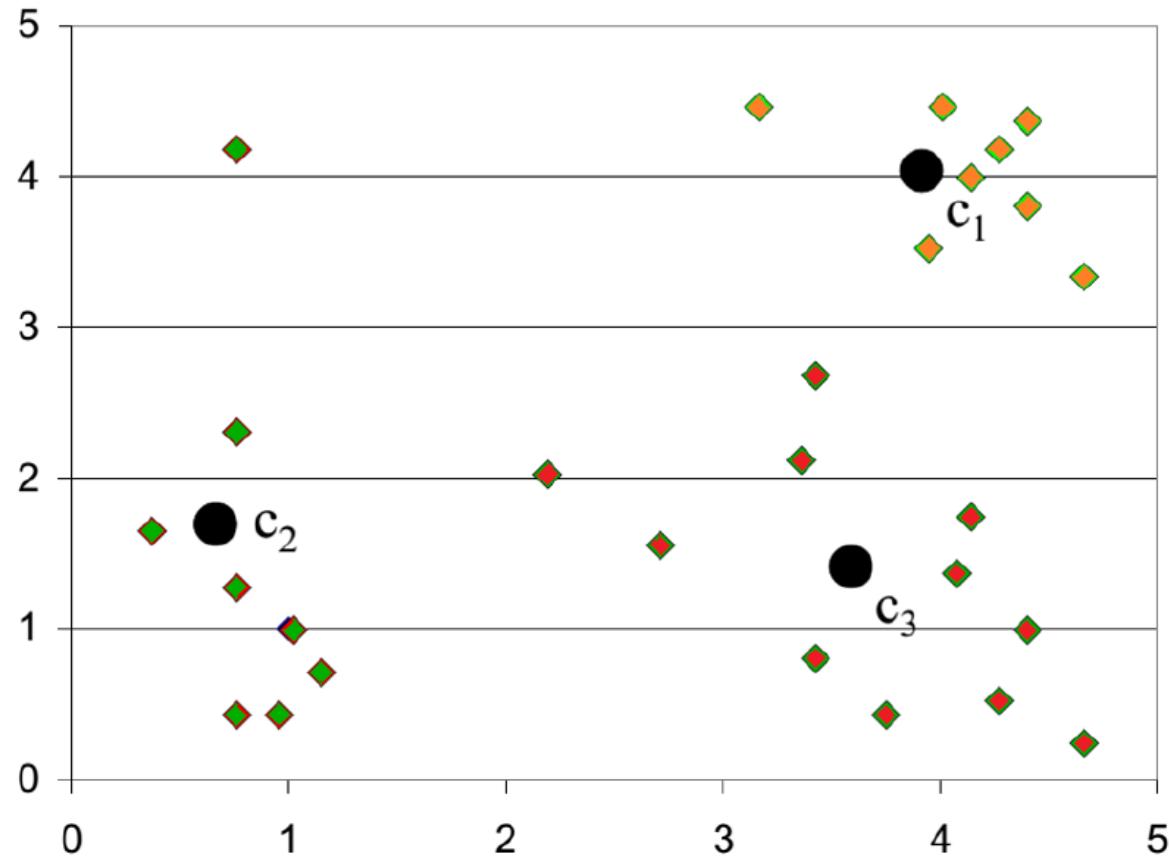
# K-Means (example)

Second iteration



# K-Means (example)

Result of second iteration



# K-Means (objective/cost function)

$$\min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \mu_i\|^2$$

The diagram illustrates the components of the K-Means objective function. It features a mathematical expression with annotations:

- Sum over clusters**: Points to the outer summation symbol  $\sum_{i=1}^k$ .
- Sum over samples in the cluster**: Points to the inner summation symbol  $\sum_{\mathbf{x} \in S_i}$ .
- Sample**: Points to a single term  $\|\mathbf{x} - \mu_i\|^2$ .
- Cluster center**: Points to the term  $\mu_i$ .

# K-Means (strengths)

- Simple: easy to understand and to implement
- K-means is the most popular clustering algorithm.
- Efficient: Time complexity:  $O(tknd)$ , where
  - $n$  is the number of data points,
  - $d$  is the number of features (dimensionality),
  - $k$  is the number of clusters,
  - $t$  is the number of iterations.

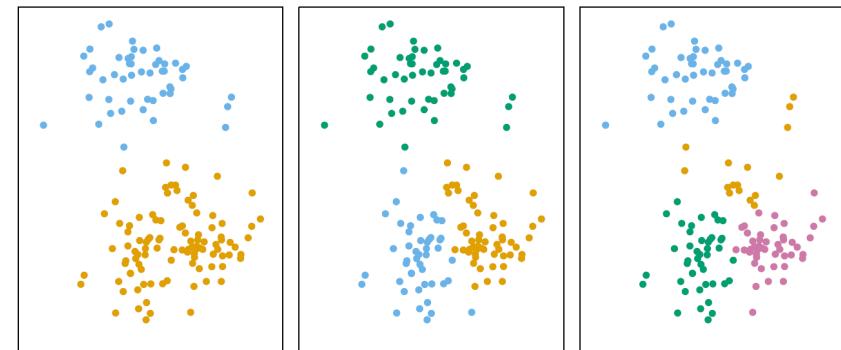
➔ Fast enough, since  $d$ ,  $k$  and  $t$  are usually small.
- Note:
  - It can end-up at a local optimum (not a convex optimization problem).

These two are the same results. Colors don't matter here (unsupervised).

# K-Means (weaknesses)

- Very sensitive to the initial centroids.
  - To address this, we can:
    - Run k-means several times; each time with different initial centroids.
    - Seed the centroids using a better method than random (e.g. choose them to be initially as far as possible from each others).
- User must manually choose  $k$ .
- Sensitive to outliers (why?)
  - Points very far away from the others.

Example: different starting values



$k = 2$

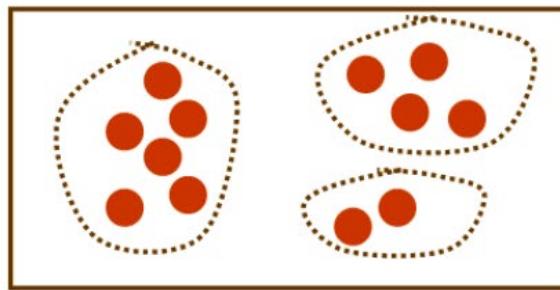
$k = 3$

$k = 4$

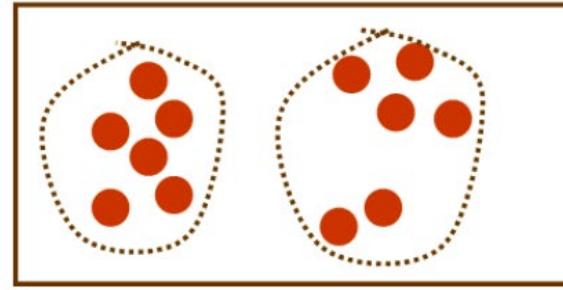
# Hierarchical clustering

# Hierarchical clustering

Up to now, considered “flat” clustering



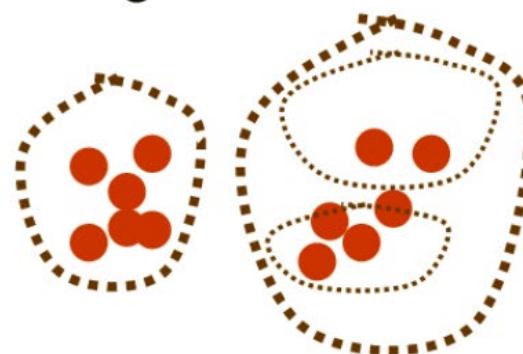
?



It's difficult to decide on the number of clusters beforehand.

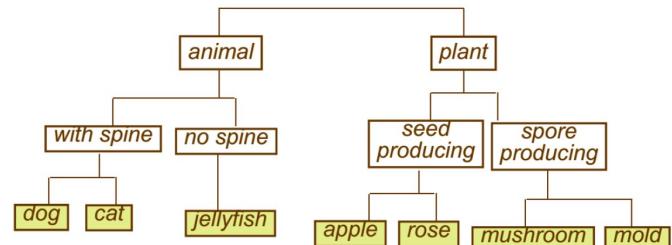
For some data, hierarchical clustering is more appropriate than “flat” clustering

Hierarchical clustering



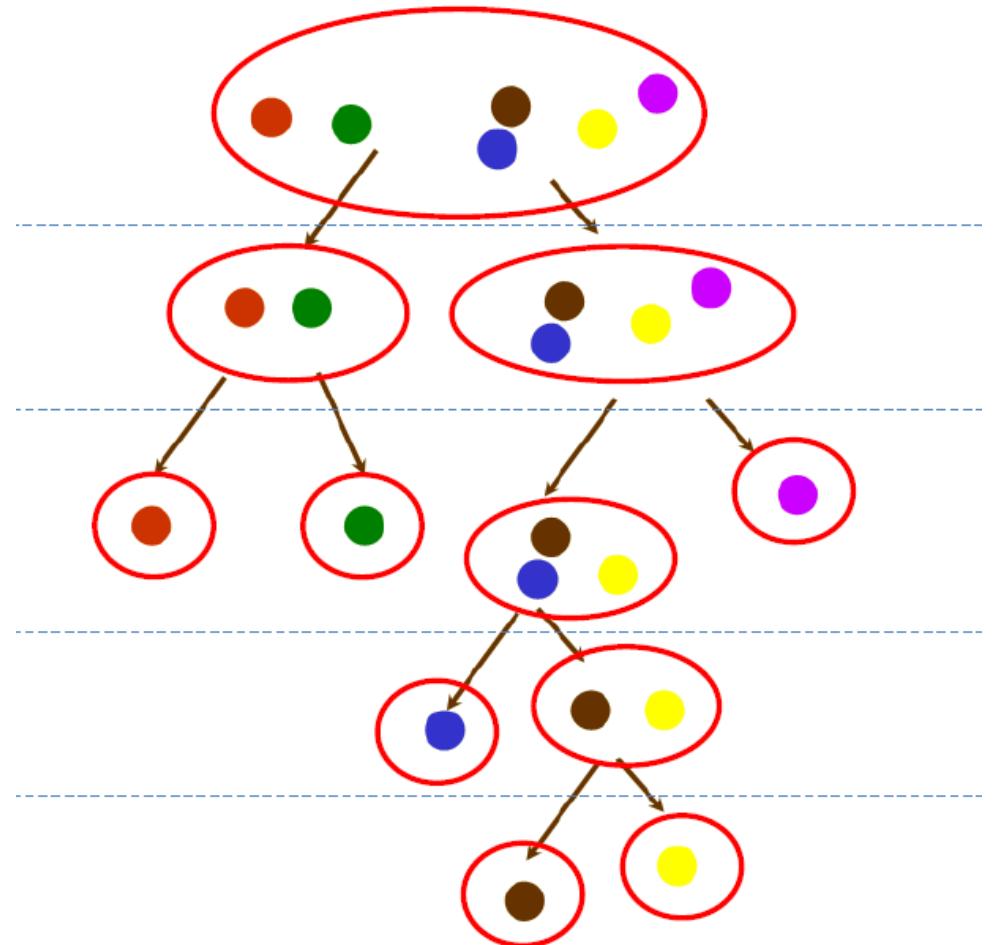
# Hierarchical clustering

- We do not know in advance how many clusters we want; in fact, we end up with a tree-like visual representation of the observations, called a **Dendrogram**, that allows us to view at once the clusterings obtained for each possible number of clusters (from **1** to **n**)
- Top-down approach:
  - start from a single cluster with all examples
  - recursively split clusters into subclusters
- Bottom-up approach:
  - start with **n** clusters of individual examples (singletons)
  - recursively aggregate pairs of clusters
- Hierarchical (agglomerative) clustering
  - Initially, each point in cluster by itself.
  - Repeatedly combine the two “nearest” clusters into one.

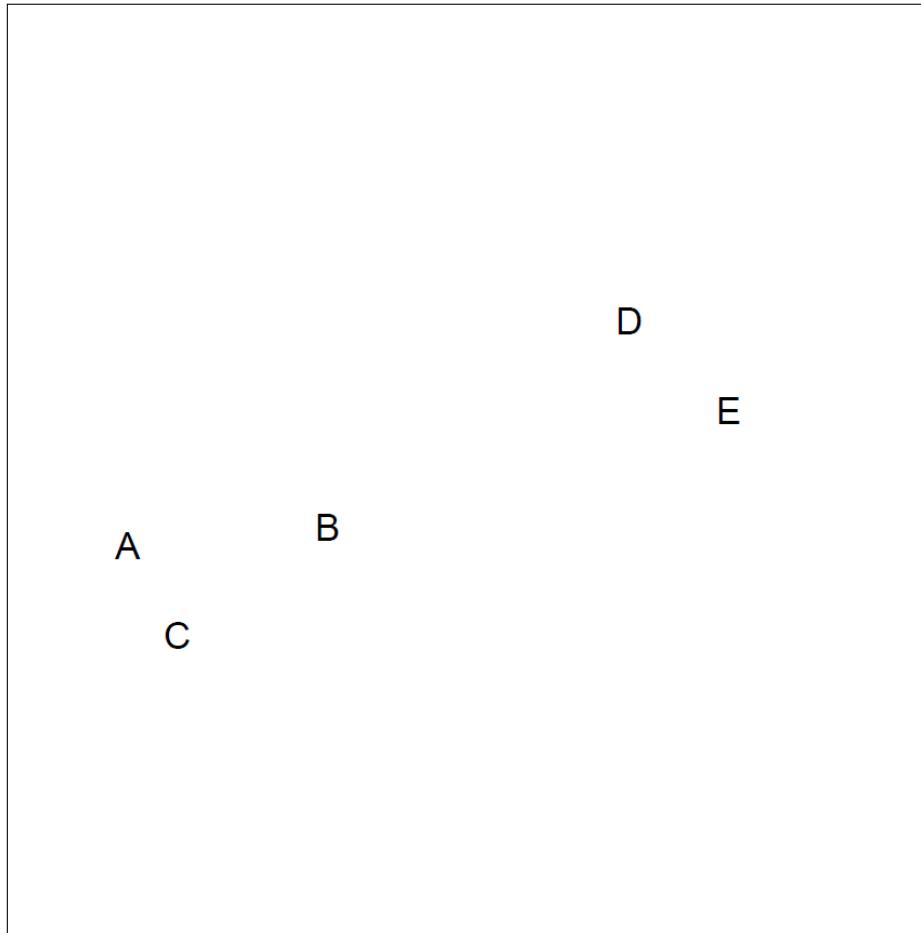


# Hierarchical clustering (divisive)

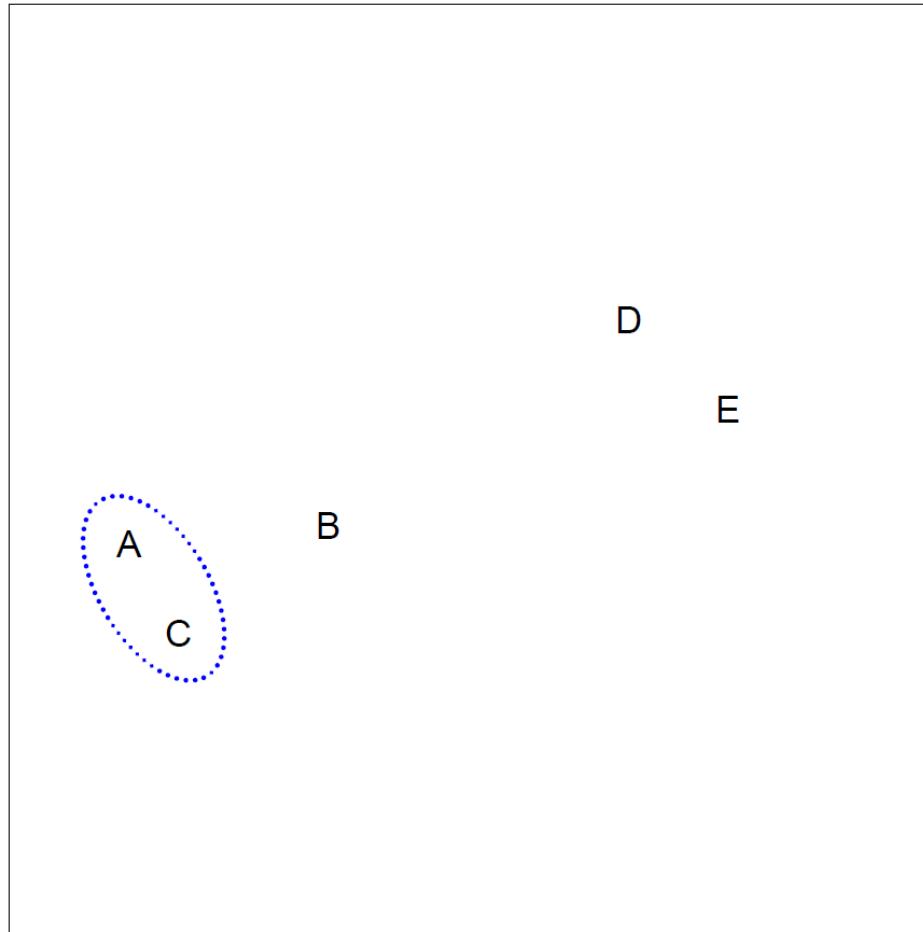
- Any flat clustering algo can be used
  - e.g. k-means with  $k=2$  repeated 5 times.



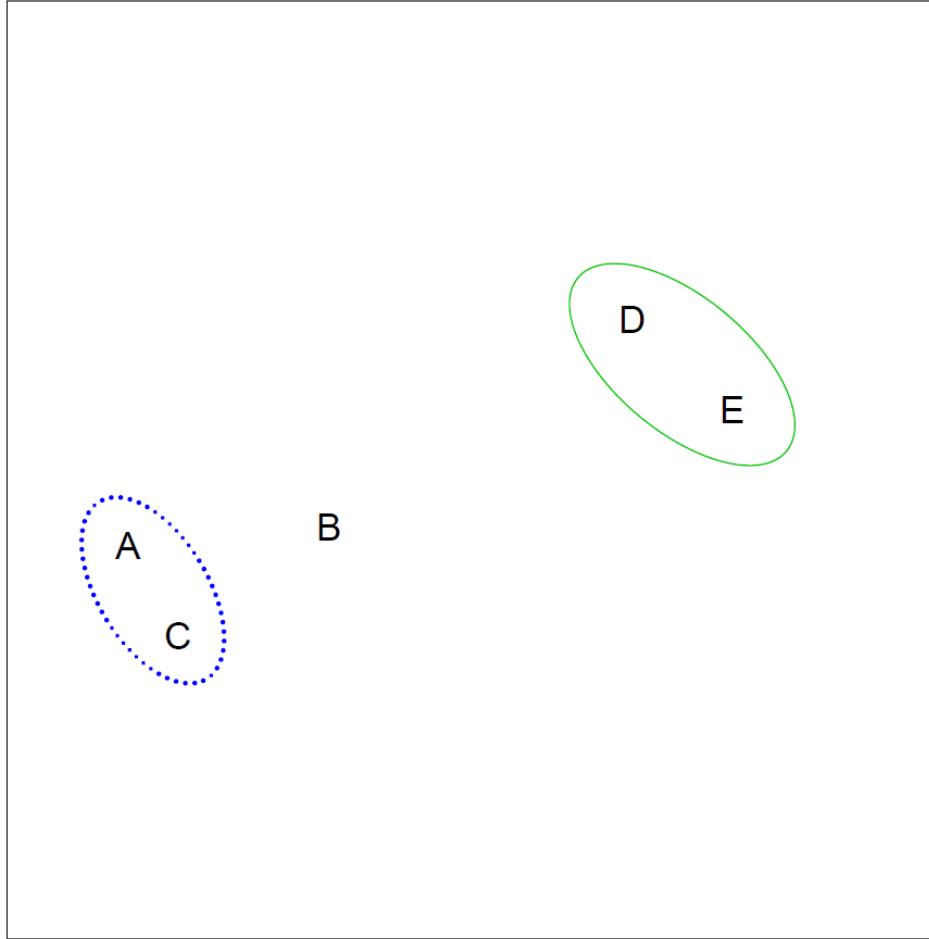
# Hierarchical clustering (agglomerative)



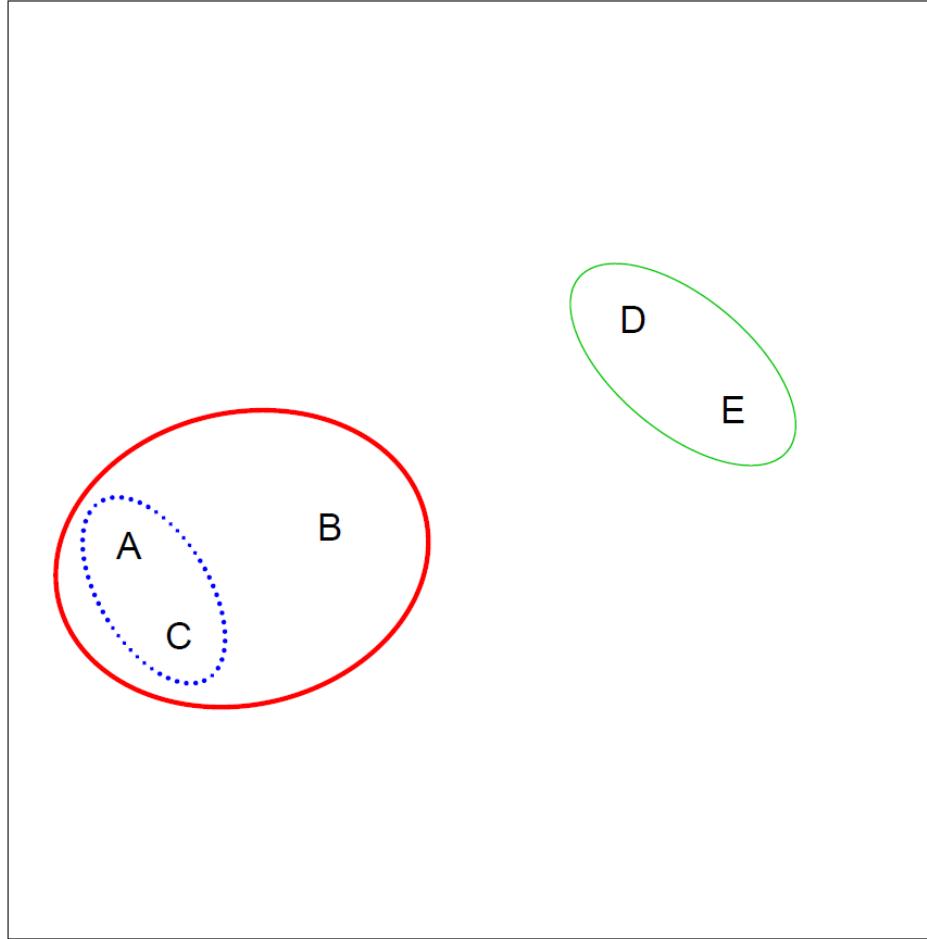
# Hierarchical clustering (agglomerative)



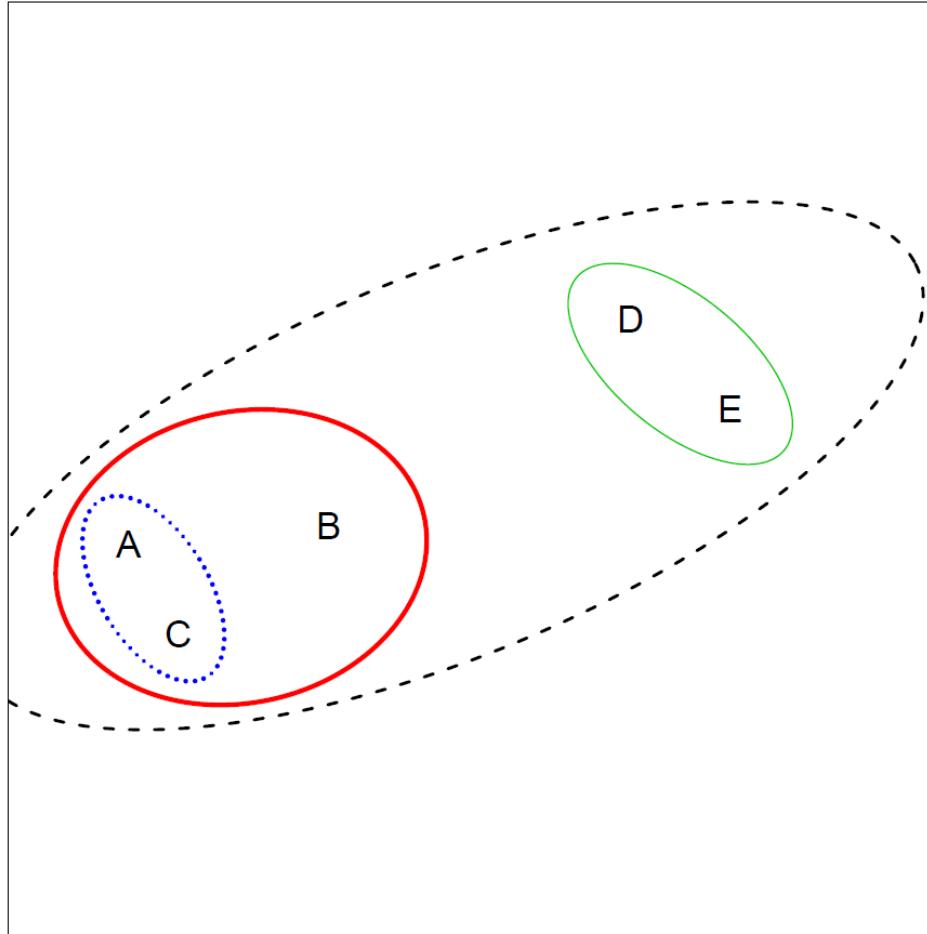
# Hierarchical clustering (agglomerative)



# Hierarchical clustering (agglomerative)

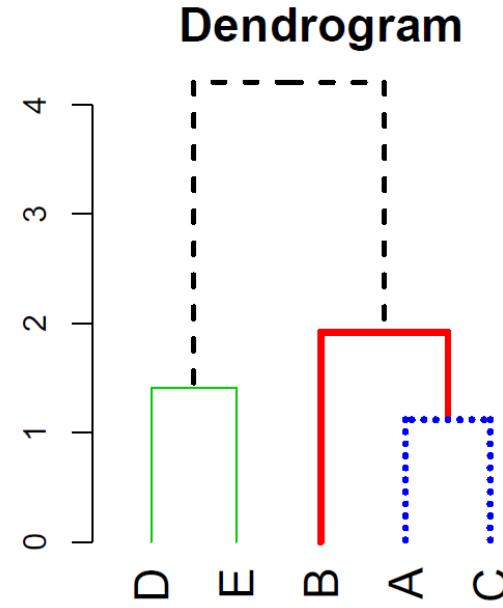
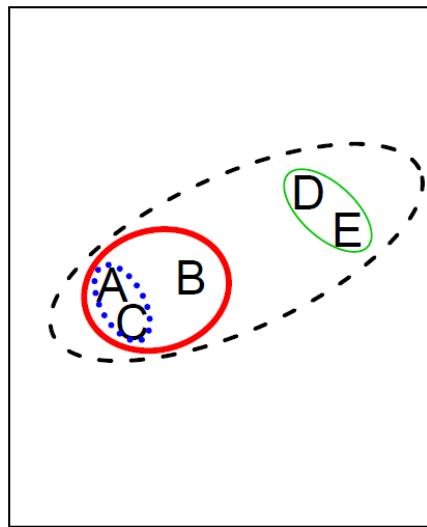


# Hierarchical clustering (agglomerative)



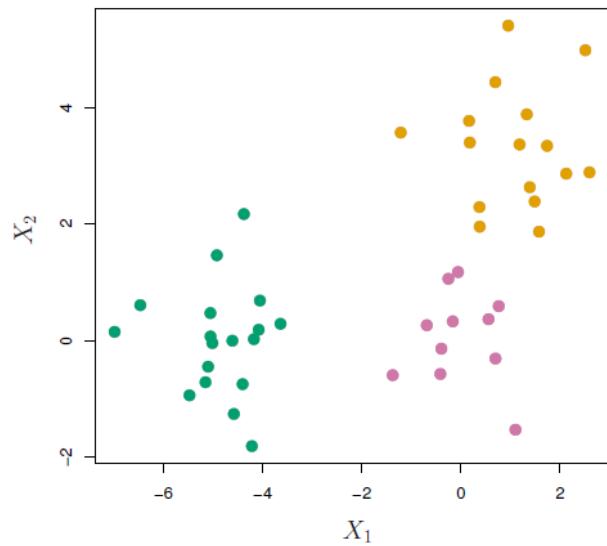
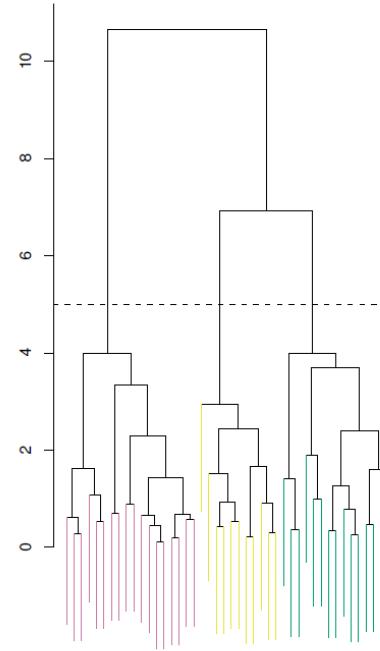
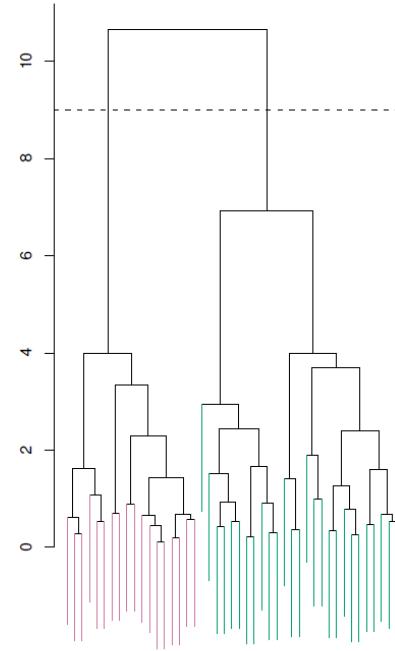
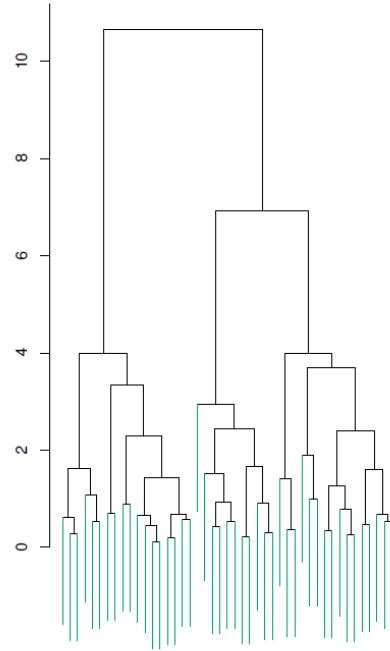
# Hierarchical clustering (agglomerative)

- Start with each point in its own cluster.
- Identify the **closest** two clusters and merge them.
- Repeat.
- Ends when all points are in a single cluster.



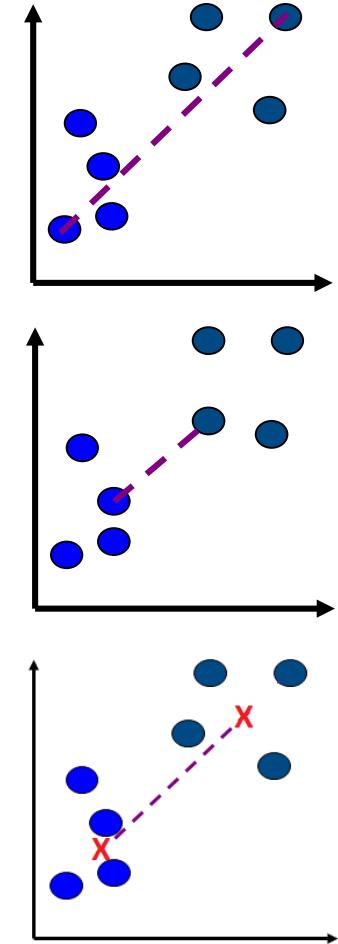
# Hierarchical clustering

- Dendrogram



# Hierarchical clustering

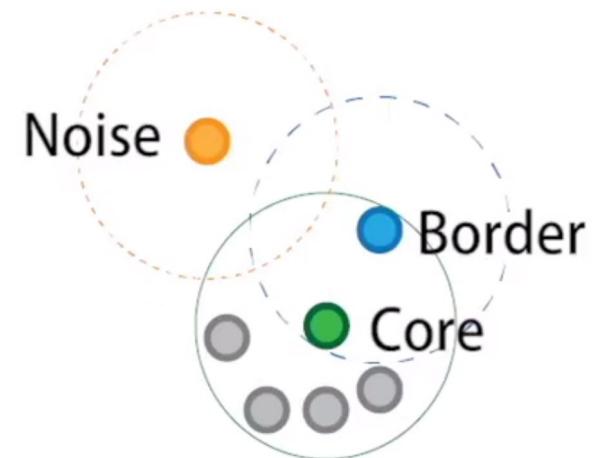
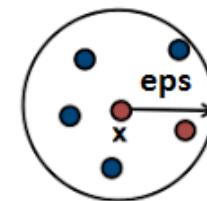
<i>Linkage</i>	<i>Description</i>
Complete	Maximal inter-cluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>largest</i> of these dissimilarities.
Single	Minimal inter-cluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>smallest</i> of these dissimilarities.
Average	Mean inter-cluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>average</i> of these dissimilarities.
Centroid	Dissimilarity between the centroid for cluster A (a mean vector of length $p$ ) and the centroid for cluster B. Centroid linkage can result in undesirable <i>inversions</i> .



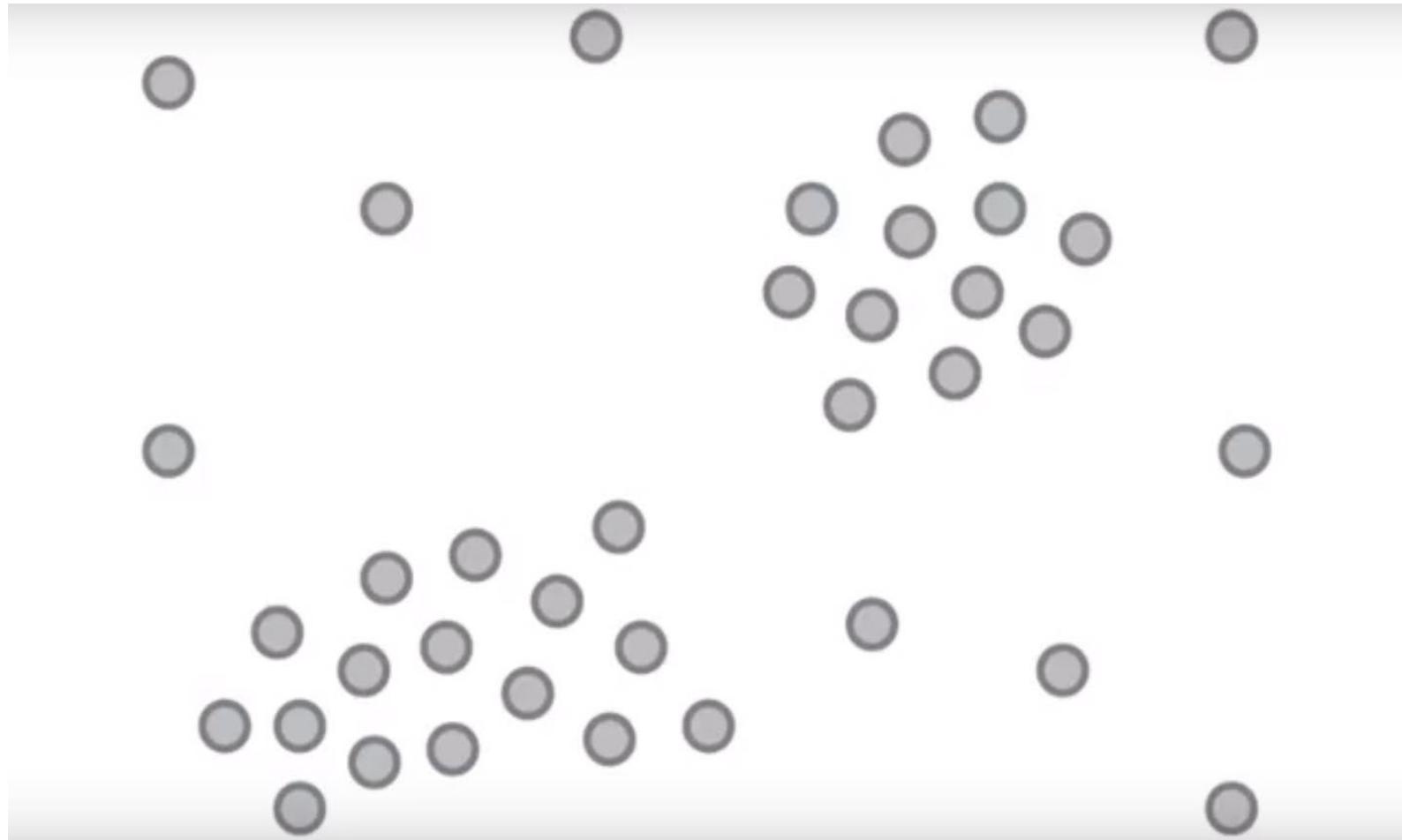
# **DBSCAN clustering method (quick overview)**

# DBSCAN (a density-based algorithm)

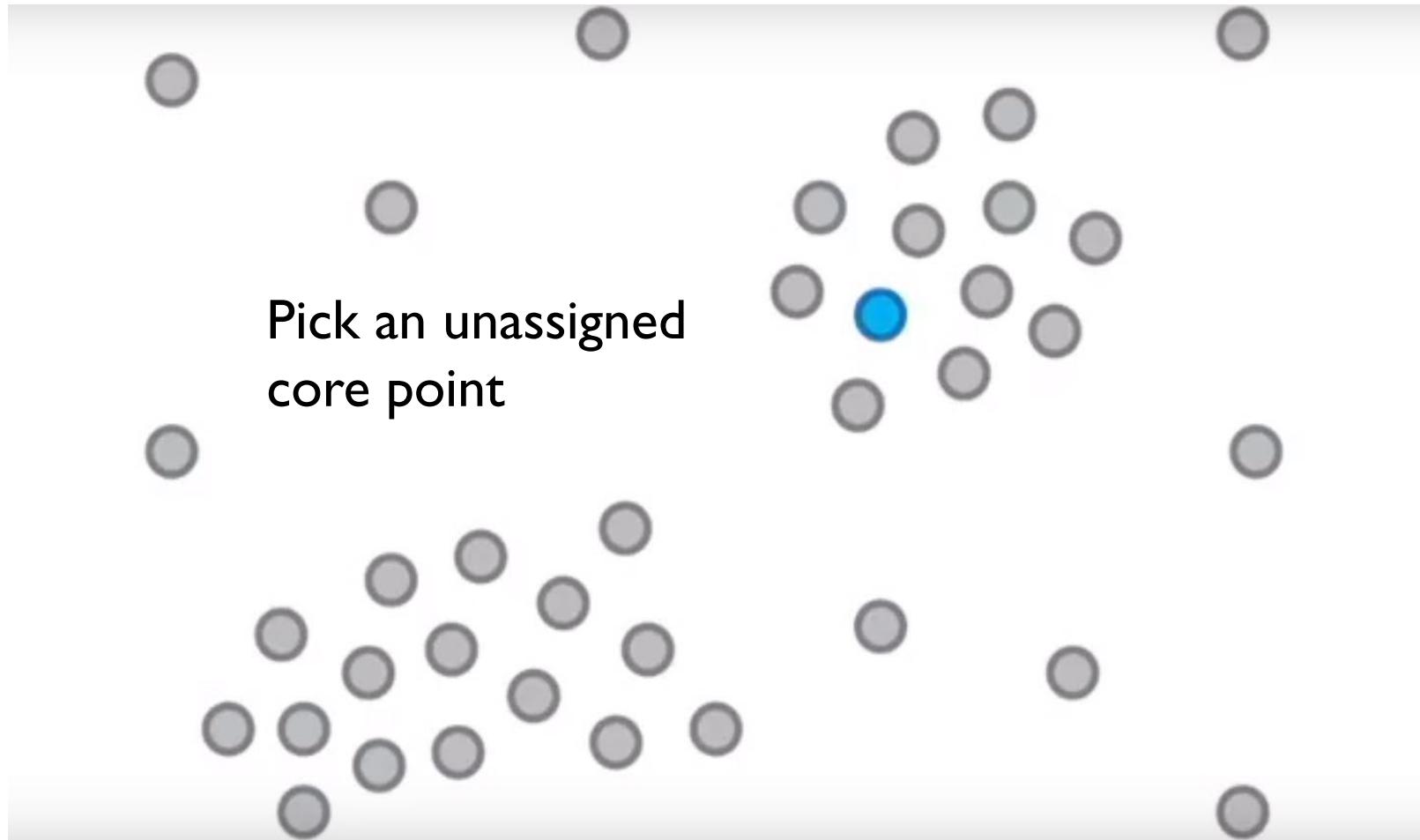
- A cluster is defined as a set of density-connected points.
- Two parameters:
  - *eps* : defines the neighborhood
    - e.g.  $N(x) = \{z \in X : \text{dist}(x, z) < \text{eps}\}$
  - *minpts*: minimum number of points the neighborhood.
- **Core points**: Point  $x$  is a core point if  $|N(x)| > \text{minpts}$
- **Border points**: Not a core point but has at least one core point in its neighborhood.
- **Noise points**: Not a core not a border point.



# DBSCAN

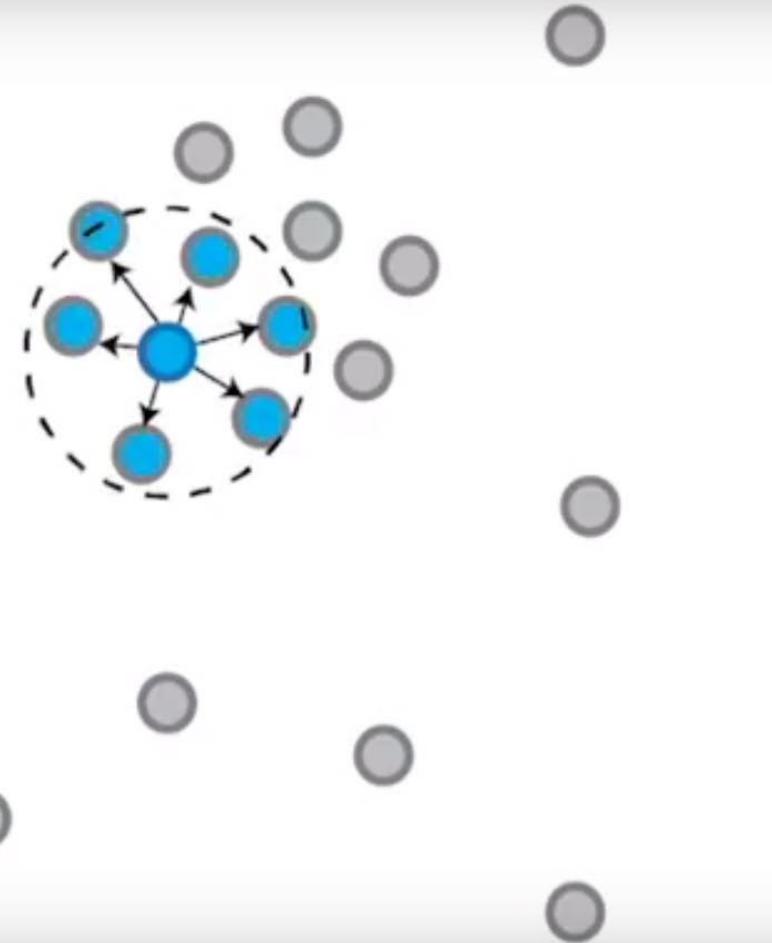


# DBSCAN



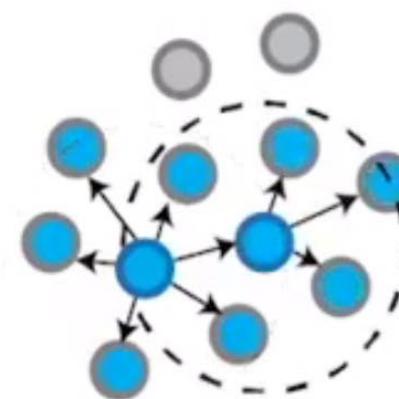
# DBSCAN

Perform a DFS



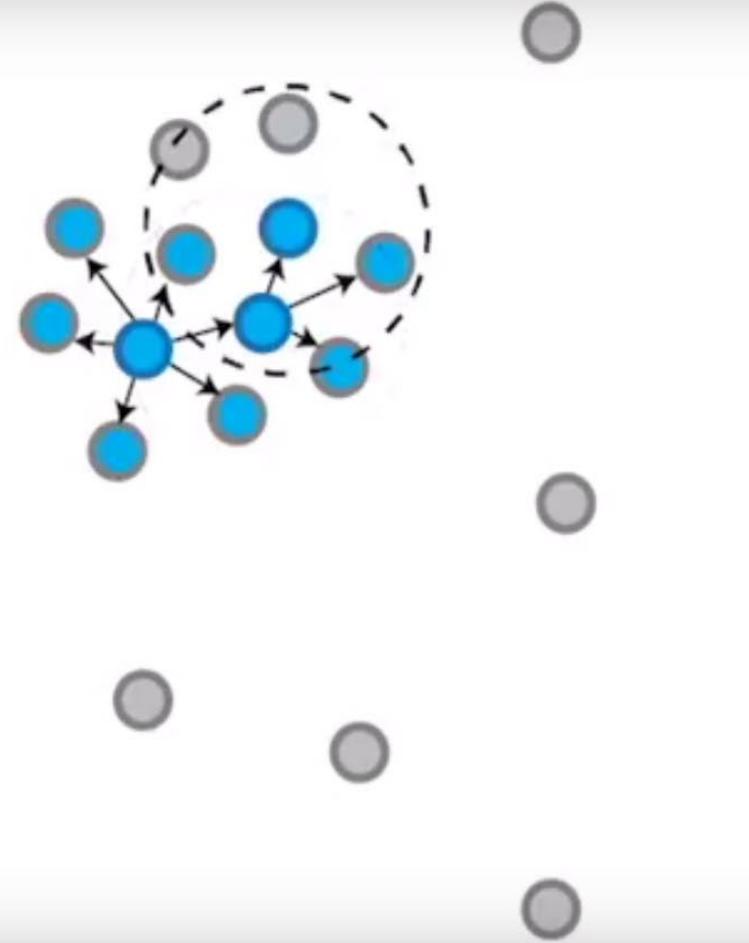
# DBSCAN

Perform a DFS

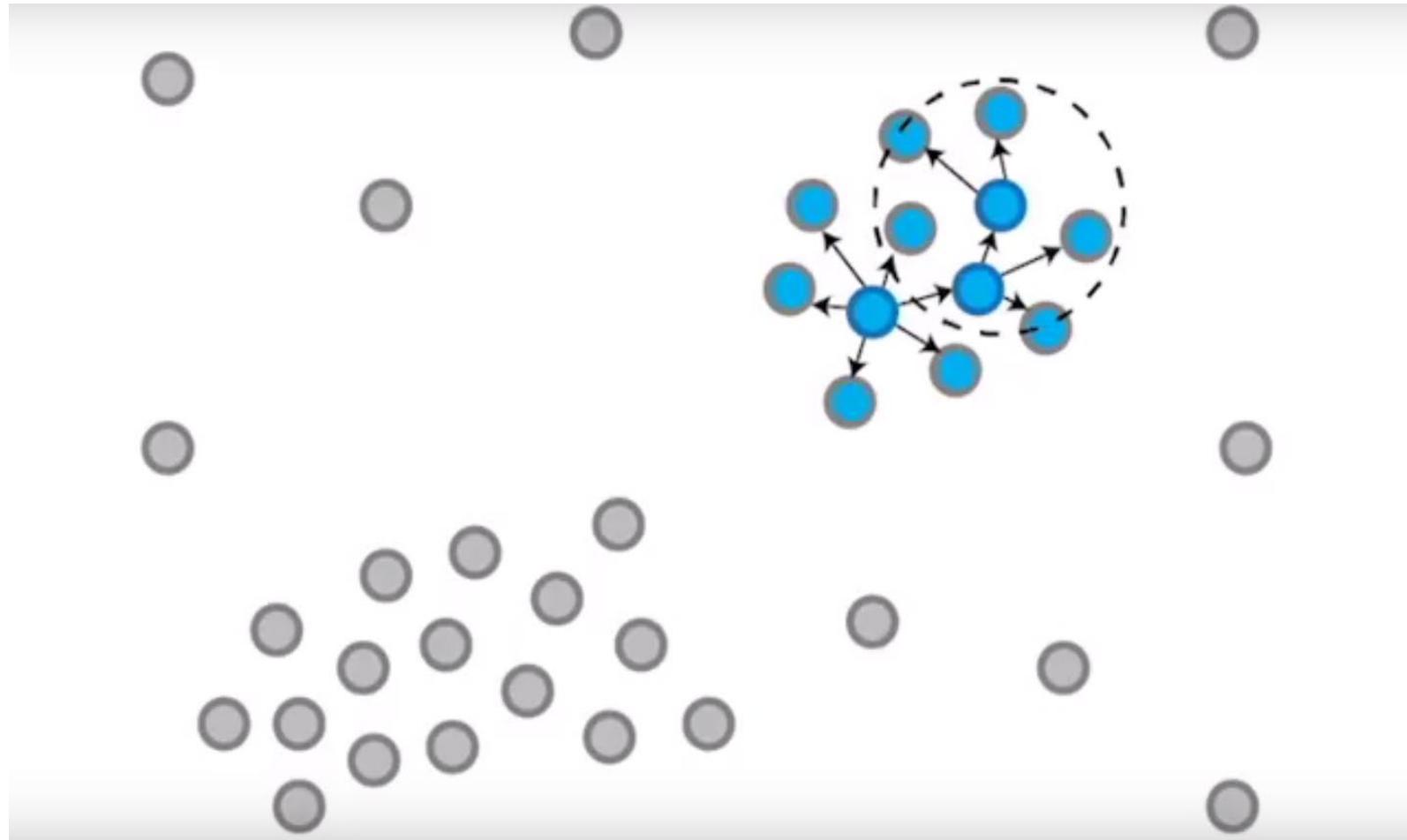


# DBSCAN

Add all points in  
the neighborhood  
to the same cluster

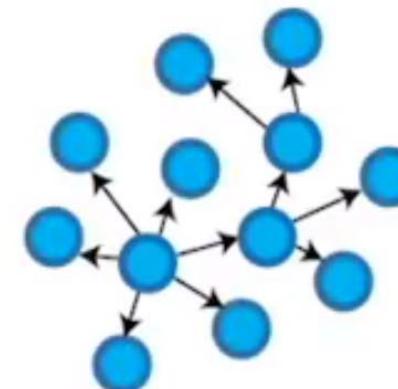


# DBSCAN

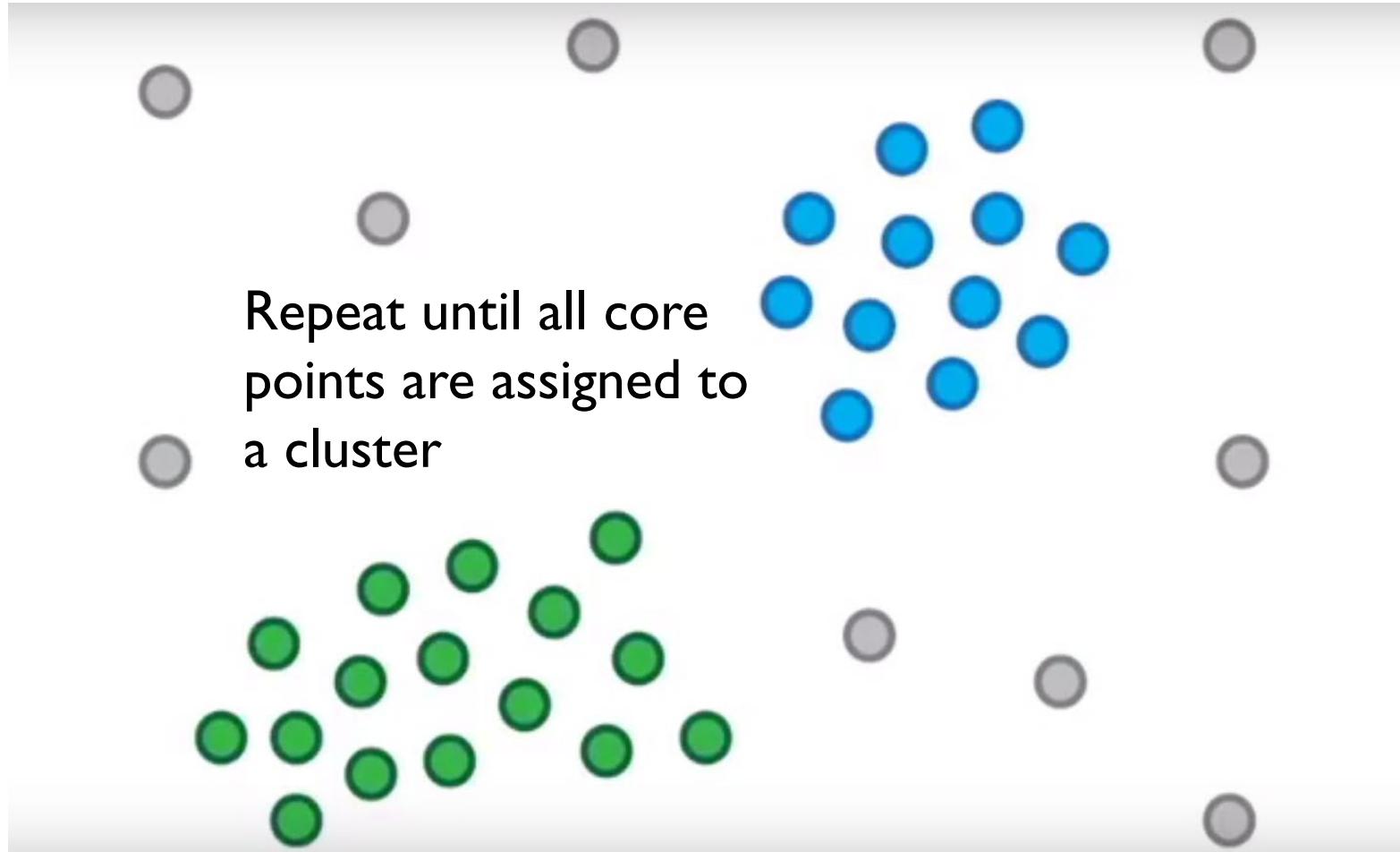


# DBSCAN

Recursively apply the search on each unexplored core point



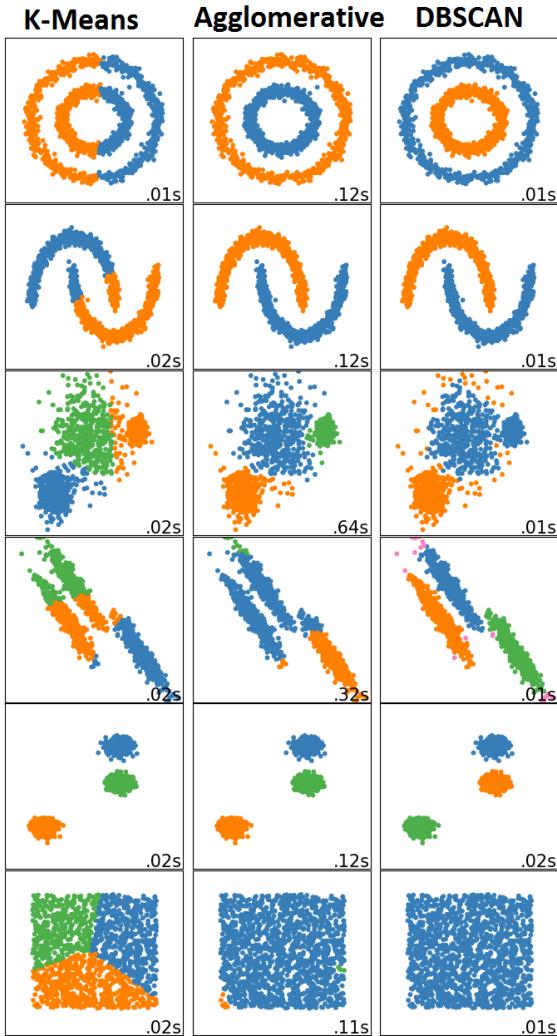
# DBSCAN



# DBSCAN

- **Strengths**
  - Can discover any number of clusters
  - Can find clusters of varying size and shape
    - change
  - Can detect and ignore outliers in the data
- **Weaknesses**
  - Very sensitive the choice of parameters (e.g. neighborhood parameter)

# Comparison



- Choosing the “best” algorithm is a challenge.
  - Every algorithm has limitations and works well with certain types of data.
  - One also needs to choose a suitable distance function and to select other parameter values.
  - Clustering is to certain extent subjective (personal preferences).
- We usually try different algorithms and analyze / explore the results ...

# Evaluating the clustering results

# Evaluation of the clustering quality

- No labels  $\rightarrow$  we do not know the “correct” clusters.
- Usually we use a measure that takes into account:
  - **Intra-cluster** cohesion (compactness):
    - Cohesion measures how close are the points within the same cluster.
  - **Inter-cluster** separation (isolation):
    - Separation means that different clusters should be far away from one another.

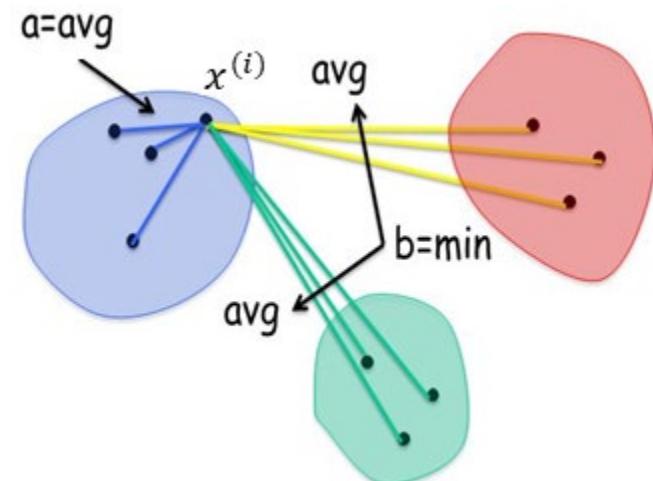
## Silhouette coefficient:

- Let  $a_i$  be the average distance from point  $x^{(i)}$  to the points in the same cluster as  $x^{(i)}$  (i.e. blue cluster in the figure).
- Let  $b_i$  be the average distance from point  $x^{(i)}$  to the points in the closest cluster to  $x^{(i)}$  (e.g. say the green cluster is closer to  $x^{(i)}$  than the red one).
- The silhouette coefficient  $s_i$  of  $x^{(i)}$  is then

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

Final silhouette coefficient :

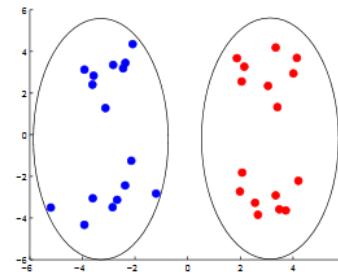
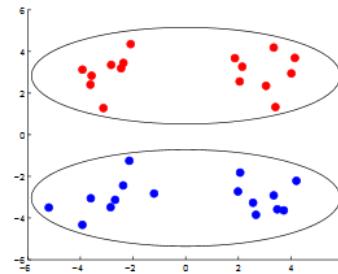
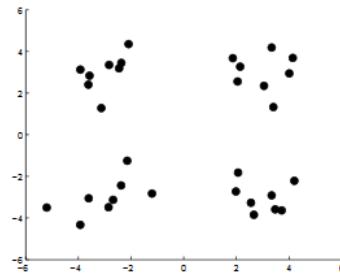
$$S = \frac{1}{n} \sum_{i=1}^n s_i$$



# **Some problems with clustering**

# Problems with clustering (subjective)

- Multiple (different) clustering solutions, for the same data

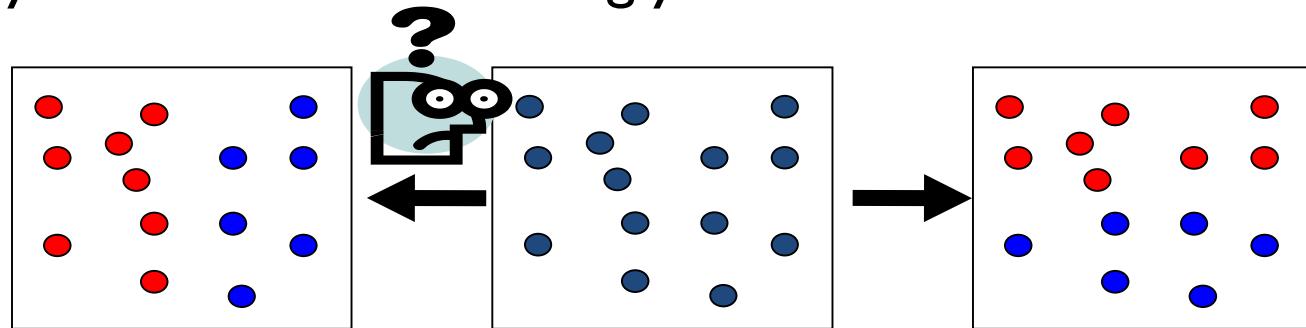


User ID	Country	Age	Gender	Blood	Heartbeat	Weight	Height	Sports	Income	Profession
1										
2										
3										
4	China	young	male							
5										
6		old								
7			female							
8										
9										
10	US	young	male							
11										
12		old	female							
13										
14										
15										

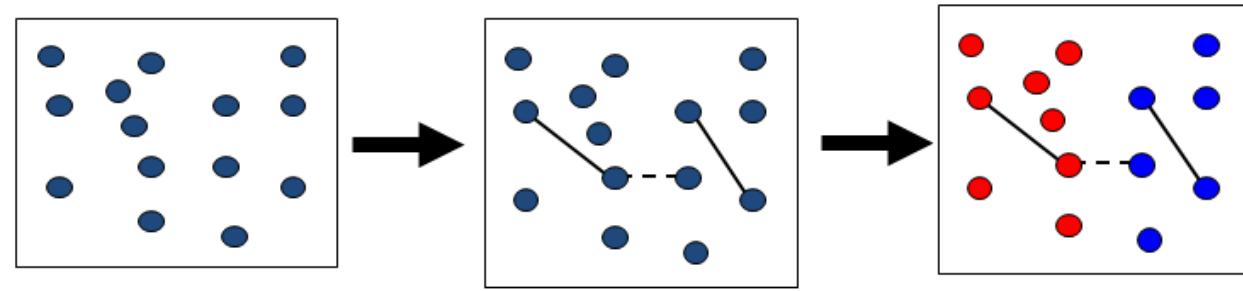
People can be clustered by different criteria

# Problems with clustering (towards semi-supervised)

- How do you tell which clustering you want?



- Constrained clustering techniques:



— Same-cluster constraint  
(must-link)

--- Different-cluster constraint  
(cannot-link)