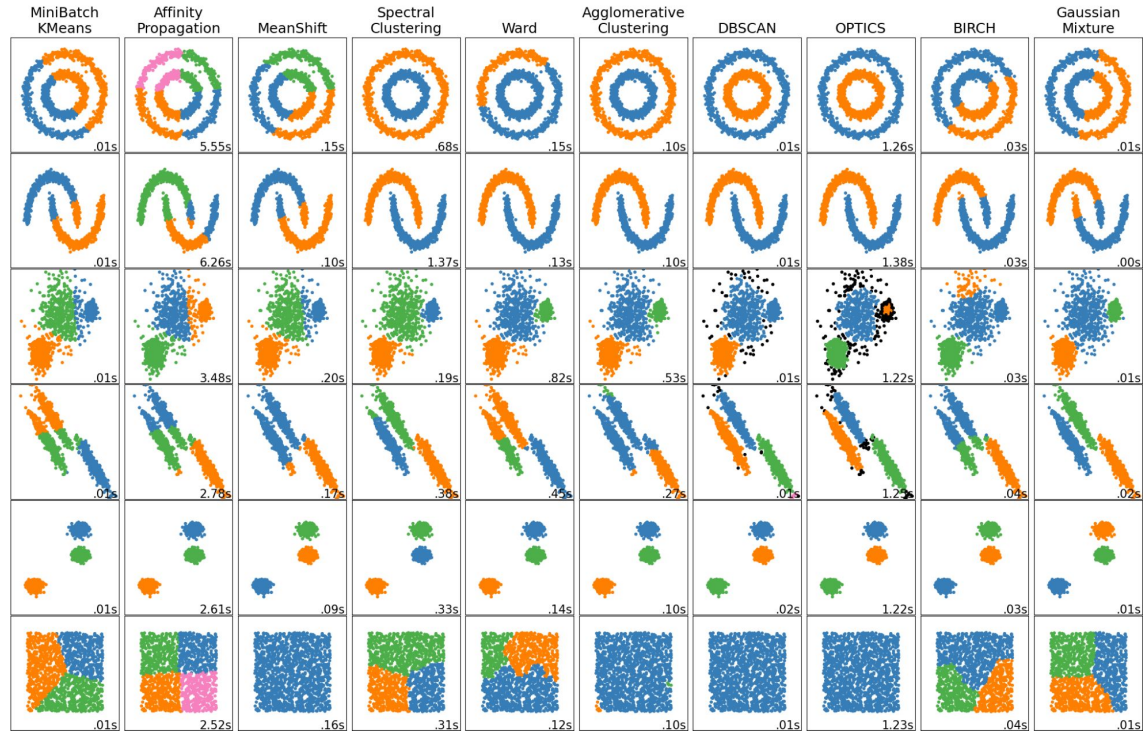


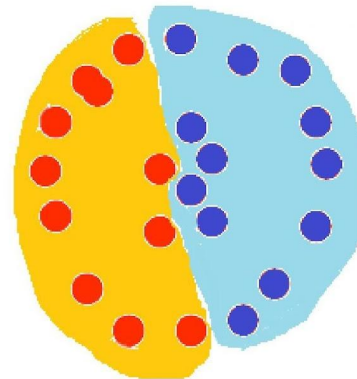
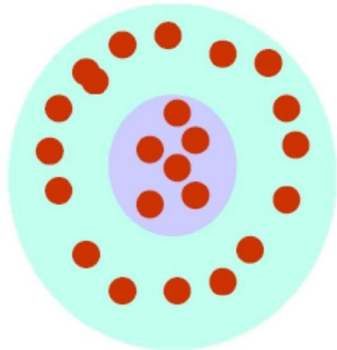
Unsupervised Learning

Clustering



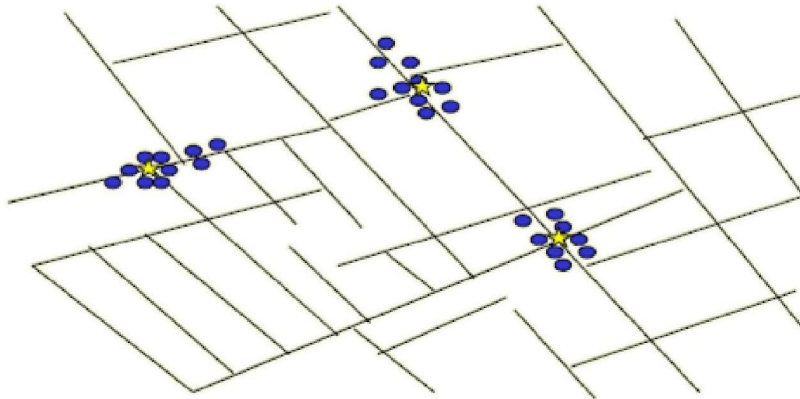
What is clustering?

- The organization of unlabeled data into similarity groups called clusters.
- A cluster is a collection of data items which are “similar” between them, and “dissimilar” to data items in other clusters.



Historic application of clustering

- John Snow, a London physician plotted the location of cholera deaths on a map during an outbreak in the 1850s.
- The locations indicated that cases were clustered around certain intersections where there were polluted wells -- thus exposing both the problem and the solution.

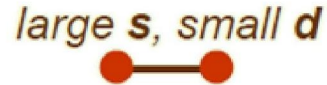


From: Nina Mishra HP Labs

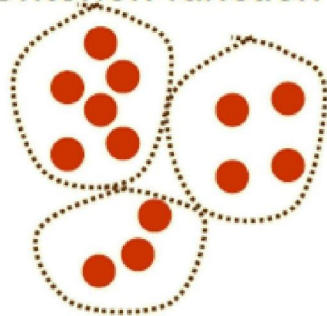
What do we need for clustering?

1. Proximity measure, *either*

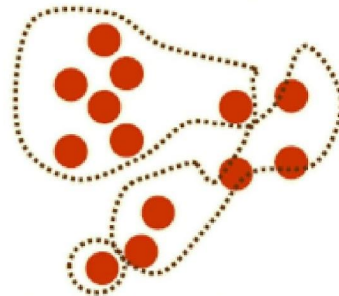
- similarity measure $s(x_i, x_k)$: large if x_i, x_k are similar
- dissimilarity(or distance) measure $d(x_i, x_k)$: small if x_i, x_k are similar



2. Criterion function to evaluate a clustering



good clustering



bad clustering

3. Algorithm to compute clustering

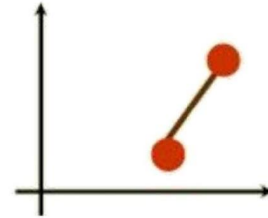
- For example, by optimizing the criterion function

Distance (dissimilarity) measures

- Euclidean distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{k=1}^d (\mathbf{x}_i^{(k)} - \mathbf{x}_j^{(k)})^2}$$

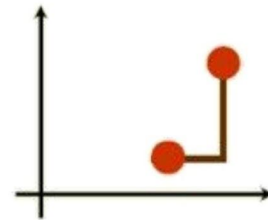
- translation invariant



- Manhattan (city block) distance

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_{k=1}^d |\mathbf{x}_i^{(k)} - \mathbf{x}_j^{(k)}|$$

- approximation to Euclidean distance, cheaper to compute



- They are special cases of **Minkowski distance**:

$$d_p(\mathbf{x}_i, \mathbf{x}_j) = \left(\sum_{k=1}^m |x_{ik} - x_{jk}|^p \right)^{\frac{1}{p}}$$

(p is a positive integer)

What properties should a distance measure have?

$D(A,B) = D(B,A)$ Symmetry

$D(A,A) = 0$ Constancy of Self-Similarity

$D(A,B) = 0$ iff $A = B$ Positivity (Separation)

$D(A,B) \leq D(A,C) + D(B,C)$ Triangular Inequality

Cluster evaluation (a hard problem)

- **Intra-cluster cohesion** (compactness):
 - Cohesion measures how near the data points in a cluster are to the cluster centroid.
 - Sum of squared error (SSE) is a commonly used measure.
- **Inter-cluster separation** (isolation):
 - Separation means that different cluster centroids should be far away from one another.
- In most applications, expert judgments are still the key

How many clusters?



- Possible approaches
 1. fix the number of clusters to k
 2. find the best clustering according to the criterion function (number of clusters may vary)

Types of clustering:

1. **Centroids-based Clustering (Partitioning methods)**
2. **Connectivity-based Clustering (Hierarchical clustering)**
3. **Distribution-based Clustering**
4. **Density-based Clustering (Model-based methods)**
5. **Fuzzy Clustering**
6. **Constraint-based (Supervised Clustering)**
7. **Competitive (SOM)**

Types of clustering:

1. **Centroids-based Clustering (Partitioning methods)**
2. **Connectivity-based Clustering (Hierarchical clustering)**
3. **Distribution-based Clustering**
4. **Density-based Clustering (Model-based methods)**
5. **Fuzzy Clustering**
6. **Constraint-based (Supervised Clustering)**
7. **Competitive (SOM)**

K-Means clustering

- K-means (MacQueen, 1967) is a **partitional clustering** algorithm
- Let the set of data points D be $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$,
where $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ir})$ is a **vector** in $X \subseteq R^r$, and r is the number of dimensions.
- The k -means algorithm partitions the given data into k clusters:
 - Each cluster has a cluster **center**, called **centroid**.
 - k is specified by the user

K-means algorithm

- Given k , the *k-means* algorithm works as follows:
 1. Choose k (random) data points (*seeds*) to be the initial *centroids*, cluster centers
 2. Assign each data point to the closest *centroid*
 3. Re-compute the *centroids* using the current cluster memberships
 4. If a convergence criterion is not met, repeat steps 2 and 3

K-means convergence (stopping) criterion

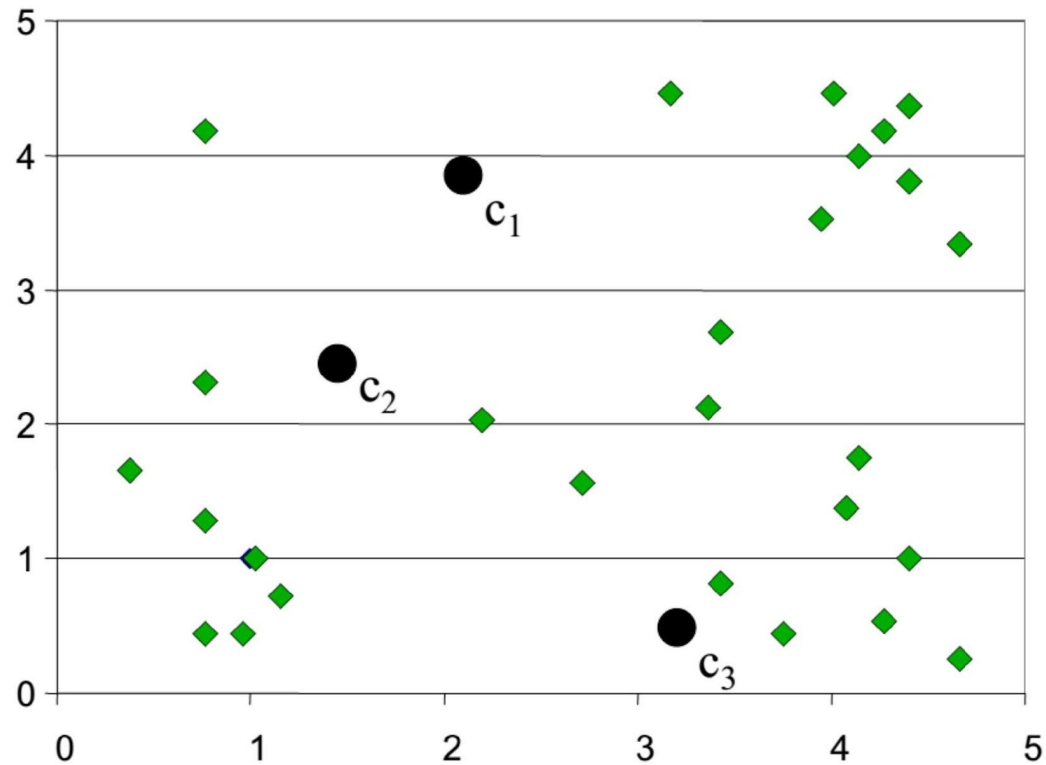
- no (or minimum) re-assignments of data points to different clusters, *or*
- no (or minimum) change of centroids, *or*
- minimum decrease in the **sum of squared error (SSE)**,

$$SSE = \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} d(\mathbf{x}, \mathbf{m}_j)^2$$

- C_j is the j th cluster,
- \mathbf{m}_j is the centroid of cluster C_j (the mean vector of all the data points in C_j),
- $d(\mathbf{x}, \mathbf{m}_j)$ is the (Euclidean) distance between data point \mathbf{x} and centroid \mathbf{m}_j .

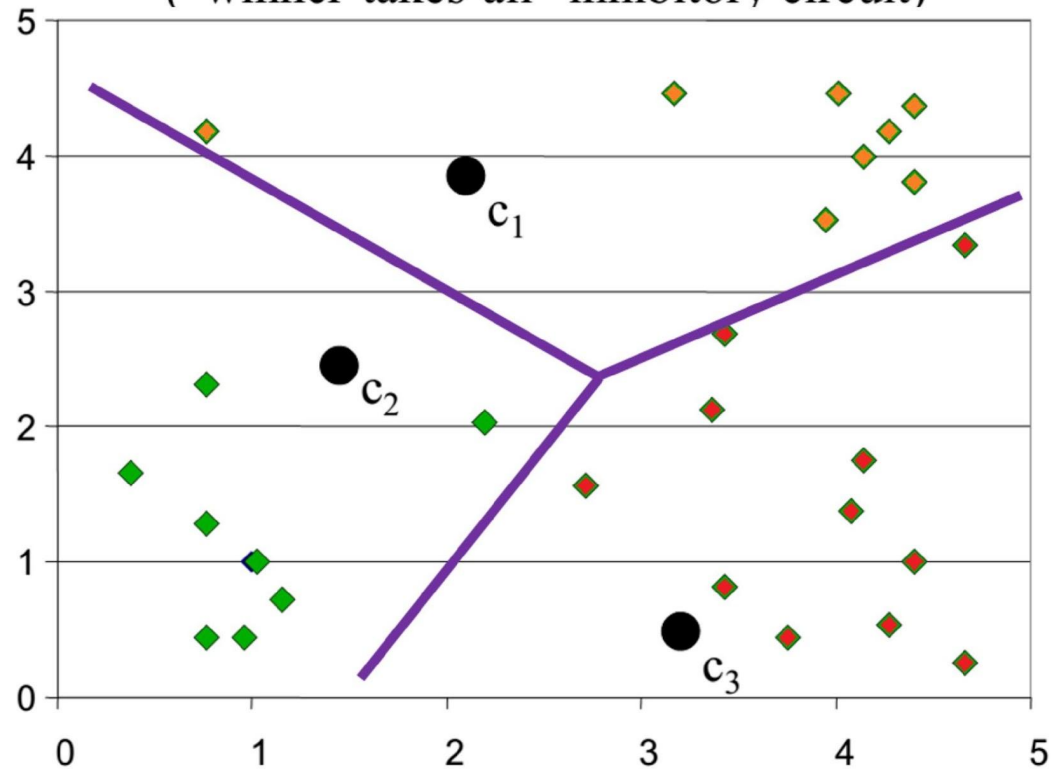
K-means clustering example: step 1

Randomly initialize the cluster centers (synaptic weights)



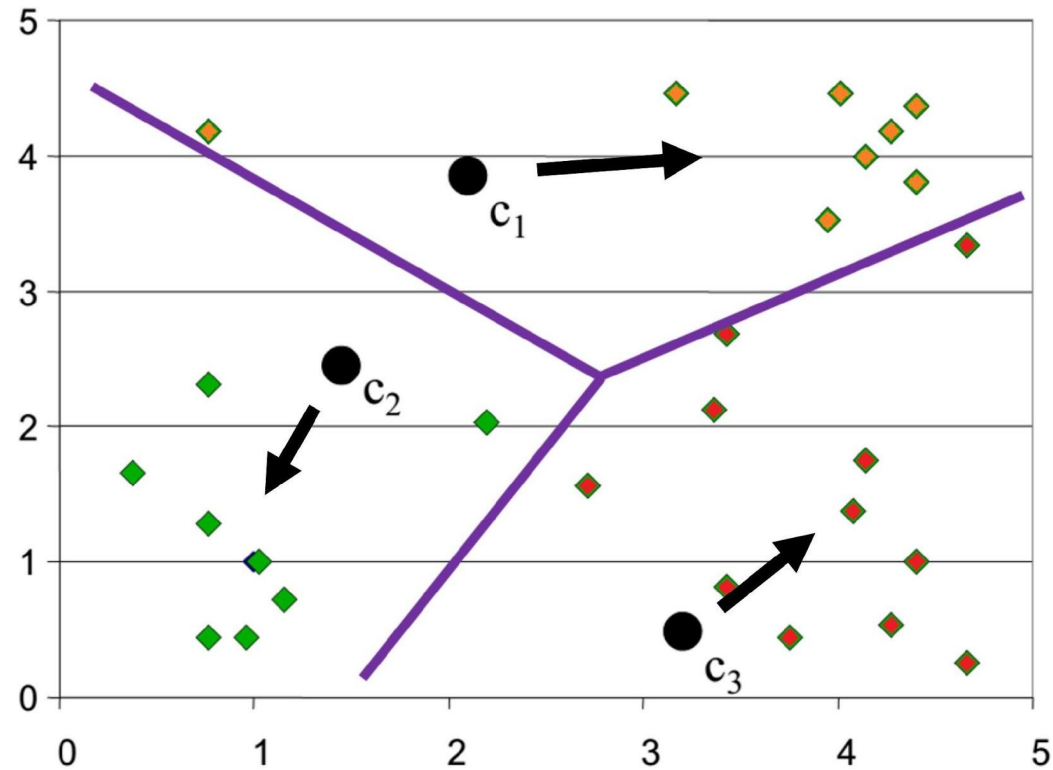
K-means clustering example – step 2

Determine cluster membership for each input
 (“winner-takes-all” inhibitory circuit)



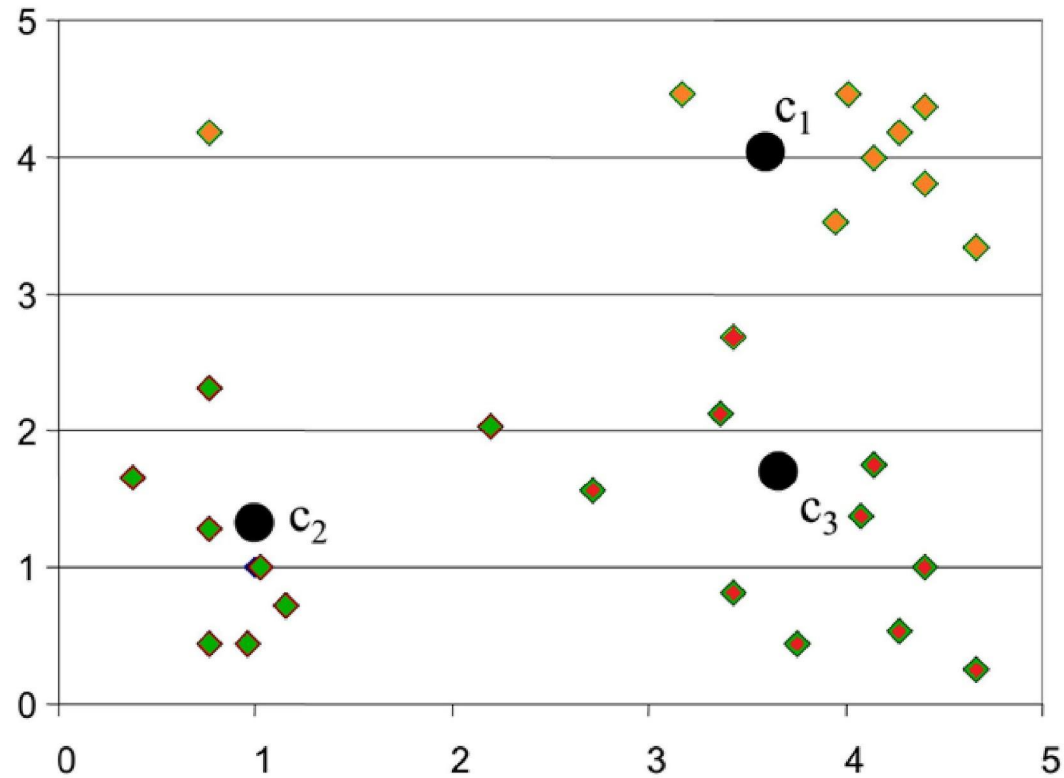
K-means clustering example – step 3

Re-estimate cluster centers (adapt synaptic weights)



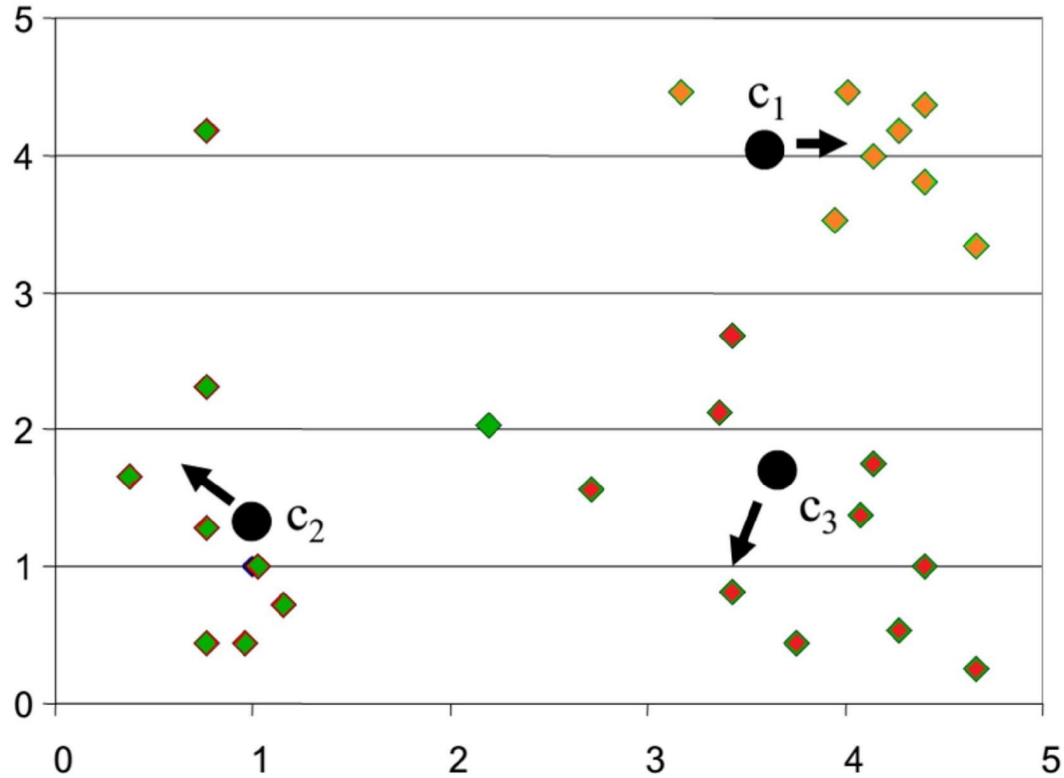
K-means clustering example

Result of first iteration



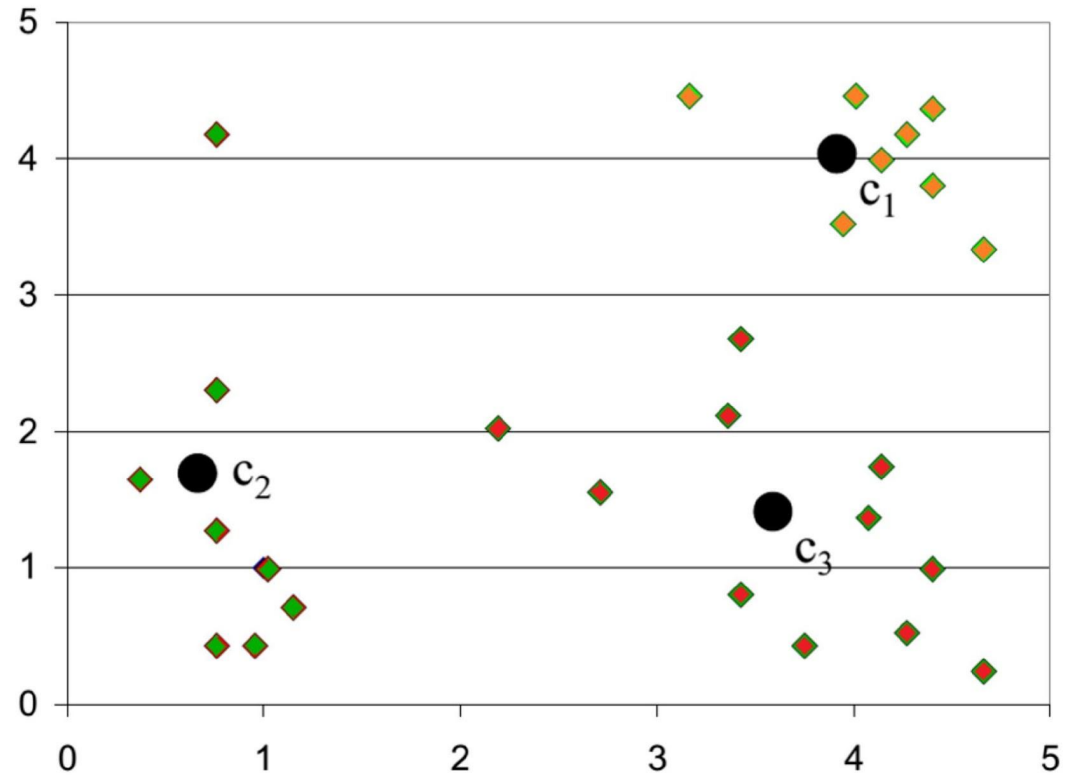
K-means clustering example

Second iteration



K-means clustering example

Result of second iteration



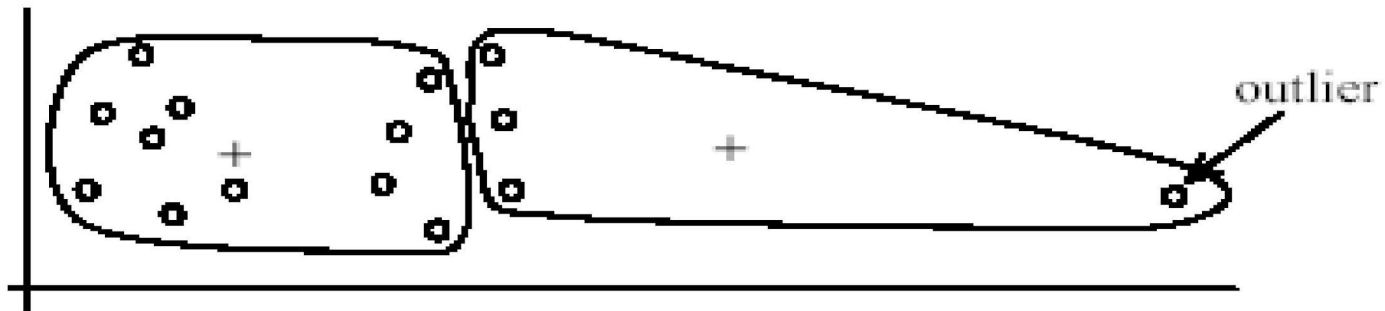
Why use K-means?

- Strengths:
 - Simple: easy to understand and to implement
 - Efficient: Time complexity: $O(tkn)$, where n is the number of data points, k is the number of clusters, and t is the number of iterations.
 - Since both k and t are small. k -means is considered a linear algorithm.
- K-means is the most popular clustering algorithm.
- Note that: it terminates at a **local optimum** if SSE is used. The **global optimum** is hard to find due to complexity.

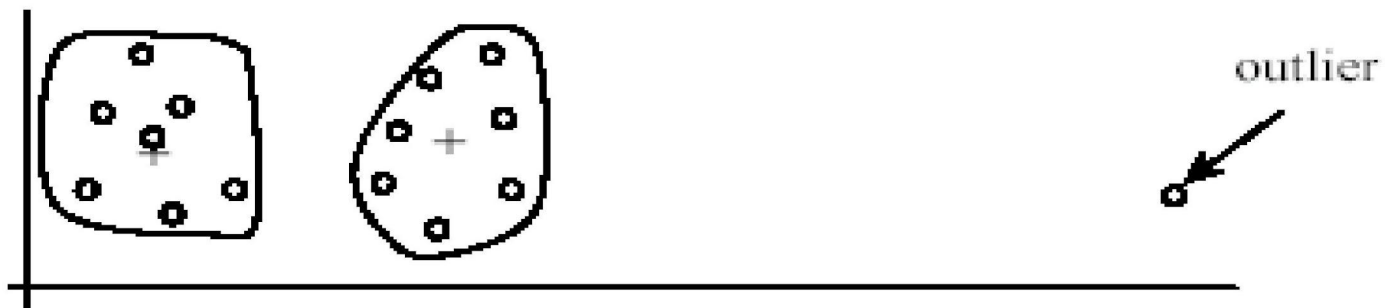
Weaknesses of K-means

- The algorithm is only applicable if the **mean** is defined.
 - For categorical data, *k*-mode - the centroid is represented by most frequent values.
- The user needs to specify *k*.
- The algorithm is sensitive to **outliers**
 - Outliers are data points that are very far away from other data points.
 - Outliers could be errors in the data recording or some special data points with very different values.

Outliers



(A): Undesirable clusters

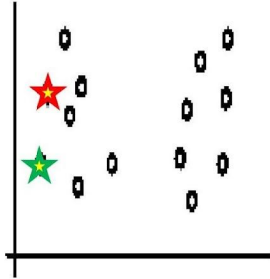


(B): Ideal clusters

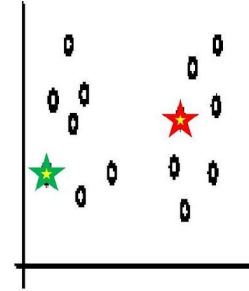
Dealing with outliers

- Remove some data points that are much further away from the centroids than other data points
 - To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.
- Perform random sampling: by choosing a small subset of the data points, the chance of selecting an outlier is much smaller
 - Assign the rest of the data points to the clusters by distance or similarity comparison, or classification

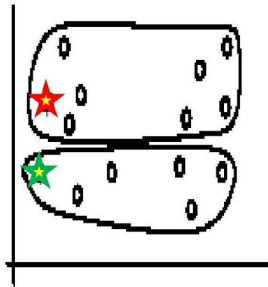
Sensitivity to initial seeds



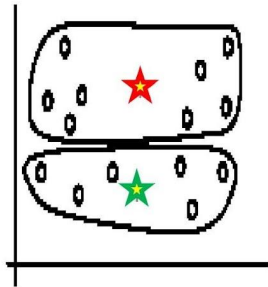
Random selection of seeds (centroids)



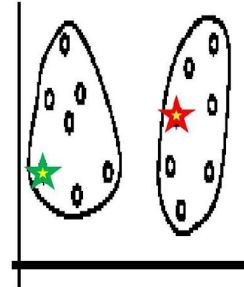
Random selection of seeds (centroids)



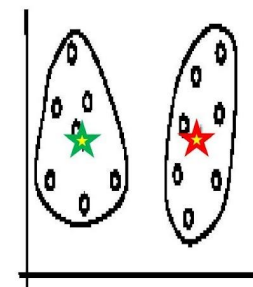
Iteration 1



Iteration 2



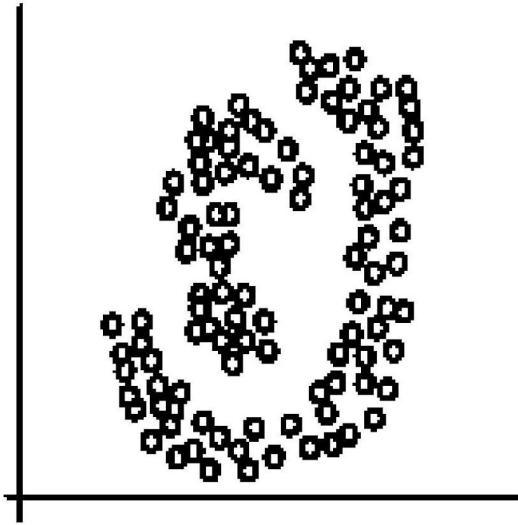
Iteration 1



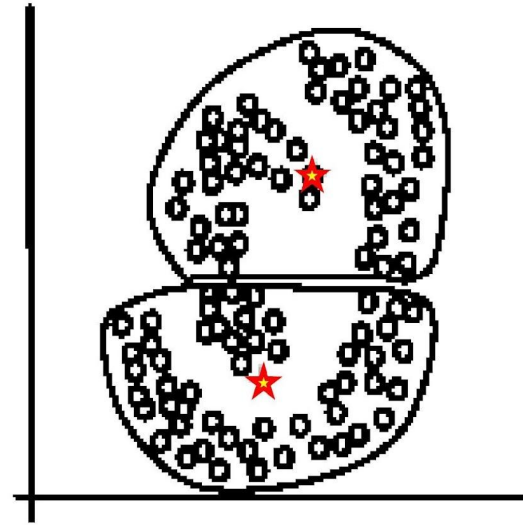
Iteration 2

Special data structures

- The k -means algorithm is not suitable for discovering clusters that are not hyper-ellipsoids (or hyper-spheres).



(A): Two natural clusters

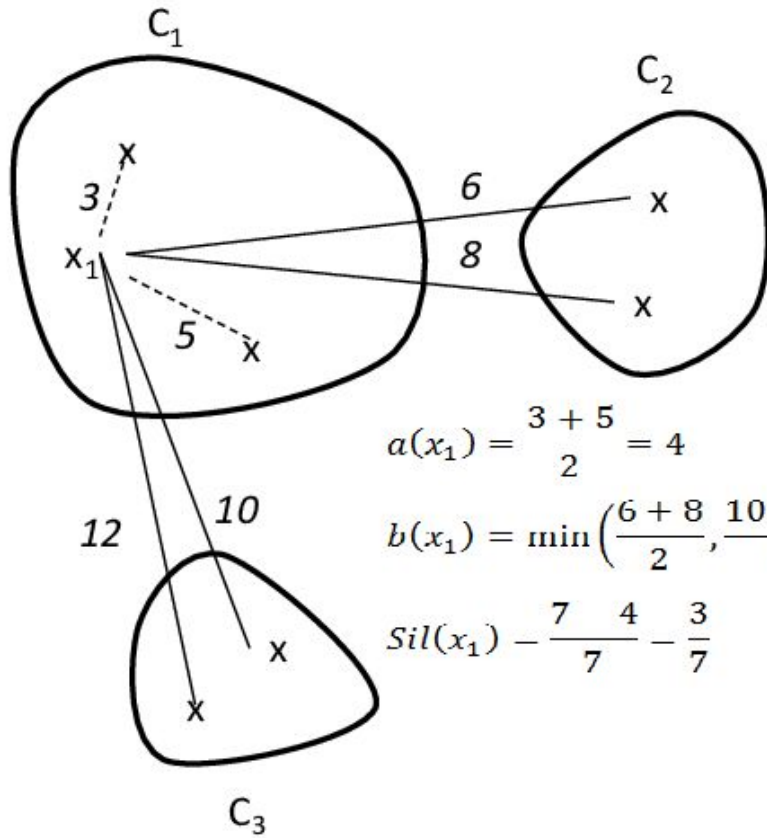


(B): k -means clusters

K-means summary

- Despite weaknesses, *k*-means is still the most popular algorithm due to its simplicity and efficiency
- No clear evidence that any other clustering algorithm performs better in general
- Comparing different clustering algorithms is a difficult task. No one knows the correct clusters!

Clustering Evaluation - Silhouette Coefficient



$$a(x_1) = \frac{3 + 5}{2} = 4$$

$$b(x_1) = \min\left(\frac{6 + 8}{2}, \frac{10 + 12}{2}\right) = 7$$

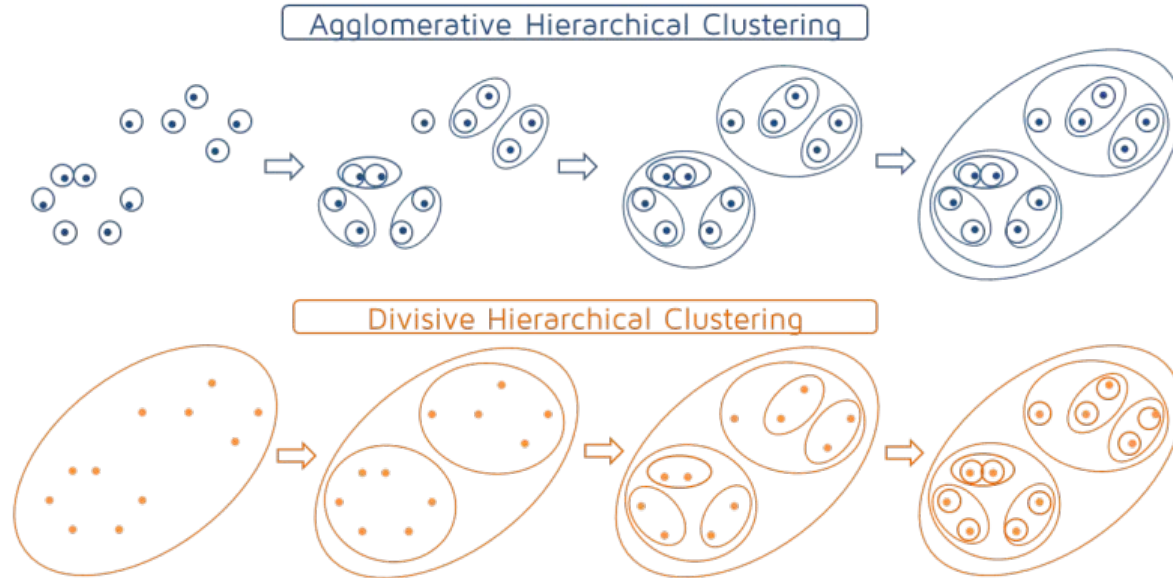
$$Sil(x_1) = \frac{7 - 4}{7} = \frac{3}{7}$$

Types of clustering:

1. Centroids-based Clustering (Partitioning methods)
2. **Connectivity-based Clustering (Hierarchical clustering)**
3. Distribution-based Clustering
4. Density-based Clustering (Model-based methods)
5. Fuzzy Clustering
6. Constraint-based (Supervised Clustering)
7. Competitive (SOM)

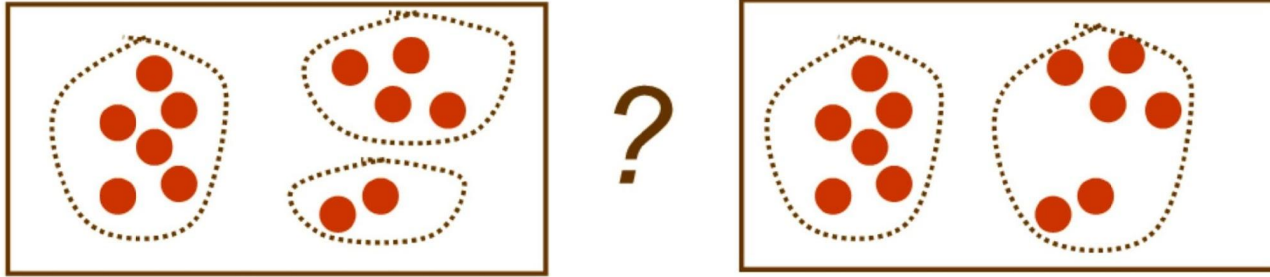
Connectivity-based Clustering (Hierarchical clustering)

- Top-down: DIANA (divisive clustering)
- Bottom-up: AGNES (agglomerative clustering)



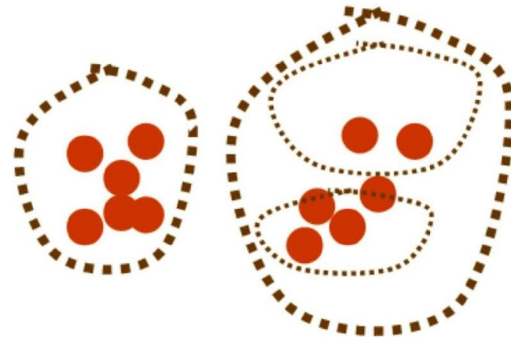
Hierarchical clustering

- Up to now, considered “flat” clustering

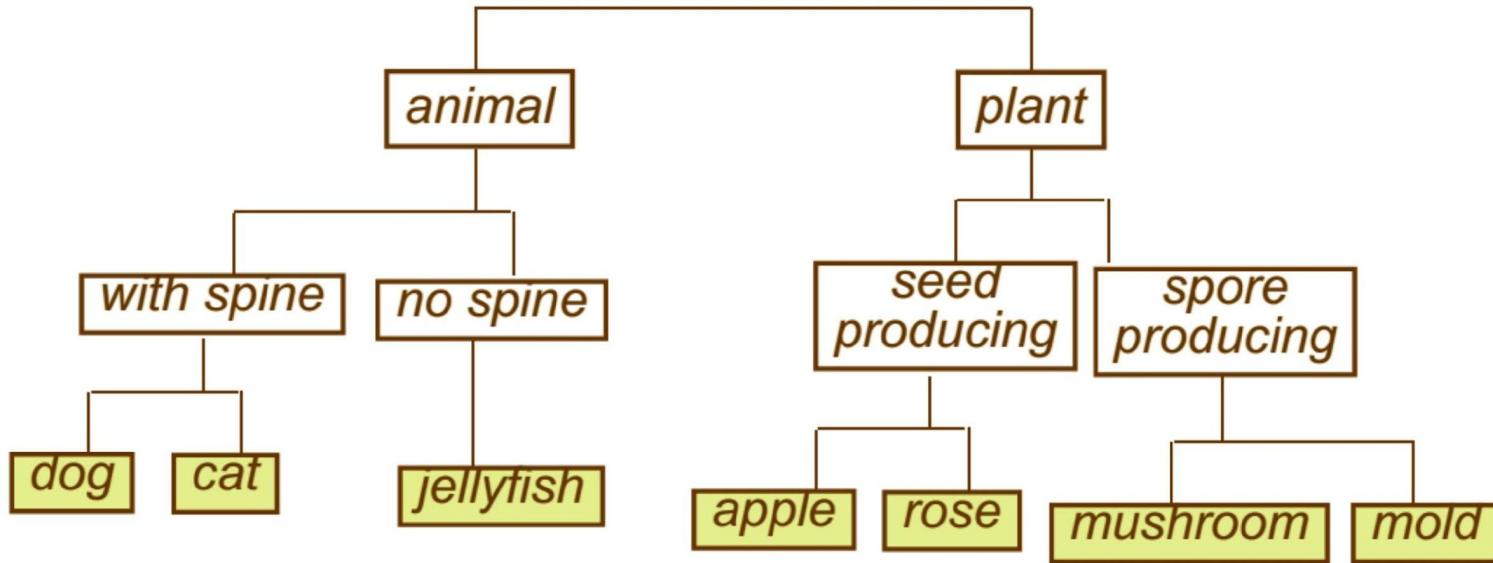


- For some data, hierarchical clustering is more appropriate than “flat” clustering

- Hierarchical clustering

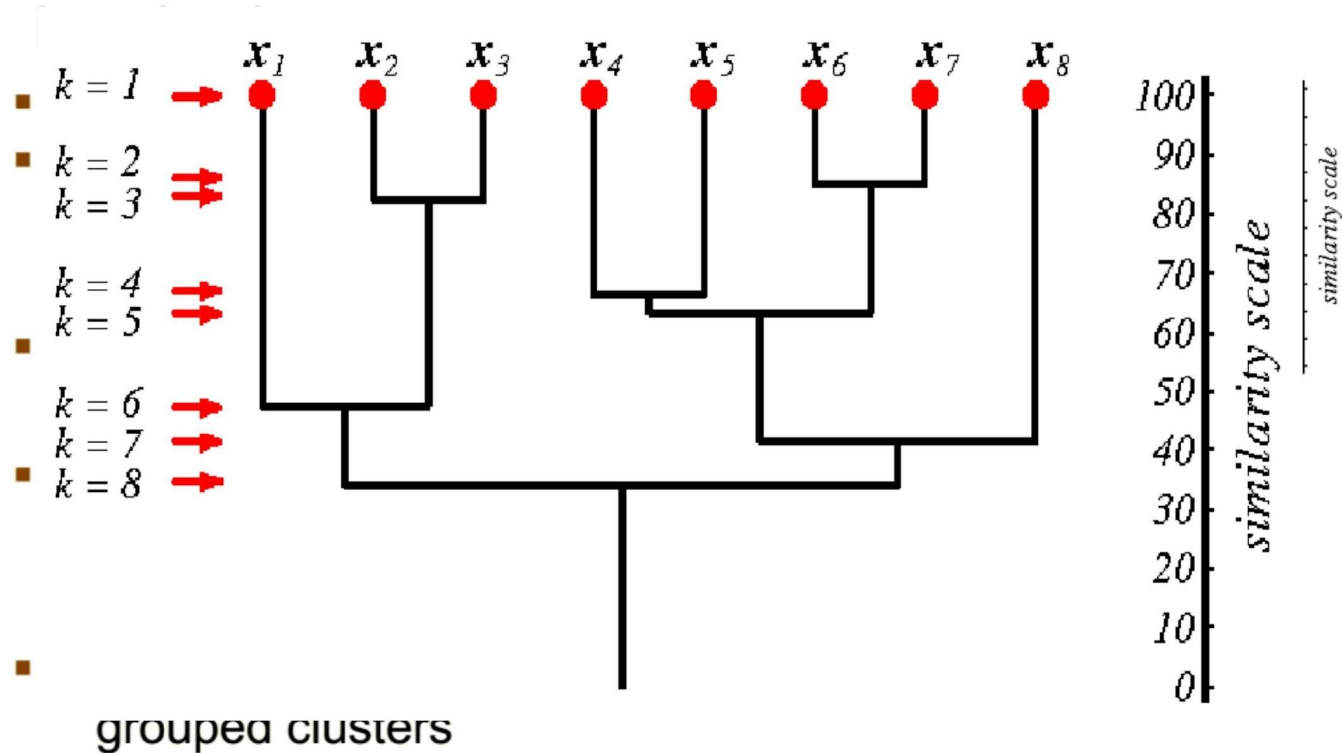


Example: biological taxonomy



A Dendrogram

- preferred way to represent a hierarchical clustering



Types of hierarchical clustering

- **Divisive (top down) clustering**

Starts with all data points in one cluster, the root, then

- Splits the root into a set of child clusters. Each child cluster is recursively divided further
- stops when only singleton clusters of individual data points remain, i.e., each cluster with only a single point

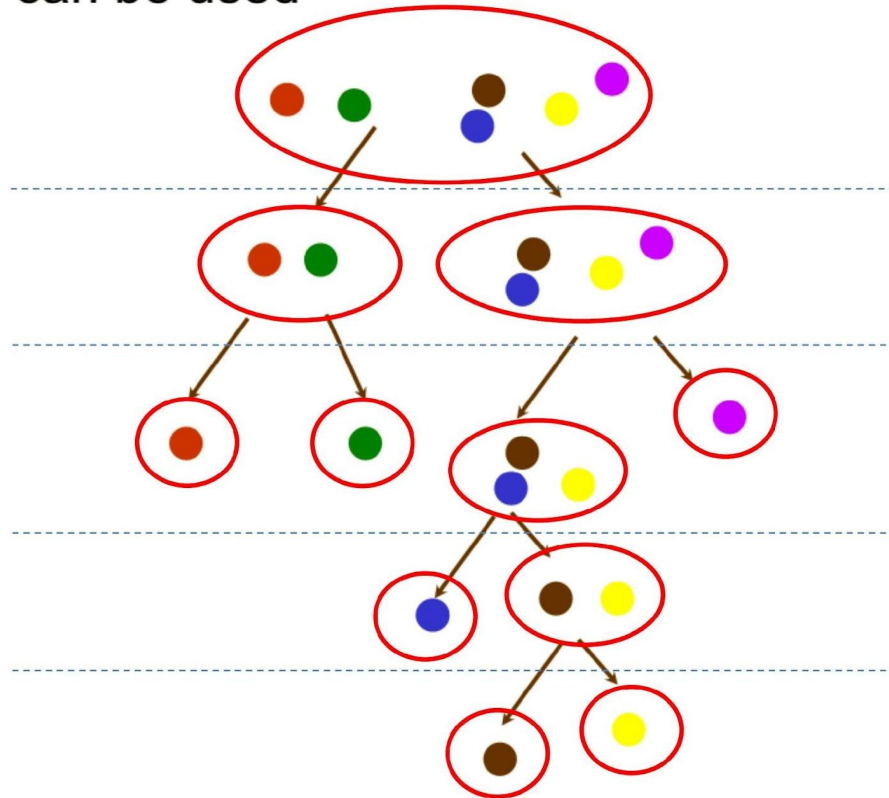
- **Agglomerative (bottom up) clustering**

The dendrogram is built from the bottom level by

- merging the most similar (or nearest) pair of clusters
- stopping when all the data points are merged into a single cluster (i.e., the root cluster).

Divisive hierarchical clustering

- Any “flat” algorithm which produces a fixed number of clusters can be used
 - set $c = 2$

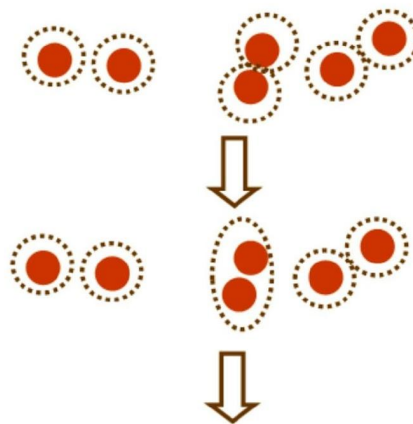


Agglomerative hierarchical clustering

initialize with each example in
singleton cluster

while there is more than **1** cluster

1. find 2 nearest clusters
2. merge them



■ Four common ways to measure cluster distance

1. minimum distance $d_{\min}(D_i, D_j) = \min_{x \in D_i, y \in D_j} \|x - y\|$

2. maximum distance $d_{\max}(D_i, D_j) = \max_{x \in D_i, y \in D_j} \|x - y\|$

3. average distance $d_{\text{avg}}(D_i, D_j) = \frac{1}{n_i n_j} \sum_{x \in D_i} \sum_{y \in D_j} \|x - y\|$

4. mean distance $d_{\text{mean}}(D_i, D_j) = \|\mu_i - \mu_j\|$

Single linkage or Nearest neighbor

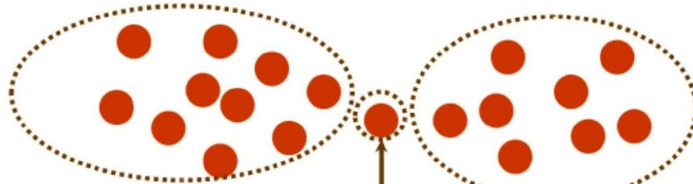
- Agglomerative clustering with minimum distance

$$d_{\min}(D_i, D_j) = \min_{x \in D_i, y \in D_j} \|x - y\|$$



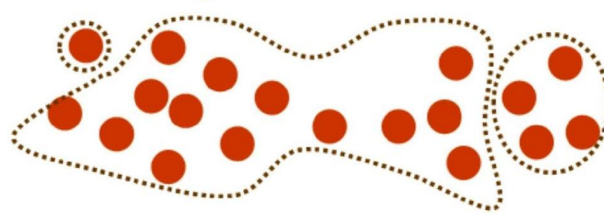
- generates minimum spanning tree
- encourages growth of elongated clusters
- disadvantage: very sensitive to noise

what we want at level with $c=3$



noisy sample

what we get at level with $c=3$

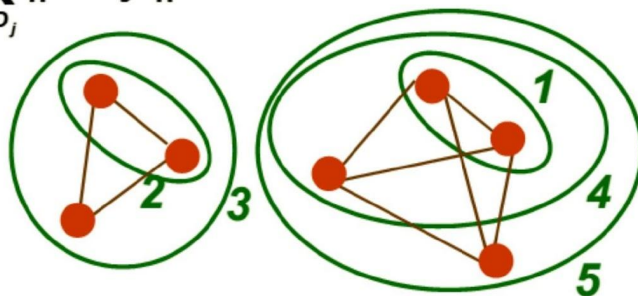


Complete linkage or Farthest neighbor

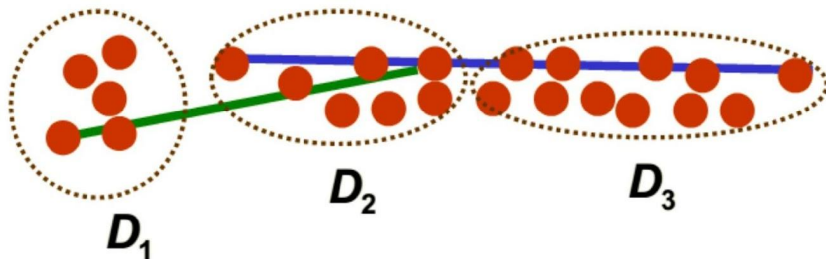
- Agglomerative clustering with maximum distance

$$d_{\max}(D_i, D_j) = \max_{x \in D_i, y \in D_j} \|x - y\|$$

- encourages compact clusters



- Does not work well if elongated clusters present



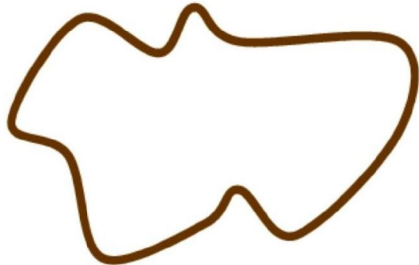
- $d_{\max}(D_1, D_2) < d_{\max}(D_2, D_3)$
- thus D_1 and D_2 are merged instead of D_2 and D_3

Divisive vs. Agglomerative

- Agglomerative is faster to compute, in general
- Divisive may be less “blind” to the global structure of the data

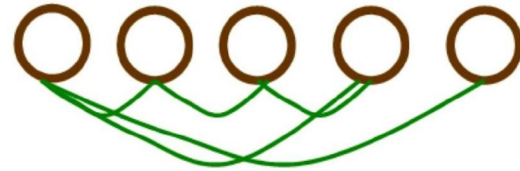
Divisive

when taking the first step (split), have access to all the data; can find the best possible split in 2 parts



Agglomerative

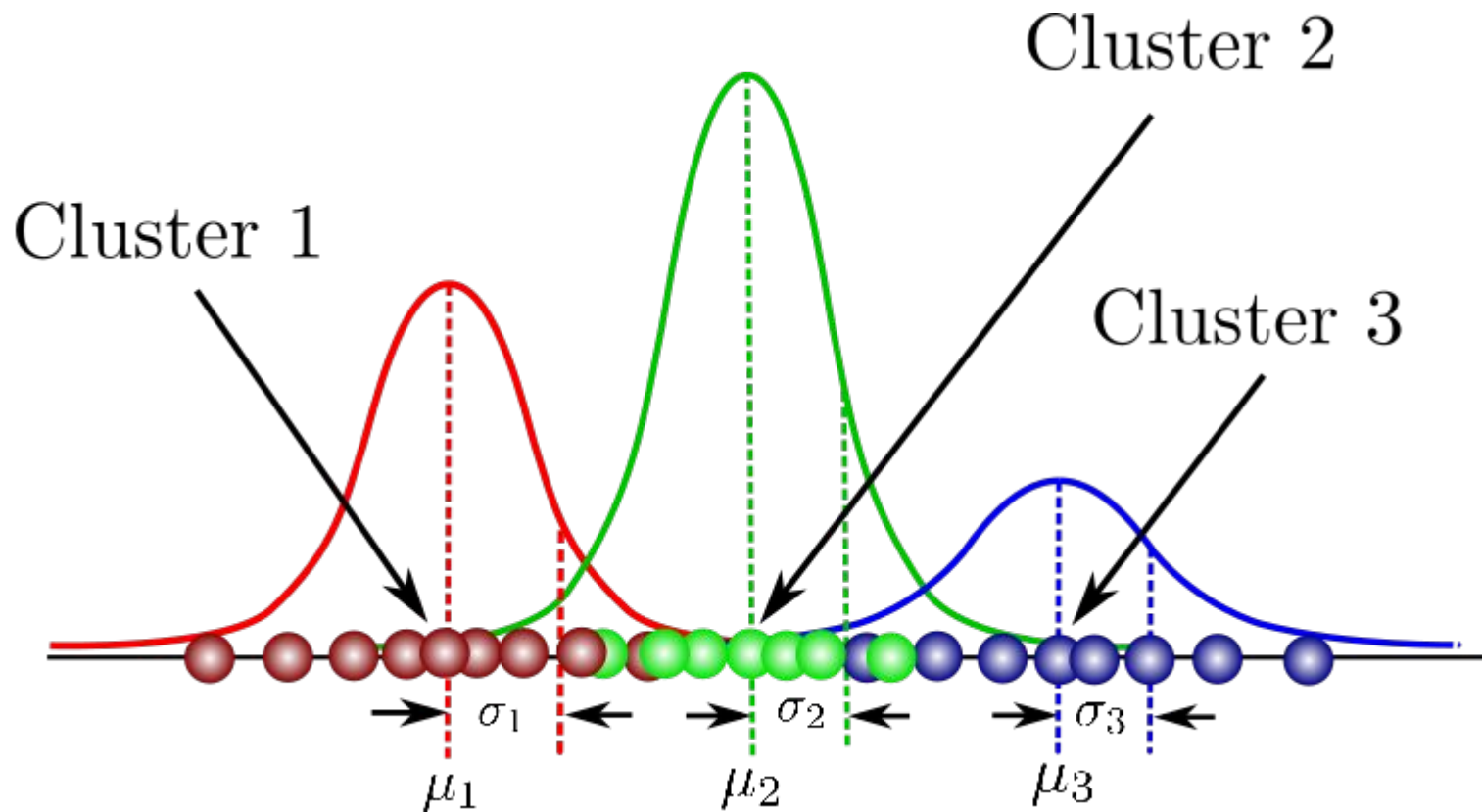
when taking the first step merging, do not consider the global structure of the data, only look at pairwise structure



Types of clustering:

1. Centroids-based Clustering (Partitioning methods)
2. Connectivity-based Clustering (Hierarchical clustering)
3. **Distribution-based Clustering**
4. Density-based Clustering (Model-based methods)
5. Fuzzy Clustering
6. Constraint-based (Supervised Clustering)
7. Competitive (SOM)

Gaussian Mixture Models (GMM)



Εύρεση παραμέτρων των γκαουσιανών με τον αλγόριθμο Expectation - Maximization (EM)

Types of clustering:

1. Centroids-based Clustering (Partitioning methods)
2. Connectivity-based Clustering (Hierarchical clustering)
3. Distribution-based Clustering
4. **Density-based Clustering (Model-based methods)**
5. Fuzzy Clustering
6. Constraint-based (Supervised Clustering)
7. Competitive (SOM)

DBSCAN

Density-Based Spatial Clustering of Applications with Noise

Basic idea of DBSCAN algorithm

The DBSCAN algorithm captures the dense regions as clusters and mainly requires two parameters for finding clusters:

1. **Epsilon:** The maximum distance (euclidean distance) between a pair of points. The two points are considered as neighbors if and only if they are separated by a distance less than or equal to epsilon.
2. **MinPoints:** The minimum number of points required to form a dense cluster.

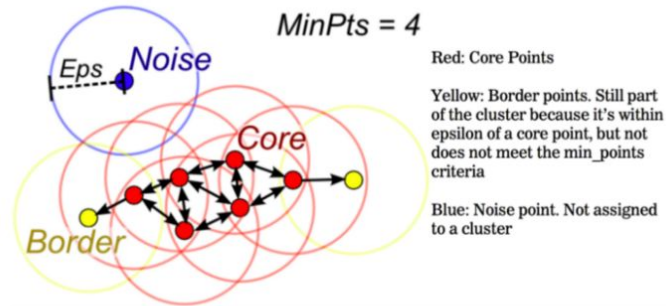
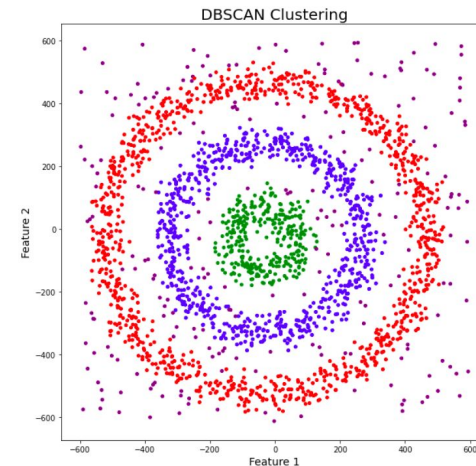


Figure 1 demonstrating density-based clustering

The MinPts = 4 means minimum 4 points are required to form a dense cluster. Also, a pair of points must be separated by a distance of less than or equal to Eps to be considered as neighbors. Based on the above two parameters, data points are classified into 3 categories as follows:

- **Core point:** A selected point is a core point if it has **at least minimum number of points (MinPts) including itself within its epsilon-neighborhood**. In figure 1, **red** points are core points that have at least MinPts=4 in their neighborhood. If we've a core point, it means it is a dense region.
- **Border point:** A selected point that is **within a neighborhood of a core point** but it itself cannot be a core point. In the figure 1, **yellow** points are identified as border points. If we've a border point, it means the point is in a vicinity or at the border of dense region.
- **Noise point:** A selected point that is neither a core point nor a border point. It means these points are **outliers** that are not associated with any dense clusters. In the figure 1, **blue** point is identified as noise point.

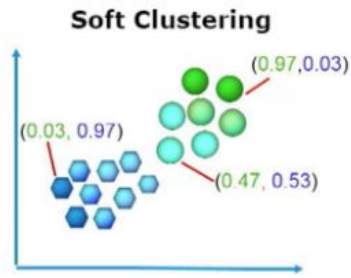


Types of clustering:

1. **Centroids-based Clustering (Partitioning methods)**
2. **Connectivity-based Clustering (Hierarchical clustering)**
3. **Distribution-based Clustering**
4. **Density-based Clustering (Model-based methods)**
5. **Fuzzy Clustering**
6. **Constraint-based (Supervised Clustering)**
7. **Competitive (SOM)**

Fuzzy C-Means

K-Means versus Fuzzy C-Means



Algorithm 1 Fuzzy C-Means Algorithm

1: Initialize $U=[u_{ij}]$ matrix, $U^{(0)}$

2: At k-step: calculate the center vectors $C^{(k)} = [c_j]$ with $U^{(k)}$

3:
$$C_i = \frac{\sum_{j=1}^n u_{ij}^m x_j}{\sum_{j=1}^n u_{ij}^m}$$

4: Update $U(k)$, $U(k+1)$

5:
$$u_{ij} = \frac{1}{\sum_{k=1}^c \left(\frac{\|x_i - c_j\|}{\|x_i - c_k\|} \right)^{\frac{2}{m-1}}}$$

6: If $\|U^{(k+1)} - U^{(k)}\| < \varepsilon$ then STOP;

7: Else

8: return to step 2.

Types of clustering:

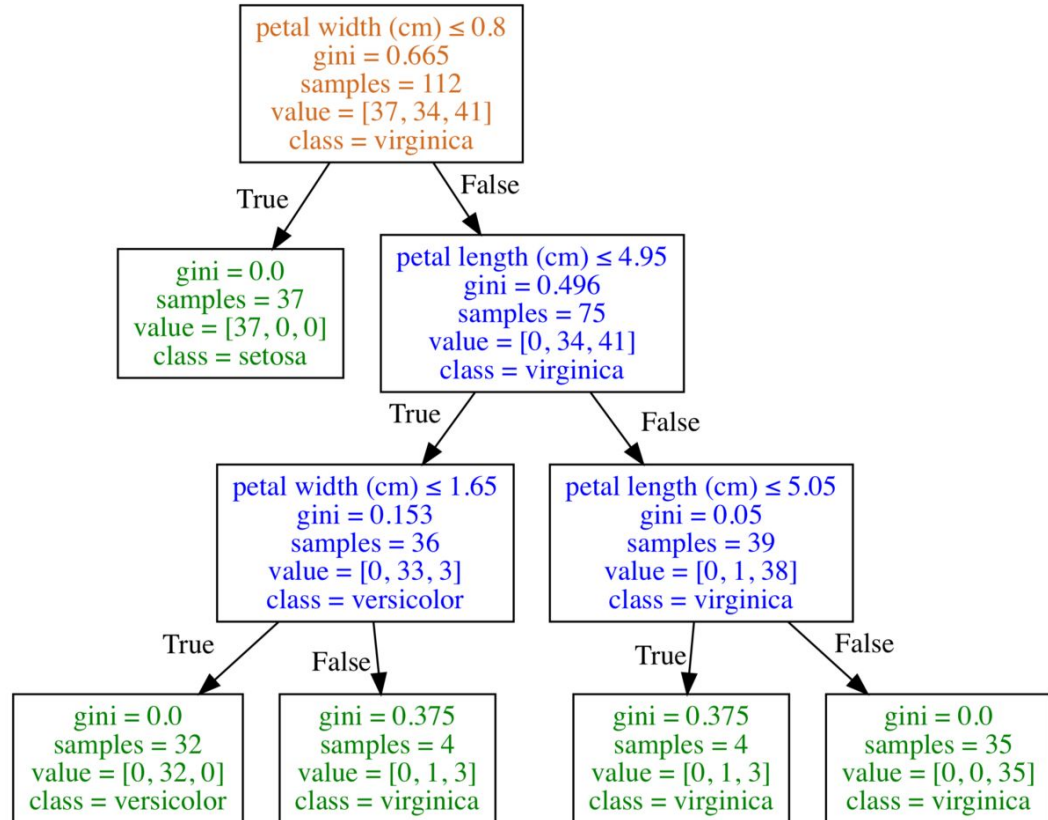
1. **Centroids-based Clustering (Partitioning methods)**
2. **Connectivity-based Clustering (Hierarchical clustering)**
3. **Distribution-based Clustering**
4. **Density-based Clustering (Model-based methods)**
5. **Fuzzy Clustering**
6. **Constraint-based (Supervised Clustering)**
7. **Competitive (SOM)**

Constrained-based (Supervised Clustering)

Decision Trees

Random Forests

Gradient Boosting (XGBoost)



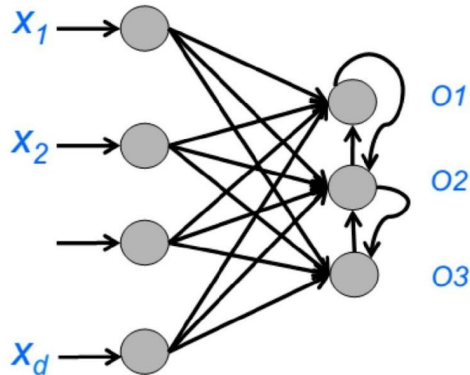
Types of clustering:

1. **Centroids-based Clustering (Partitioning methods)**
2. **Connectivity-based Clustering (Hierarchical clustering)**
3. **Distribution-based Clustering**
4. **Density-based Clustering (Model-based methods)**
5. **Fuzzy Clustering**
6. **Constraint-based (Supervised Clustering)**
7. **Competitive (SOM)**

Competitive learning

A form of unsupervised training where output units are said to be in competition for input patterns

- During training, the output unit that provides the highest activation to a given input pattern is declared the winner and is moved closer to the input pattern, whereas the rest of the neurons are left unchanged
- This strategy is also called winner-take-all since only the winning neuron is updated
 - Output units may have lateral inhibitory connections so that a winner neuron can inhibit others by an amount proportional to its activation level

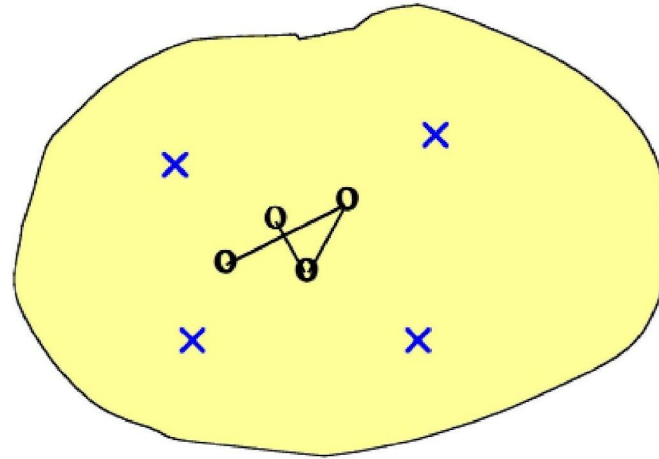


Competitive learning algorithm: Kohonen Self Organization Maps (K-SOM)

- ◆ Initialize the units to have random weights
- ◆ Repeat
 - ◆ Find the weight vector which is closest to the presented input vector. Call this the winner or the winning vector.
 - ◆ Modify the winner so as to move closer to the input vector
 - ◆ modifying weights so as to make them more similar to the values in the input vector.

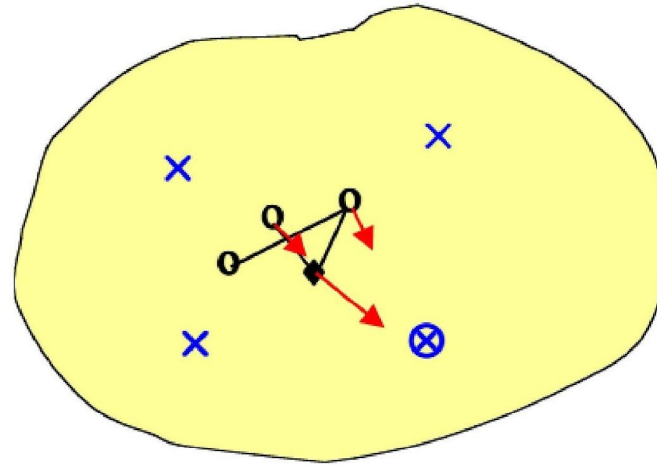
K-SOM example

- Four input data points (crosses) in 2D space.
- Four output nodes in a discrete 1D output space (mapped to 2D as circles).
- Random initial weights start the output nodes at random positions.



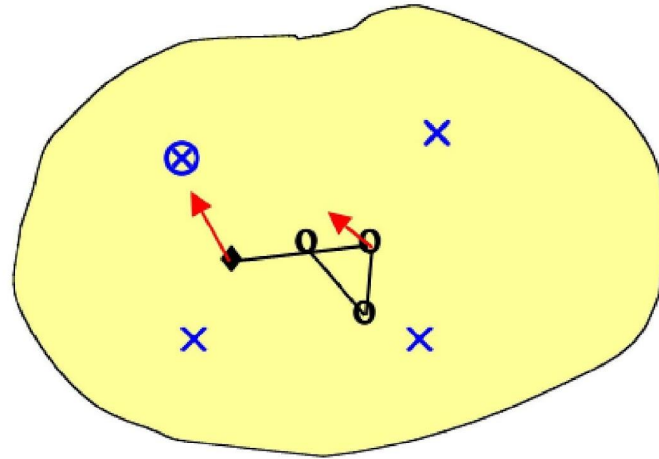
K-SOM example

- Randomly pick one input data point for training (cross in circle).
- The closest output node is the winning neuron (solid diamond).
- This winning neuron is moved towards the input data point, while its two neighbors move also by a smaller increment (arrows).



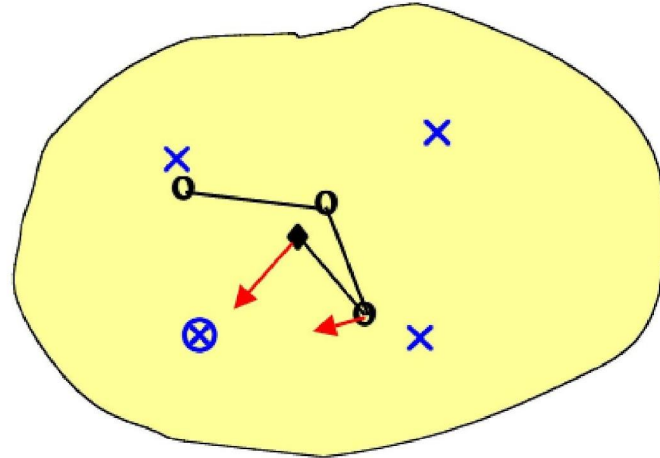
K-SOM example

- Randomly pick another input data point for training (cross in circle).
- The closest output node is the new winning neuron (solid diamond).
- This winning neuron is moved towards the input data point, while its single neighboring neuron move also by a smaller increment (arrows).



K-SOM example

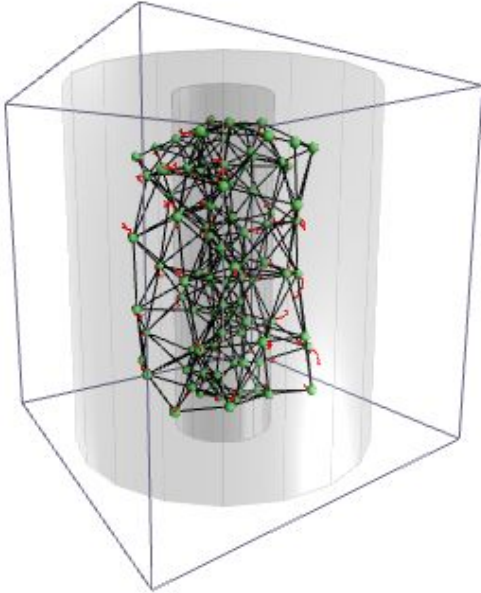
- Continue to randomly pick data points for training, and move the winning neuron and its neighbors (by a smaller increment) towards the training data points.
- Eventually, the whole output grid unravels itself to represent the input space.



DemoGNG.js

by Dr. Bernd Fritzke (fritzke@web.de)

Simulator for *competitive learning methods* and *self-organizing neural networks*

Growing Neural Gas	 <p>model:GNG-U t=377000 nodes=100 edges=381</p>	Self-Organizing Map
Growing Neural Gas with Utility		Growing Grid
Neural Gas		Competitive Hebbian Learning
Neural Gas + Competitive Hebbian Learning		LBG
Hard Competitive Learning		LBG with Utility

[full-page simulator](#), [manual](#).

Hebbian vs. Competitive learning

- H networks are used to extract information globally from the input space.
- H networks requires all weights to be updated at each epoch.
- H networks implement associative memory while C networks are selectors – only one can win!
- C networks are used to clusters similar inputs.
- C networks compete for resources.
- C networks, only the winner's weight is updated each epoch.

Note: epoch – one complete presentation of the input data to the network being trained.

Summary

- Clustering has a long history and still is in active research
 - There are a huge number of clustering algorithms, among them: Density based algorithm, Sub-space clustering, Scale-up methods, Neural networks based methods, Fuzzy clustering, Co-clustering ...
 - More are still coming every year
- Clustering is hard to evaluate, but very useful in practice
- Clustering is highly application dependent (and to some extent subjective)
- Competitive learning in neuronal networks performs clustering analysis of the input data

Clustering Method	Description	Advantages	Disadvantages	Algorithms
Partitioning methods	Based on centroids and data points are assigned into a cluster based on its proximity to the cluster centroid	Easy to implement, faster processing, can work on larger data, easy to interpret the outputs	We need to specify the number of centroids apriori, clusters that get created are of inconsistent sizes and densities, Affected by noise and outliers	k-means, k-medians, k-modes
Hierarchical Clustering	Based on top-to-bottom hierarchy of the data points to create clusters.	Easy to implement, the number of clusters need not be specified apriori, dendrograms are easy to interpret.	Cluster assignment is strict and cannot be undone, high time complexity, cannot work for a larger dataset	DIANA, AGNES, hclust etc.
Distribution-based Clustering	Based on the probability distribution of the data, clusters are derived from various metrics like mean, variance etc.	Number of clusters need not be specified apriori, works on real-time data, metrics are easy to understand and tune	Complex algorithm and slow, cannot be scaled to larger data	Gaussian Mixed Models, DBCLASD
Density-based Clustering (Model-based methods)	Based on density of the data points, also known as model based clustering	Can handle noise and outliers, need not specify number of clusters in the start, clusters that are created are highly homogenous, no restrictions on cluster shapes.	Complex algorithm and slow, cannot be scaled to larger data	DENCAST, DBSCAN
Fuzzy Clustering	Based on Partitioning Approach but data points can belong to more than one cluster	Can work on highly overlapped data, a higher rate of convergence	We need to specify the number of centroids apriori, Affected by noise and outliers, Slow algorithm and cannot be scaled	Fuzzy C Means, Rough k means
Constraint Based (Supervised Clustering)	Clustering is directed and controlled by user constraints	Creates a perfect decision boundary, can automatically determine the outcome classes based on constraints, future data can be classified based on the training boundaries	Overfitting, high level of misclassification errors, cannot be trained on larger datasets	Decision Trees, Random Forest, Gradient Boosting