

# Cluster Analysis

---

Ying Shen, SSE, Tongji University

# Cluster analysis

- Cluster analysis groups data objects based only on the attributes in the data.
- The main objective is that
  - ❖ The objects within a group be similar to one another and
  - ❖ They are different from the objects in the other groups.

# Cluster analysis

- Cluster analysis is important in the following areas:
  - ❖ Biology
  - ❖ Information retrieval
  - ❖ Medicine
  - ❖ Business

# Cluster analysis

- Cluster analysis provides an abstraction from individual data objects to the clusters in which those data objects reside.
- Some clustering techniques characterize each cluster in terms of a cluster prototype.
- The prototype is a data object that is representative of the other objects in the cluster.

# Different types of clusterings

- We consider the following types of clusterings
  - ❖ Partitional versus hierarchical
  - ❖ Exclusive versus fuzzy
  - ❖ Complete versus partial

# Partitional versus hierarchical

- A partitional clustering is a division of the set of data objects into subsets (clusters).
- A hierarchical clustering is a set of nested clusters that are organized as a tree.
- Each node (cluster) in the tree (except for the leaf nodes) is the union of its children (sub-clusters).
- The root of the tree is the cluster containing all the objects.
- Often, but not always, the leaves of the tree are singleton clusters of individual data objects.

# Partitional versus hierarchical

- The following figures form a hierarchical (nested) clustering with 1, 2, 4 and 6 clusters on each level.
- A hierarchical clustering can be viewed as a sequence of partitional clusterings.
- A partitional clustering can be obtained by taking any member of that sequence, i.e. by cutting the hierarchical tree at a certain level.

# Partitional versus hierarchical



(a) Original points.



(b) Two clusters.



(c) Four clusters.



(d) Six clusters.



## Exclusive versus fuzzy

- In an exclusive clustering, each object is assigned to a single cluster.
- However, there are many situations in which a point could reasonably be placed in more than one cluster.

# Exclusive versus fuzzy

- In a fuzzy clustering, every object belongs to every cluster with a membership weight that is between
  - ❖ 0 (absolutely does not belong) and
  - ❖ 1 (absolutely belongs).
- This approach is useful for avoiding the arbitrariness of assigning an object to only one cluster when it is close to several.
- A fuzzy clustering can be converted to an exclusive clustering by assigning each object to the cluster in which its membership value is the highest.

# Complete versus partial

- A complete clustering assigns every object to a cluster.
- A partial clustering does not assign every object to a cluster.
- The motivation of partial clustering is that some objects in a data set may not belong to well-defined groups.
- Instead, they may represent noise or outliers.

# K-means

- K-means is a prototype-based clustering technique which creates a one-level partitioning of the data objects.
- Specifically, K-means defines a prototype in terms of the centroid of a group of points.
- K-means is typically applied to objects in a continuous  $n$ -dimensional space.

# K-means

- The basic K-means algorithm is summarized below
  1. Select K points as initial centroids
  2. Repeat
    - a. Form K clusters by assigning each point to its closest centroid.
    - b. Recompute the centroid of each cluster.
  3. Until centroids do not change.

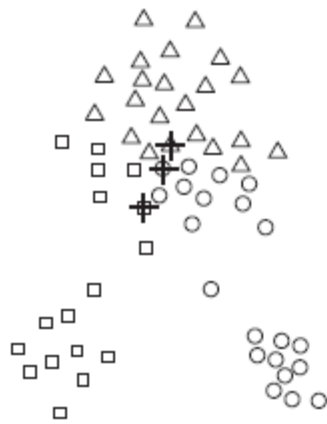
# K-means

- We first choose  $K$  initial centroids, where  $K$  is a user-defined parameter, namely, the number of clusters desired.
- Each point is then assigned to the closest centroid.
- Each collection of points assigned to a centroid is a cluster.
- The centroid of each cluster is then updated based on the points assigned to the cluster.
- We repeat the assignment and update steps until the centroids remain the same.

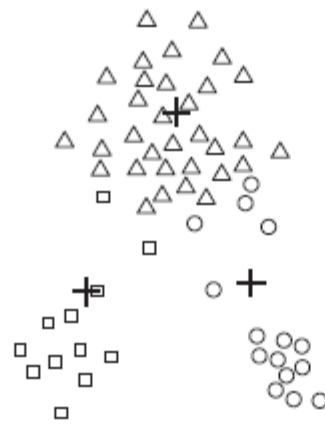
# K-means

- These steps are illustrated in the following figures.
- Starting from three centroids, the final clusters are found in four assignment-update steps.

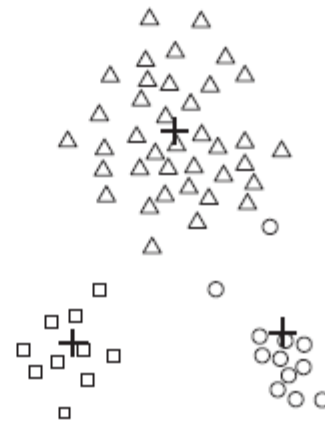
# K-means



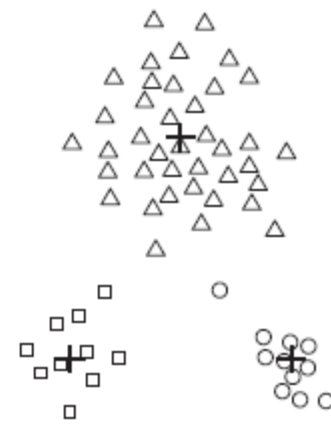
(a) Iteration 1.



(b) Iteration 2.



(c) Iteration 3.



(d) Iteration 4.



# K-means

- Each sub-figure shows
  - ❖ The centroids at the start of the iteration and
  - ❖ The assignment of the points to those centroids.
- The centroids are indicated by the “+” symbol.
- All points belonging to the same cluster have the same marker shape.

# K-means

- In the first step, points are assigned to the initial centroids, which are all in the largest group of points.
- After points are assigned to a centroid, the centroid is then updated.
- In the second step
  - ❖ Points are assigned to the updated centroids and
  - ❖ The centroids are updated again.

# K-means

- We can observe that two of the centroids move to the two small groups of points at the bottom of the figures.
- When the K-means algorithm terminates, the centroids have identified the natural groupings of points.

# Proximity measure

- To assign a point to the closest centroid, we need a proximity measure that quantifies the notion of “closest”.
- Euclidean ( $L_2$ ) distance is often used for data point in Euclidean space.

# Proximity measure

- The goal of the clustering is typically expressed by an objective function.
- Consider data whose proximity measure is Euclidean distance.
- For our objective function, which measures the quality of a clustering, we can use the sum of the squared error (SSE).

# Proximity measure

- We calculate the Euclidean distance of each data point to its closest centroid.
- We then compute the total sum of the squared distances, which is also known as the sum of the squared error (SSE).
- A small value of SSE means that the prototypes (centroids) of this clustering are a better representation of the points in their cluster.

# Proximity measure

- The SSE is defined as follows:

$$SSE = \sum_{i=1}^K \sum_{\mathbf{x} \in C_i} d(\mathbf{x}, \mathbf{c}_i)^2$$

- In this equation
  - ❖  $\mathbf{x}$  is a data object.
  - ❖  $C_i$  is the  $i$ -th cluster.
  - ❖  $\mathbf{c}_i$  is the centroid of cluster  $C_i$ .
  - ❖  $d$  is the Euclidean ( $L_2$ ) distance between two objects in Euclidean space.

# Proximity measure

- It can be shown that the mean of the data points in the cluster minimizes the SSE of the cluster.
- The centroid (mean) of the  $i$ -th cluster is defined as

$$\mathbf{c}_i = \frac{1}{m_i} \sum_{\mathbf{x} \in C_i} \mathbf{x}$$

- In this equation,  $m_i$  is the number of objects in the  $i$ -th cluster



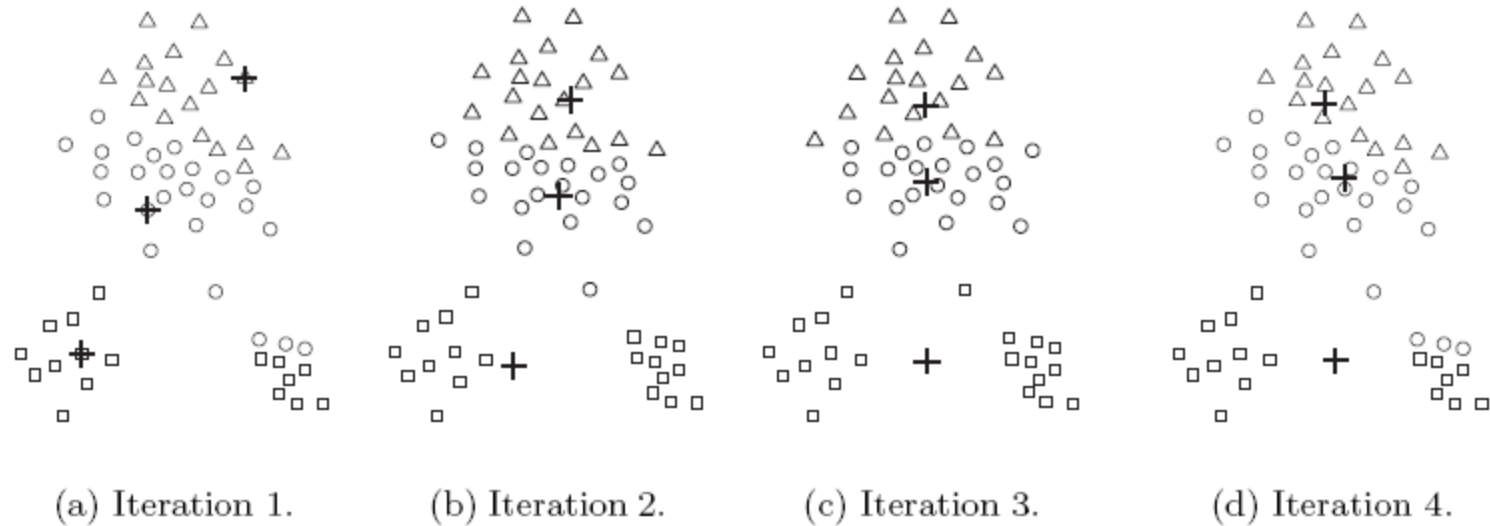
# Proximity measure

- Steps 2a and 2b of the K-means algorithm attempt to minimize the SSE.
- Step 2a forms clusters by assigning points to their nearest centroid, which minimizes the SSE for the given set of centroids.
- Step 2b recomputes the centroids so as to further minimize the SSE.

# Choosing initial centroids

- Choosing the proper initial centroids is the key step of the basic K-means procedure.
- A common approach is to choose the initial centroids randomly.
- Randomly selected initial centroids may be poor choices.
- This is illustrated in the following figures.

# Choosing initial centroids



# Choosing initial centroids

- One technique that is commonly used to address the problem of choosing initial centroids is to perform multiple runs.
- Each run uses a different set of randomly chosen initial centroids.
- We then choose the set of clusters with the minimum SSE.

# Outliers

- When the Euclidean distance is used, outliers can influence the clusters that are found.
- When outliers are present, the resulting cluster centroids may not be as representative as they otherwise would be.
- The SSE will be higher as well.
- Because of this, it is often useful to discover outliers and eliminate them beforehand.

# Outliers

- To identify the outliers, we can keep track of the contribution of each point to the SSE.
- We then eliminate those points with unusually high contributions to the SSE.
- We may also want to eliminate small clusters, since they frequently represent groups of outliers.

# Post-processing

- Two post-processing strategies that decrease the SSE by increasing the number of clusters are
  - ❖ Split a cluster
    - The cluster with the largest SSE is usually chosen.
  - ❖ Introduce a new cluster centroid
    - Often the point that is farthest from its associated cluster center is chosen.
  - ❖ We can determine this if we keep track of the contribution of each point to the SSE.

# Post-processing

- Two post-processing strategies that decrease the number of clusters, while trying to minimize the increase in total SSE, are
  - ❖ Disperse a cluster
    - This is accomplished by removing the centroid that corresponds to the cluster.
    - The points in that cluster are then re-assigned to other clusters.
    - The cluster that is dispersed should be the one that increases the total SSE the least.
  - ❖ Merge two clusters
    - We can merge the two clusters that result in the smallest increase in total SSE.



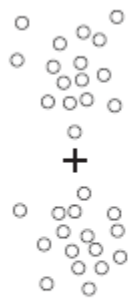
# Bisecting K-means

- Bisecting K-means algorithm is an extension of the basic K-means algorithm.
- The main steps of the algorithm are described as follows
  - ❖ To obtain K clusters, split the set of all points into two clusters.
  - ❖ Select one of these clusters to split.
  - ❖ Continue the process until K clusters have been produced.

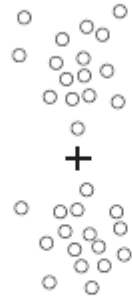
# Bisecting K-means

- There are a number of different ways to choose which cluster to split.
  - ❖ We can choose the largest cluster at each step.
  - ❖ We can also choose the one with the largest SSE.
  - ❖ We can also use a criterion based on both size and SSE.
- Different choices result in different clusters.
- We often refine the resulting clusters by using their centroids as the initial centroids for the basic K-means algorithm.
- The bisecting K-means algorithm is illustrated in the following figure.

# Bisecting K-means



(a) Iteration 1.



(b) Iteration 2.



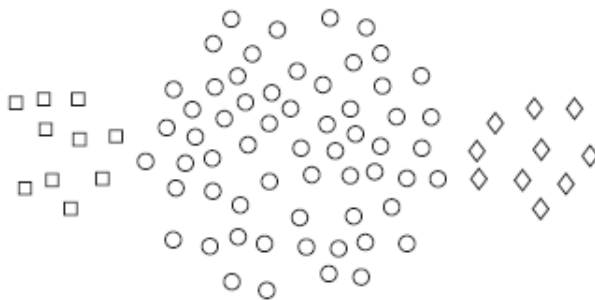
(c) Iteration 3.

# Limitations of K-means

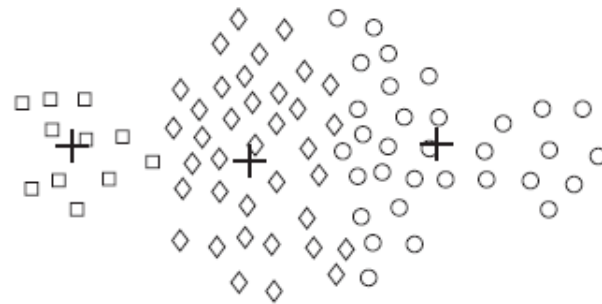
- K-means and its variations have a number of limitations with respect to finding different types of clusters.
- In particular, K-means has difficulty detecting clusters with non-spherical shapes or widely different sizes or densities.
- This is because K-means is designed to look for globular clusters of similar sizes and densities, or clusters that are well separated.
- This is illustrated in the following examples.

# Limitations of K-means

- In this example, K-means cannot find the three natural clusters because one of the clusters is much larger than the other two.
- As a result, the largest cluster is divided into sub-clusters.
- At the same time, one of the smaller clusters is combined with a portion of the largest cluster.



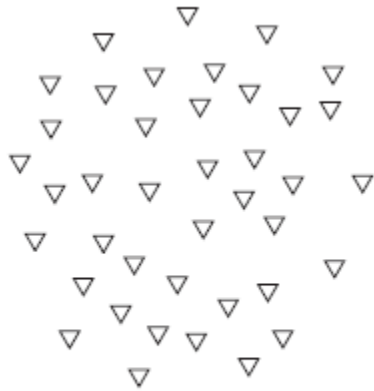
(a) Original points.



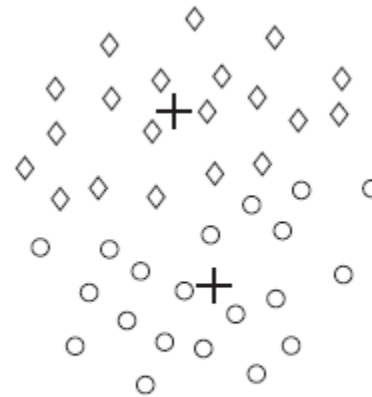
(b) Three K-means clusters.

# Limitations of K-means

- In this example, K-means fails to find the three natural clusters.
- This is because the two smaller clusters are much denser than the largest cluster.



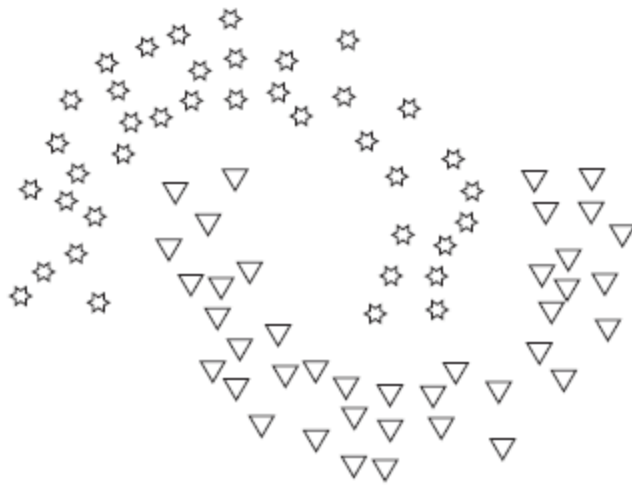
(a) Original points.



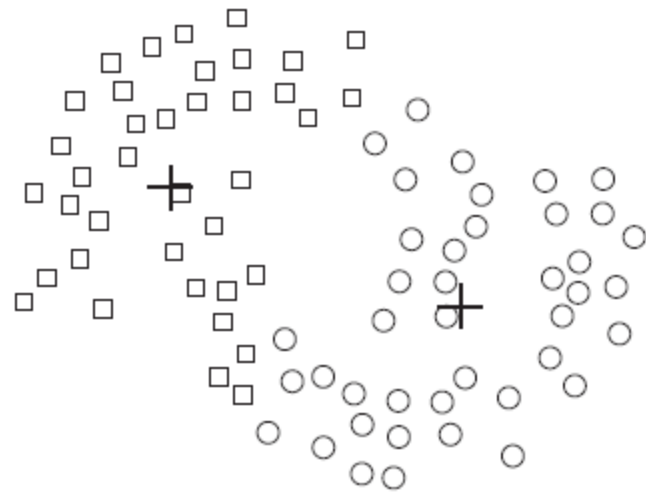
(b) Three K-means clusters.

# Limitations of K-means

- In this example, K-means finds two clusters that mix portions of the two natural clusters.
- This is because the shape of the natural clusters is not globular.



(a) Original points.



(b) Two K-means clusters.

# Hierarchical clustering

- A hierarchical clustering is a set of nested clusters that are organized as a tree.
- There are two basic approaches for generating a hierarchical clustering
  - ❖ Agglomerative
  - ❖ Divisive



# Hierarchical clustering

- In agglomerative hierarchical clustering, we start with the points as individual clusters.
- At each step, we merge the closest pair of clusters.
- This requires defining a notion of cluster proximity.

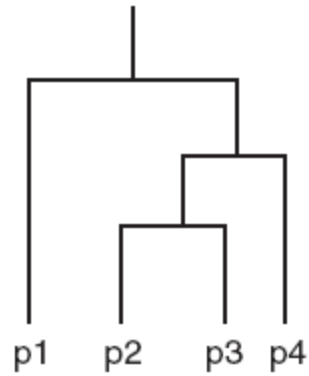
# Hierarchical clustering

- In divisive hierarchical clustering, we start with one, all-inclusive cluster.
- At each step, we split a cluster.
- This process continues until only singleton clusters of individual points remain.
- In this case, we need to decide
- Which cluster to split at each step and
- How to do the splitting.

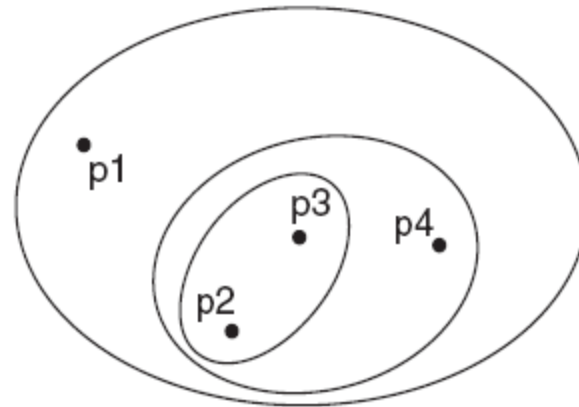
# Hierarchical clustering

- A hierarchical clustering is often displayed graphically using a tree-like diagram called the dendrogram.
- The dendrogram displays both
  - ❖ the cluster-subcluster relationships and
  - ❖ the order in which the clusters are merged (agglomerative) or split (divisive).
- For sets of 2-D points, a hierarchical clustering can also be graphically represented using a nested cluster diagram.

# Hierarchical clustering



(a) Dendrogram.



(b) Nested cluster diagram.

# Hierarchical clustering

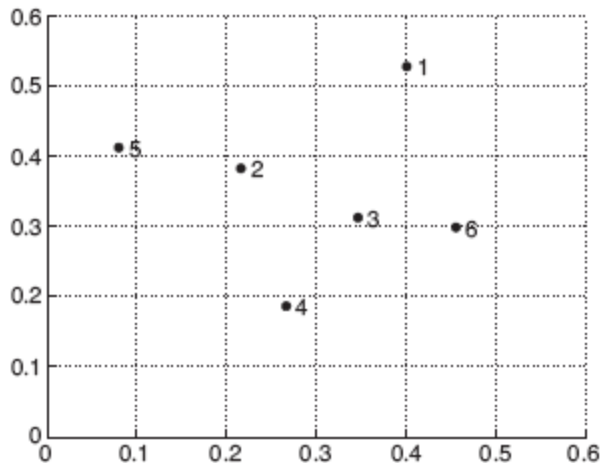
- The basic agglomerative hierarchical clustering algorithm is summarized as follows
  - ❖ Compute the proximity matrix.
  - ❖ Repeat
    - Merge the closest two clusters
    - Update the proximity matrix to reflect the proximity between the new cluster and the original clusters.
  - ❖ Until only one cluster remains

# Hierarchical clustering

- Different definitions of cluster proximity leads to different versions of hierarchical clustering.
- These versions include
  - ❖ Single link or MIN
  - ❖ Complete link or MAX
  - ❖ Group average
  - ❖ Ward's method

# Hierarchical clustering

- We consider the following set of data points.
- The Euclidean distance matrix for these data points is shown in the following slide.



Point	<i>x</i> Coordinate	<i>y</i> Coordinate
p1	0.40	0.53
p2	0.22	0.38
p3	0.35	0.32
p4	0.26	0.19
p5	0.08	0.41
p6	0.45	0.30

# Hierarchical clustering

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00



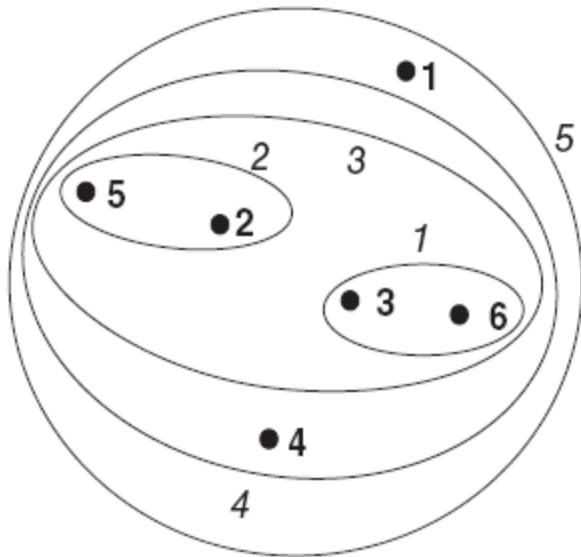
# Single link

- We now consider the single link or MIN version of hierarchical clustering.
- In this case, the proximity of two clusters is defined as the minimum of the distance between any two points in the two different clusters.
- This technique is good at handling non-elliptical shapes.
- However, it is sensitive to noise and outliers.

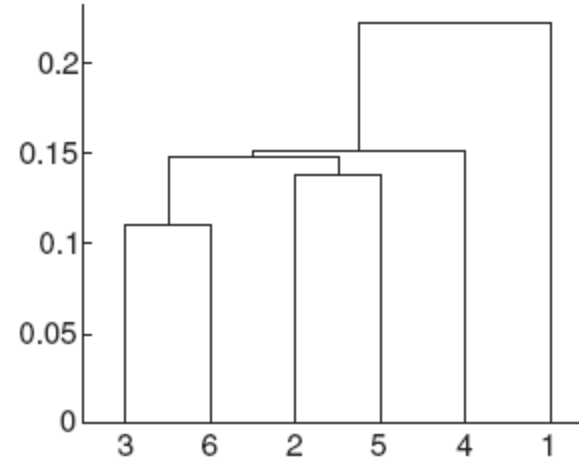
# Single link

- The following figure shows the result of applying single link technique to our example data.
- The left figure shows the nested clusters as a sequence of nested ellipses.
- The numbers associated with the ellipses indicate the order of the clustering.
- The right figure shows the same information in the form of a dendrogram.
- The height at which two clusters are merged in the dendrogram reflects the distance of the two clusters.

# Single link



(a) Single link clustering.



(b) Single link dendrogram.

# Single link

- For example, we see that the distance between points 3 and 6 is 0.11.
- That is the height at which they are joined into one cluster in the dendrogram.
- As another example, the distance between clusters {3,6} and {2,5} is

$$\begin{aligned}d(\{3,6\}, \{2,5\}) &= \min(d(3,2), d(6,2), d(3,5), d(6,5)) \\ &= \min(0.15, 0.25, 0.28, 0.39) \\ &= 0.15\end{aligned}$$

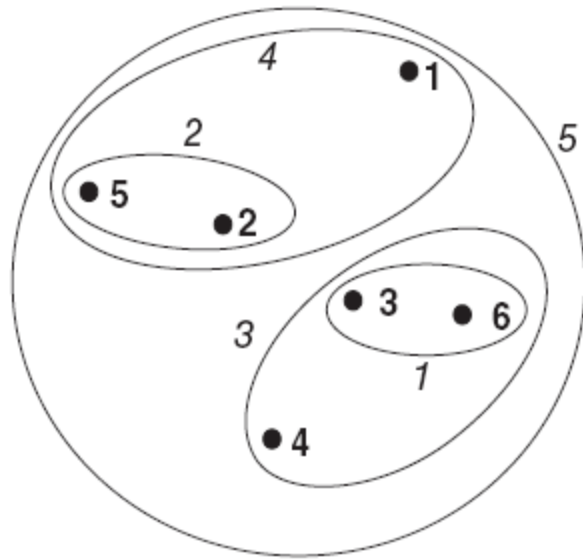
# Complete link

- We now consider the complete link or MAX version of hierarchical clustering.
- In this case, the proximity of two clusters is defined as the maximum of the distance between any two points in the two different clusters.
- Complete link is less susceptible to noise and outliers.
- However, it tends to produce clusters with globular shapes.

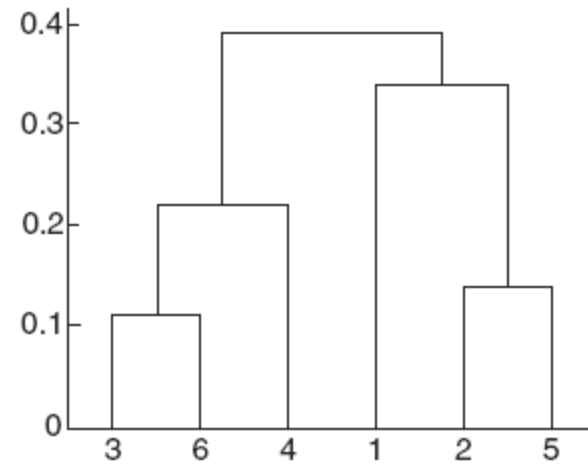
# Complete link

- The following figure shows the results of applying the complete link approach to our sample data points.
- As with single link, points 3 and 6 are merged first.
- However,  $\{3,6\}$  is merged with  $\{4\}$ , instead of  $\{2,5\}$  or  $\{1\}$ .

# Complete link



(a) Complete link clustering.



(b) Complete link dendrogram.

# Complete link

- This can be explained by the following calculations

$$\begin{aligned}d(\{3,6\}, \{4\}) &= \max(d(3,4), d(6,4)) \\ &= \max(0.15, 0.22) \\ &= 0.22\end{aligned}$$

$$\begin{aligned}d(\{3,6\}, \{2,5\}) &= \max(d(3,2), d(6,2), d(3,5), d(6,5)) \\ &= \max(0.15, 0.25, 0.28, 0.39) \\ &= 0.39\end{aligned}$$

$$\begin{aligned}d(\{3,6\}, \{1\}) &= \max(d(3,1), d(6,1)) \\ &= \max(0.22, 0.23) \\ &= 0.23\end{aligned}$$



# Group average

- We now consider the group average version of hierarchical clustering.
- In this case, the proximity of two clusters is defined as the average pairwise proximity among all pairs of points in the different clusters.
- This is an intermediate approach between the single and complete link approaches.

## Group average

- We consider two clusters  $C_i$  and  $C_j$ , which are of sizes  $m_i$  and  $m_j$  respectively.
- The distance between the two clusters can be expressed by the following equation

$$d(C_i, C_j) = \frac{\sum_{\mathbf{x} \in C_i, \mathbf{y} \in C_j} d(\mathbf{x}, \mathbf{y})}{m_i m_j}$$

# Group average

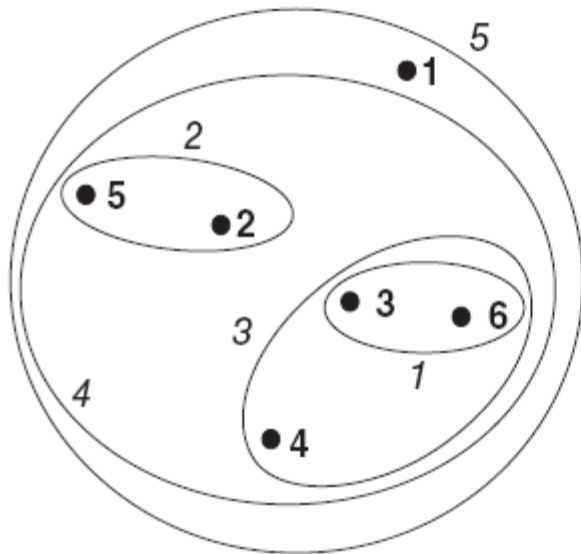
- The following figure shows the results of applying the group average to our sample data.
- The distance between some of the clusters are calculated as follows:

$$d(\{3,6,4\}, \{1\}) = \frac{0.22 + 0.37 + 0.23}{3 \times 1} = 0.27$$

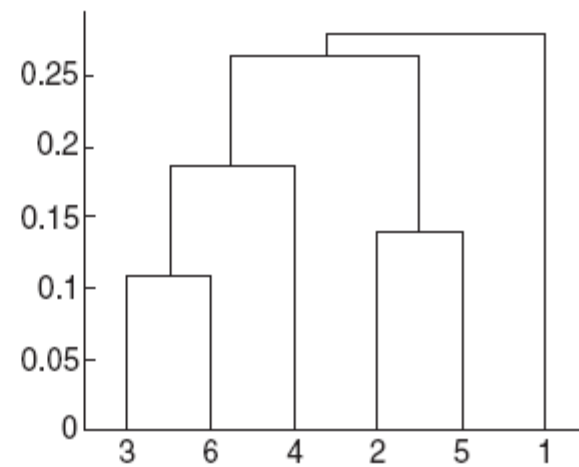
$$d(\{2,5\}, \{1\}) = \frac{0.24 + 0.34}{2 \times 1} = 0.29$$

$$d(\{3,6,4\}, \{2,5\}) = \frac{0.15 + 0.28 + 0.25 + 0.39 + 0.20 + 0.29}{3 \times 2} = 0.26$$

# Group average



(a) Group average clustering.



(b) Group average dendrogram.

## Group average

- We observe that  $d(\{3,6,4\},\{2,5\})$  is smaller than  $d(\{3,6,4\},\{1\})$  and  $d(\{2,5\},\{1\})$ .
- As a result,  $\{3,6,4\}$  and  $\{2,5\}$  are merged at the fourth stage.

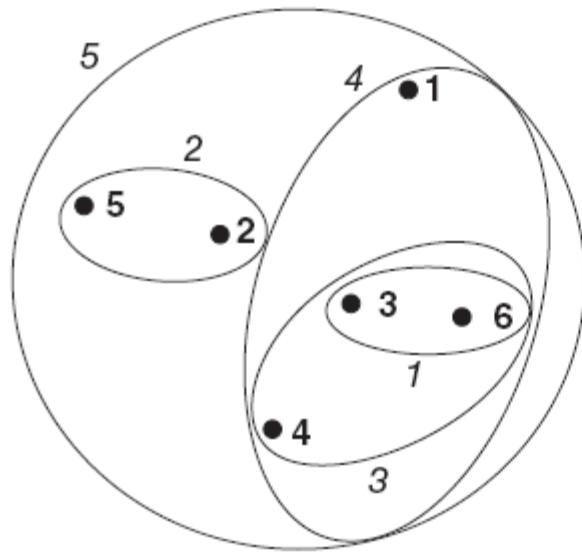
# Ward's method

- We now consider Ward's method for hierarchical clustering.
- In this case, the proximity between two clusters is defined as the increase in the sum of the squared error that results when they are merged.
- Thus, this method uses the same objective function as k-means clustering.

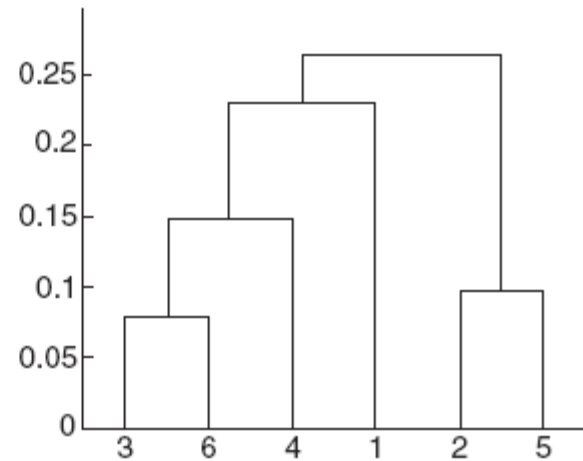
# Ward's method

- The following figure shows the results of applying Ward's method to our sample data.
- The clustering that is produced is different from those produced by single link, complete link and group average.

# Ward's method



(a) Ward's clustering.



(b) Ward's dendrogram.



## Key issues

- Hierarchical clustering is effective when the underlying application requires the creation of a multi-level structure.
- However, they are expensive in terms of their computational and storage requirements.
- In addition, once a decision is made to merge two clusters, it cannot be undone at a later time.