# Road map

- <span style="color:red">Basic concepts</span>
- K-means algorithm
- Representation of clusters
- Hierarchical clustering
- Distance functions
- Data standardization
- Summary

# Supervised learning vs. unsupervised learning

- Supervised learning: discover patterns in the data that relate data attributes with a target (class) attribute.
  - These patterns are then utilized to predict the values of the target attribute in future data instances.
- Unsupervised learning: The data have no target attribute.
  - We want to explore the data to find some intrinsic structures in them.

# Clustering

- Clustering is a technique for finding <span style="color:red">similarity groups</span> in data, called **clusters**. I.e.,
  - it groups data instances that are similar to (near) each other in one cluster and data instances that are very different (far away) from each other into different clusters.
- Clustering is often called an **unsupervised learning** task as no class values denoting an *a priori* grouping of the data instances are given, which is the case in supervised learning.
- Due to historical reasons, clustering is often considered synonymous with unsupervised learning.
  - In fact, association rule mining is also unsupervised
- This chapter focuses on clustering.

# An illustration

- The data set has three natural groups of data points, i.e., 3 natural clusters.

# What is clustering for?

- Let us see some real-life examples
- Example 1: groups people of similar sizes together to make "small" , "medium" and "large" T-Shirts.
  - Tailor-made for each person: too expensive
  - One-size-fits-all: does not fit all.
- Example 2: In marketing, segment customers according to their similarities
  - To do targeted marketing.

# What is clustering for? (cont…)

- Example 3: Given a collection of text documents, we want to organize them according to their content similarities,
  - To produce a topic hierarchy
- In fact, clustering is one of the most utilized data mining techniques.
  - It has a long history, and used in almost every field, e.g., medicine, psychology, botany, sociology, biology, archeology, marketing, insurance, libraries, etc.
  - In recent years, due to the rapid increase of online documents, text clustering becomes important.

# Aspects of clustering

- A clustering algorithm
  - Partitional clustering
  - Hierarchical clustering
  - ...

- A distance (similarity, or dissimilarity) function

- Clustering quality
  - Inter-clusters distance $\Rightarrow$ maximized
  - Intra-clusters distance $\Rightarrow$ minimized

- The quality of a clustering result depends on the algorithm, the distance function, and the application.

# Road map

- Basic concepts
- <span style="color:red">K-means algorithm</span>
- Representation of clusters
- Hierarchical clustering
- Distance functions
- Data standardization
- Summary

# K-means clustering

- K-means is a partitional clustering algorithm
- Let the set of data points (or instances) $D$ be

  $\{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_n\}$,

  where $\mathbf{x}_i = (x_{i1}, x_{i2}, \cdots, x_{ir})$ is a vector in a real-valued space $X \subseteq R^r$, and $r$ is the number of attributes (dimensions) in the data.

- 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) Randomly choose *k* 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, go to 2).

# K-means algorithm – (cont ···)

**Algorithm** $k$-means($k$, $D$)

1   Choose $k$ data points as the initial centroids (cluster centers)
2   **repeat**
3       **for** each data point $\mathbf{x} \in D$ do
4           compute the distance from $\mathbf{x}$ to each centroid;
5           assign $\mathbf{x}$ to the closest centroid          // a centroid represents a cluster
6       **endfor**
7       re-compute the centroids using the current cluster memberships
8   **until** the stopping criterion is met

# Stopping/convergence criterion

1. no (or minimum) re-assignments of data points to different clusters,

2. no (or minimum) change of centroids, or

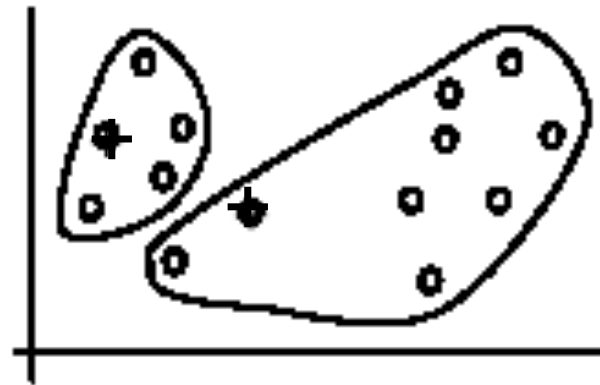3. minimum decrease in the **sum of squared error** (SSE),

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

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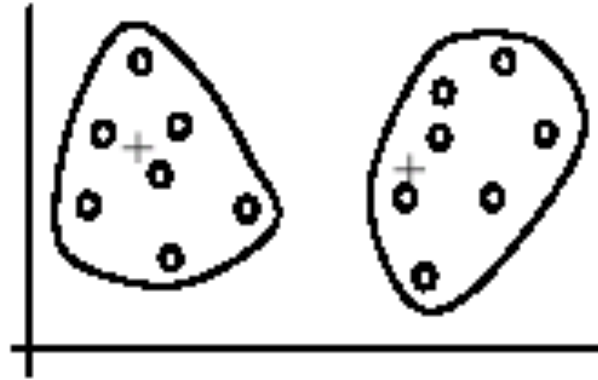
# An example



(A). Random selection of $k$ centers
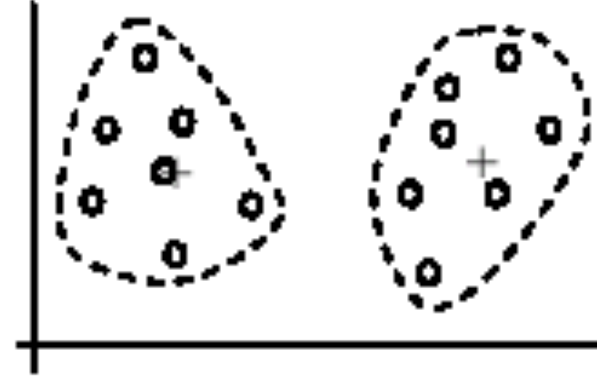
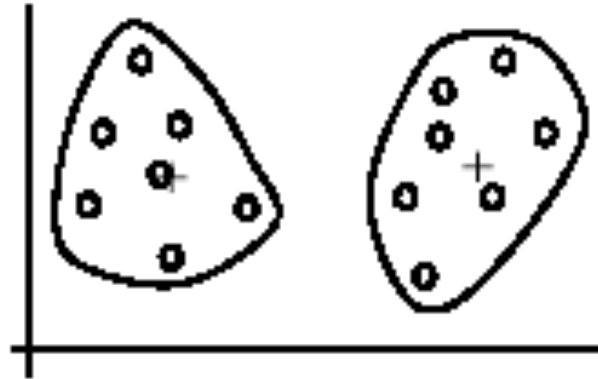*Iteration* 1: (B). Cluster assignment
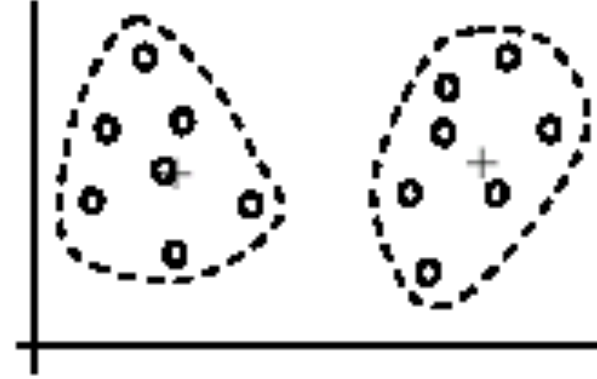
(C). Re-compute centroids

# An example (cont …)



Iteration 2: (D). Cluster assignment

(E). Re-compute centroids

Iteration 3: (F). Cluster assignment

(G). Re-compute centroids

# An example distance function

The $k$-means algorithm can be used for any application data set where the **mean** can be defined and computed. In the **Euclidean space**, the mean of a cluster is computed with:

$$\mathbf{m}_j = \frac{1}{|C_j|} \sum_{\mathbf{x}_i \in C_j} \mathbf{x}_i \qquad (2)$$

where $|C_j|$ is the number of data points in cluster $C_j$. The distance from one data point $\mathbf{x}_i$ to a mean (centroid) $\mathbf{m}_j$ is computed with

$$dist(\mathbf{x}_i, \mathbf{m}_j) = \| \mathbf{x}_i - \mathbf{m}_j \| \qquad (3)$$

$$= \sqrt{(x_{i1} - m_{j1})^2 + (x_{i2} - m_{j2})^2 + \ldots + (x_{ir} - m_{jr})^2}$$

# A disk version of *k*-means

- K-means can be implemented with data on disk
  - In each iteration, it scans the data once.
    - as the centroids can be computed incrementally
- It can be used to cluster large datasets that do not fit in main memory
- We need to control the number of iterations
  - In practice, a limited is set (< 50).
- Not the best method. There are other scale-up algorithms, e.g., BIRCH.

# A disk version of k-means (cont …)

**Algorithm** disk-$k$-means($k$, $D$)

1    Choose $k$ data points as the initial centriods $\mathbf{m}_j$, $j = 1, \ldots, k$;
2    **repeat**
3        initialize $\mathbf{s}_j = \mathbf{0}$, $j = 1, \ldots, k$;               // $\mathbf{0}$ is a vector with all 0's
4        initialize $n_j = 0$, $j = 1, \ldots, k$;               // $n_j$ is the number points in cluster $j$
5        **for** each data point $\mathbf{x} \in D$ **do**
6            $j = \arg\min_j dist(\mathbf{x}, \mathbf{m}_j)$;
7            assign $\mathbf{x}$ to the cluster $j$;
8            $\mathbf{s}_j = \mathbf{s}_j + \mathbf{x}$;
9            $n_j = n_j + 1$;
10       **endfor**
11       $\mathbf{m}_i = \mathbf{s}_j / n_j$, $i = 1, \ldots, k$;
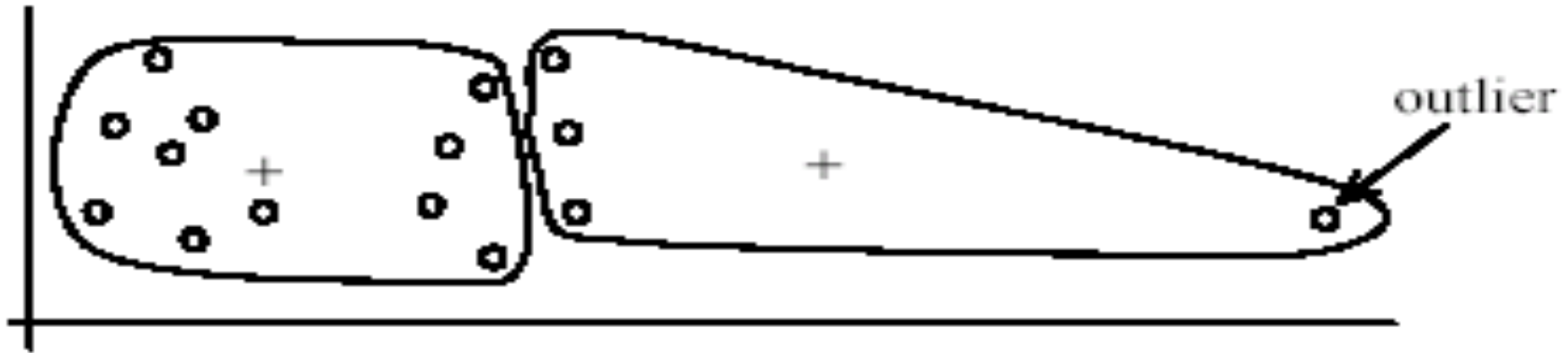12   **until** the stopping criterion is met

# Strengths of 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.
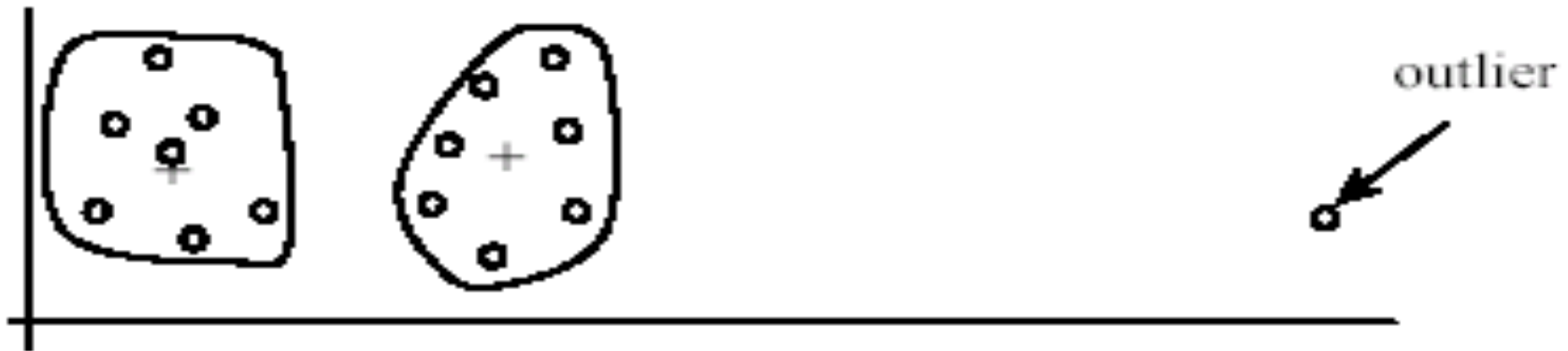
# Weaknesses of k-means

- The algorithm is only applicable if the <span style="color:red">mean</span> 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.

# Weaknesses of k-means: Problems with outliers

outlier

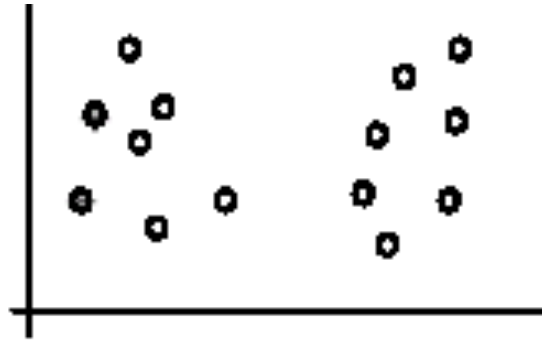(A): Undesirable clusters

outlier

(B): Ideal clusters
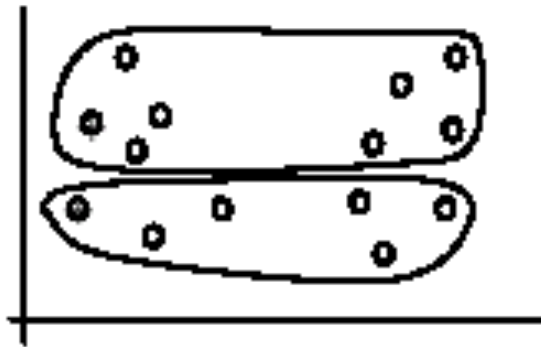
# Weaknesses of k-means: To deal with outliers

- One method is to remove some data points in the clustering process 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.
- Another method is to perform random sampling. Since in sampling we only choose a small subset of the data points, the chance of selecting an outlier is very small.
  - Assign the rest of the data points to the clusters by distance or similarity comparison, or classification
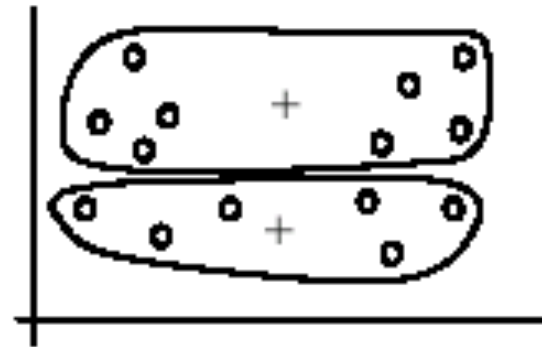
# Weaknesses of k-means (cont …)

- The algorithm is sensitive to initial seeds.
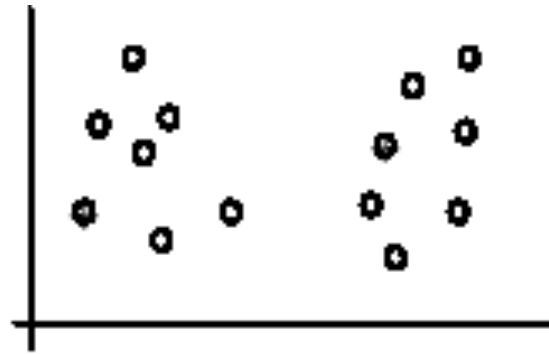


(A). Random selection of seeds (centroids)

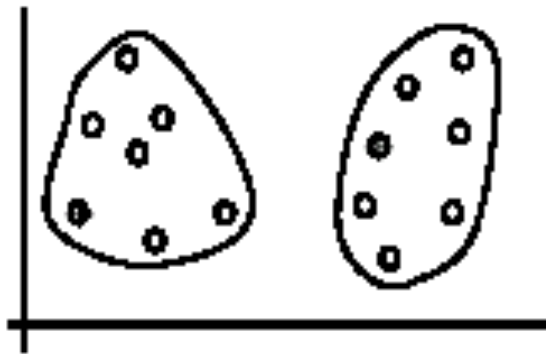(B). Iteration 1

(C). Iteration 2

# Weaknesses of k-means (cont …)

- If we use different seeds: good results



(A). Random selection of $k$ seeds (centroids)

(B). Iteration 1

(C). Iteration 2

There are some methods to help choose good seeds
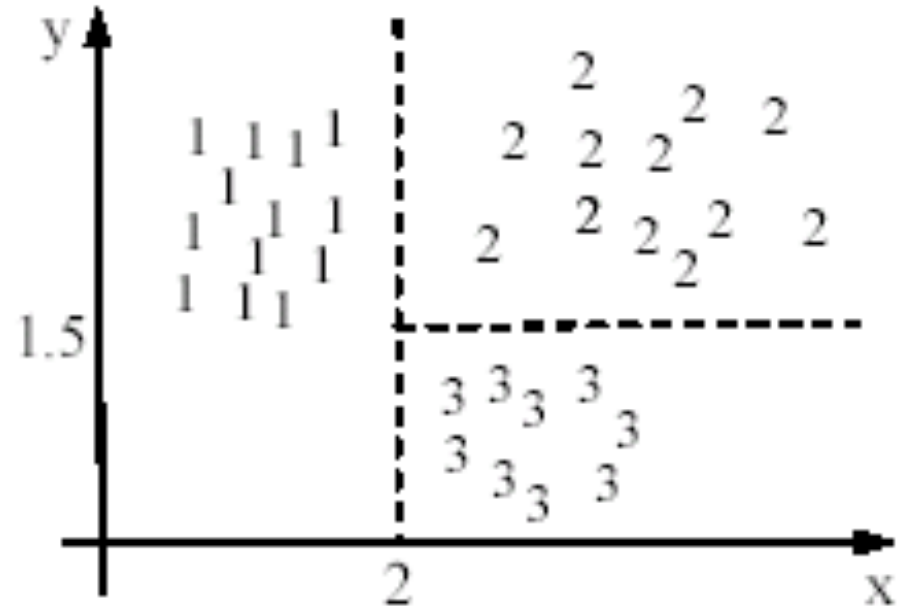
# Road map

- Basic concepts
- K-means algorithm
- <span style="color:red">Representation of clusters</span>
- Hierarchical clustering
- Distance functions
- Data standardization
- Summary

# Common ways to represent clusters

- Use the centroid of each cluster to represent the cluster.
  - compute the radius and
  - standard deviation of the cluster to determine its spread in each dimension

  - The centroid representation alone works well if the clusters are of the hyper-spherical shape.
  - If clusters are elongated or are of other shapes, centroids are not sufficient

# Using classification model

- All the data points in a cluster are regarded to have the same class label, e.g., the cluster ID.
  - run a supervised learning algorithm on the data to find a classification model.



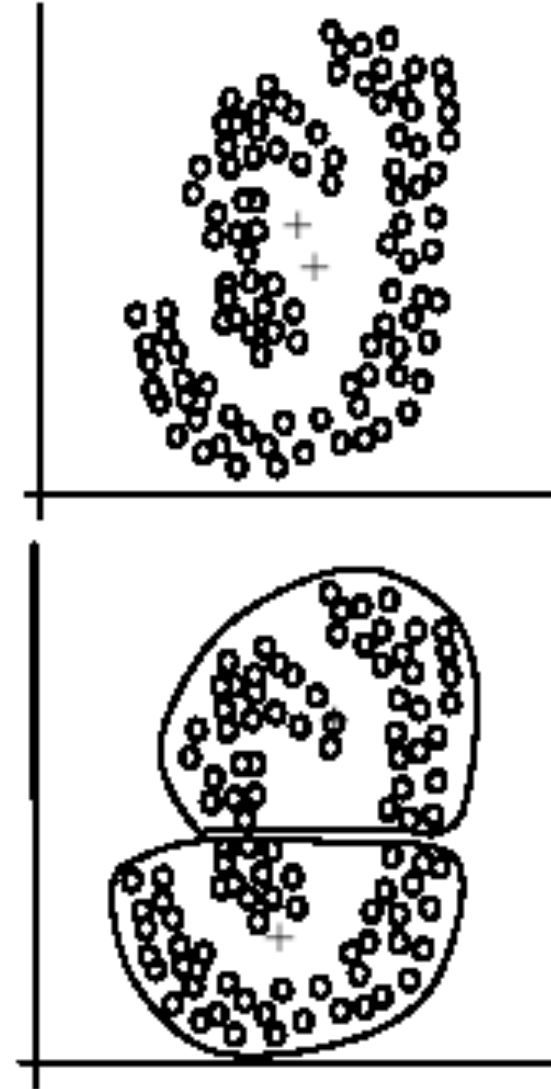$x \leq 2 \rightarrow$ cluster 1
$x > 2, y > 1.5 \rightarrow$ cluster 2
$x > 2, y \leq 1.5 \rightarrow$ cluster 3

# Use frequent values to represent cluster

- This method is mainly for clustering of categorical data (e.g., *k*-modes clustering).

- Main method used in text clustering, where a small set of frequent words in each cluster is selected to represent the cluster.

# Clusters of arbitrary shapes

- Hyper-elliptical and hyper-spherical clusters are usually easy to represent, using their centroid together with spreads.

- Irregular shape clusters are hard to represent. They may not be useful in some applications.
  - Using centroids are not suitable (upper figure) in general
  - K-means clusters may be more useful (lower figure), e.g., for making 2 size T-shirts.
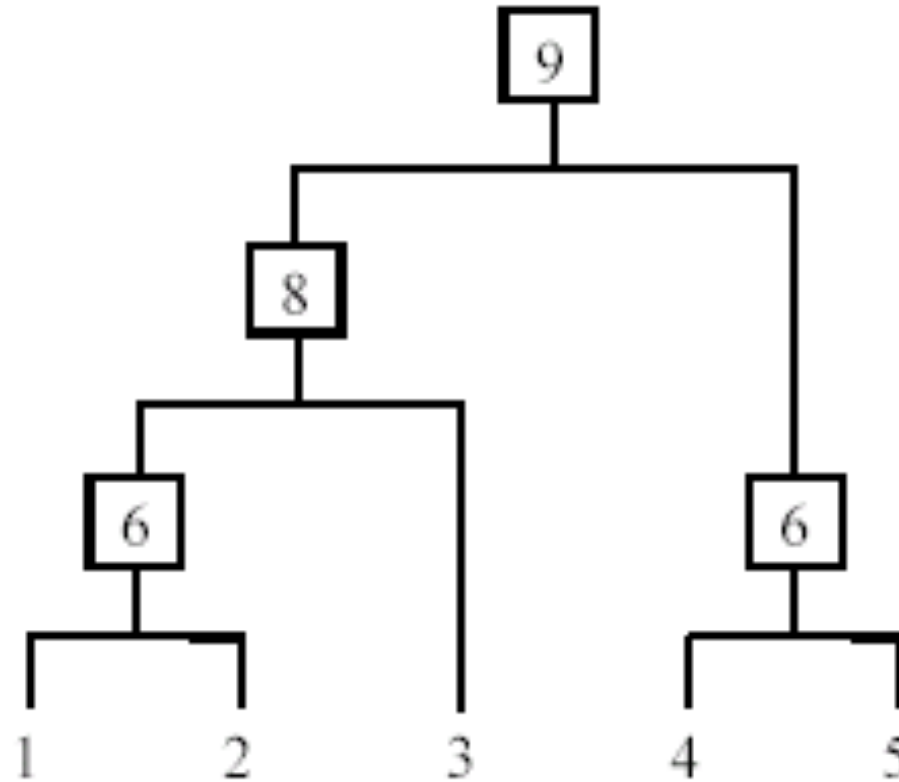
# Road map

- Basic concepts
- K-means algorithm
- Representation of clusters
- <span style="color:red">Hierarchical clustering</span>
- Distance functions
- Data standardization
- Summary

# Hierarchical Clustering

- Produce a nested sequence of clusters, a tree, also called Dendrogram.

# Types of hierarchical clustering

- Agglomerative (bottom up) clustering: It builds the dendrogram (tree) from the bottom level, and
  - merges the most similar (or nearest) pair of clusters
  - stops when all the data points are merged into a single cluster (i.e., the root cluster).
- Divisive (top down) clustering: It starts with all data points in one cluster, the root.
  - 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
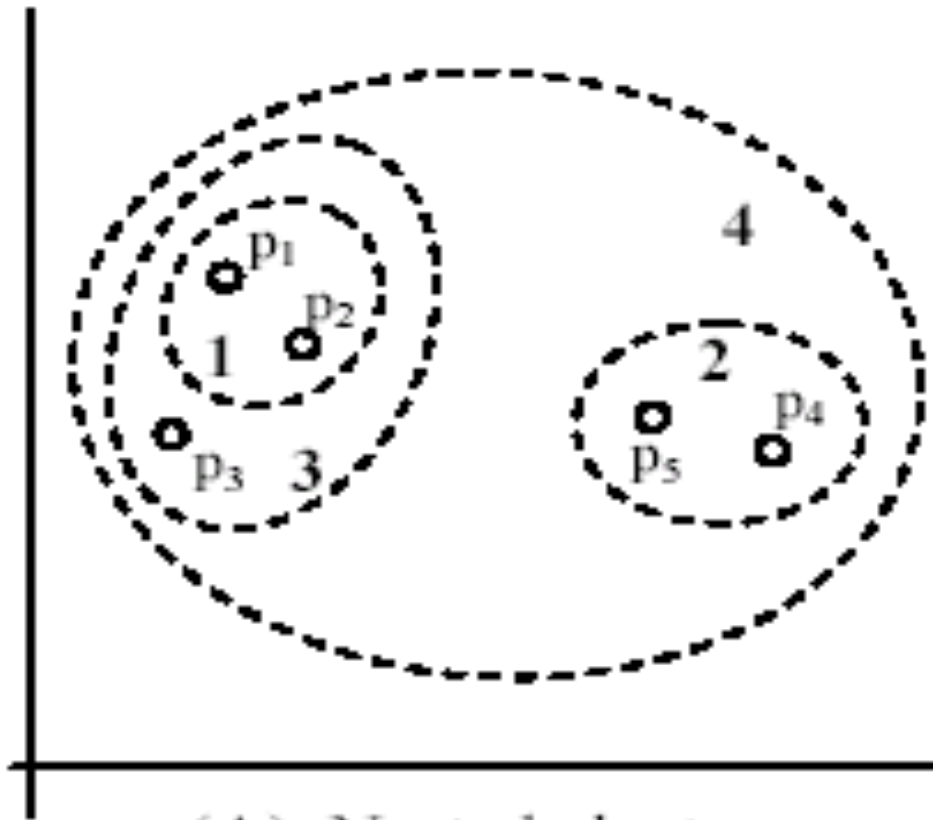
# Agglomerative clustering

It is more popular then divisive methods.

- At the beginning, each data point forms a cluster (also called a node).

- Merge nodes/clusters that have the least distance.

- Go on merging

- Eventually all nodes belong to one cluster

# An example: working of the algorithm



(A). Nested clusters

(B) Dendrogram

# Road map

- Basic concepts
- K-means algorithm
- Representation of clusters
- Hierarchical clustering
- <span style="color:red">Distance functions</span>
- Data standardization
- Summary

# Distance functions

- Key to clustering. "similarity" and "dissimilarity" can also commonly used terms.

- There are numerous distance functions for
  - Different types of data
    - Numeric data
  - Different specific applications

# Distance functions for numeric attributes

- Most commonly used functions are
  - Euclidean distance and
  - Manhattan (city block) distance
- We denote distance with: $dist(\mathbf{x}_i, \mathbf{x}_j)$, where $\mathbf{x}_i$ and $\mathbf{x}_j$ are data points (vectors)
- They are special cases of Minkowski distance. h is positive integer.

$$dist(\mathbf{x}_i, \mathbf{x}_j) = ((x_{i1} - x_{j1})^h + (x_{i2} - x_{j2})^h + ... + (x_{ir} - x_{jr})^h)^{\frac{1}{h}}$$

# Euclidean distance and Manhattan distance

- If $h = 2$, it is the <span style="color:red">Euclidean distance</span>

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \ldots + (x_{ir} - x_{jr})^2}$$

- If $h = 1$, it is the <span style="color:red">Manhattan distance</span>

$$dist(\mathbf{x}_i, \mathbf{x}_j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \ldots + |x_{ir} - x_{jr}|$$

- <span style="color:red">Weighted Euclidean distance</span>

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{w_1(x_{i1} - x_{j1})^2 + w_2(x_{i2} - x_{j2})^2 + \ldots + w_r(x_{ir} - x_{jr})^2}$$

# Road map

- Basic concepts
- K-means algorithm
- Representation of clusters
- Hierarchical clustering
- Distance functions
- <span style="color:red">Data standardization</span>
- Summary

# Data standardization

- In the Euclidean space, standardization of attributes is recommended so that all attributes can have equal impact on the computation of distances.

- Consider the following pair of data points
  - $\mathbf{x}_i$: (0.1, 20) and $\mathbf{x}_j$: (0.9, 720).

$$dist(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(0.9 - 0.1)^2 + (720 - 20)^2} = 700.000457,$$

- The distance is almost completely dominated by (720-20) = 700.

- Standardize attributes: to force the attributes to have a common value range

# Interval-scaled attributes

- Their values are real numbers following a linear scale.
  - The difference in Age between 10 and 20 is the same as that between 40 and 50.
  - The key idea is that intervals keep the same importance through out the scale
- Two main approaches to standardize interval scaled attributes, **range** and **z-score**. *f* is an attribute

$$range(x_{if}) = \frac{x_{if} - \min(f)}{\max(f) - \min(f)},$$

# Interval-scaled attributes (cont …)

- Z-score: transforms the attribute values so that they have a mean of zero and a **mean absolute deviation** of 1. The mean absolute deviation of attribute $f$, denoted by $s_f$, is computed as follows

$$s_f = \frac{1}{n}\left(\mid x_{1f} - m_f \mid + \mid x_{2f} - m_f \mid + ... + \mid x_{nf} - m_f \mid\right),$$

$$m_f = \frac{1}{n}\left(x_{1f} + x_{2f} + ... + x_{nf}\right),$$

Z-score:
$$z(x_{if}) = \frac{x_{if} - m_f}{s_f}.$$

# Summary

- Clustering is has along history and still active
  - There are a huge number of clustering algorithms
  - More are still coming every year.
- We only introduced several main algorithms. There are many others, e.g.,
  - density based algorithm, sub-space clustering, scale-up methods, neural networks based methods, fuzzy clustering, co-clustering, etc.
- Clustering is hard to evaluate, but very useful in practice. This partially explains why there are still a large number of clustering algorithms being devised every year.
- Clustering is highly application dependent and to some extent subjective.