# Data Mining

## Week 7: Clustering

**Pabitra Mitra**

Computer Science and Engineering, IIT Kharagpur

# Clustering

- Unsupervised method
- Exploratory Data Analysis

- Useful in many applications like market segment analysis
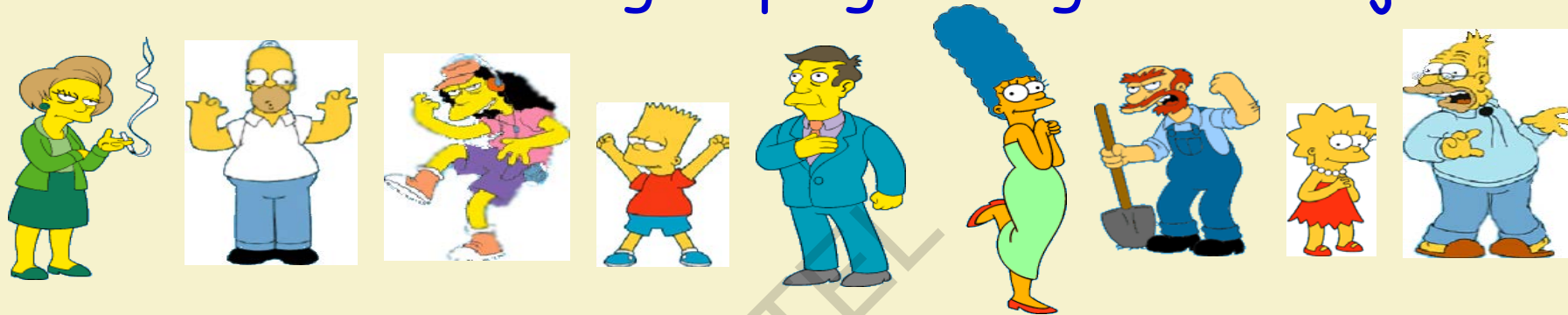
# What is clustering?

- Organizing data into classes such that there is

    - high intra-class similarity

    - low inter-class similarity

- Finding the class labels and the number of classes directly from the data (in contrast to classification).

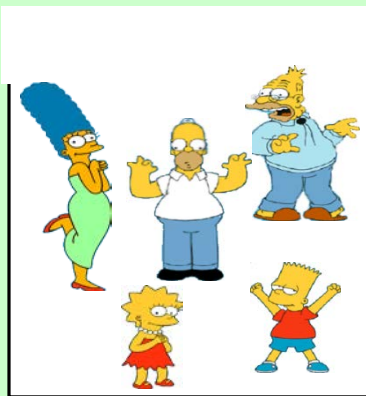- More informally, finding natural groupings among objects.

# What is a natural grouping among these objects?

# What is a natural grouping among these objects?

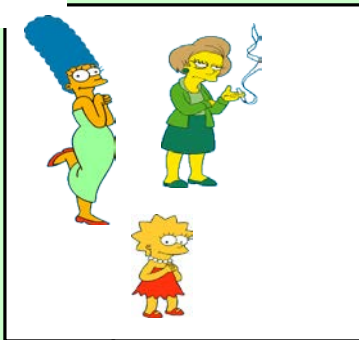

## Clustering is subjective



Simpson's Family          School Employees                    Females          Males
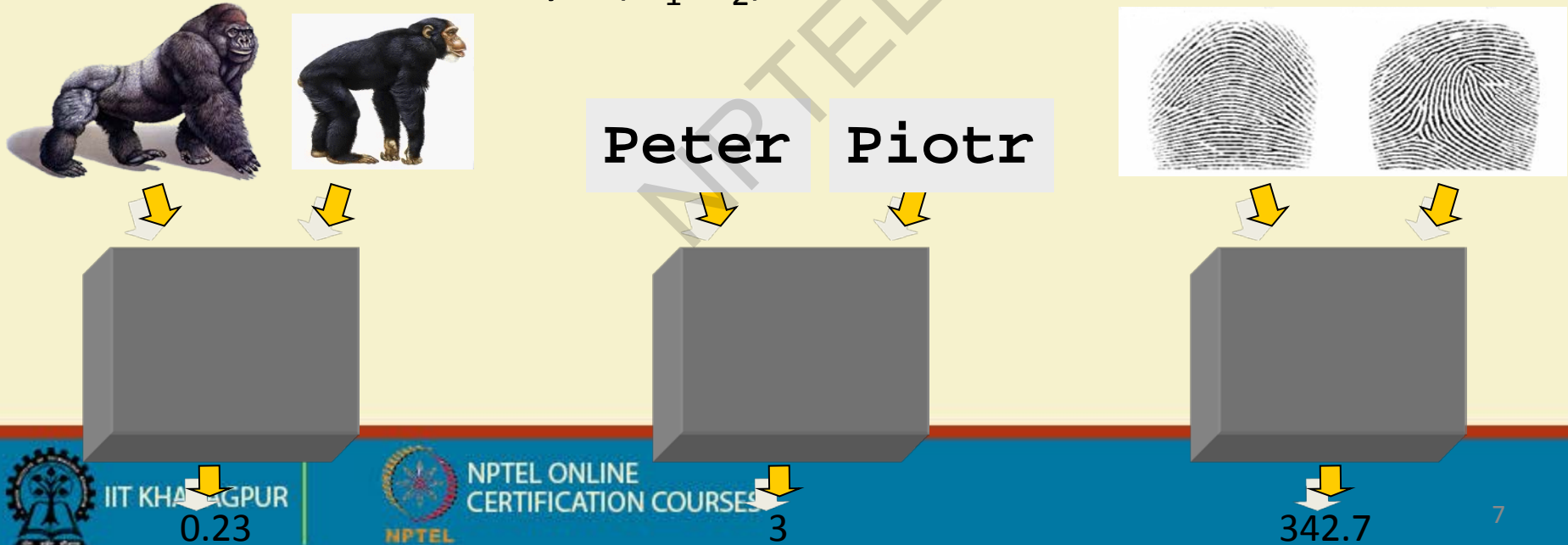
# What is similarity?

The quality or state of being similar; likeness; resemblance; as, a similarity of features.
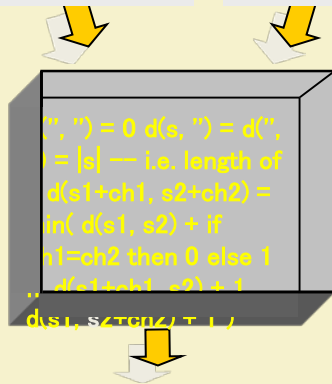


Similarity is hard to define.

# Defining distance measures

**Definition**: Let $O_1$ and $O_2$ be two objects from the universe of possible objects. The distance (dissimilarity) between $O_1$ and $O_2$ is a real number denoted by $D(O_1,O_2)$
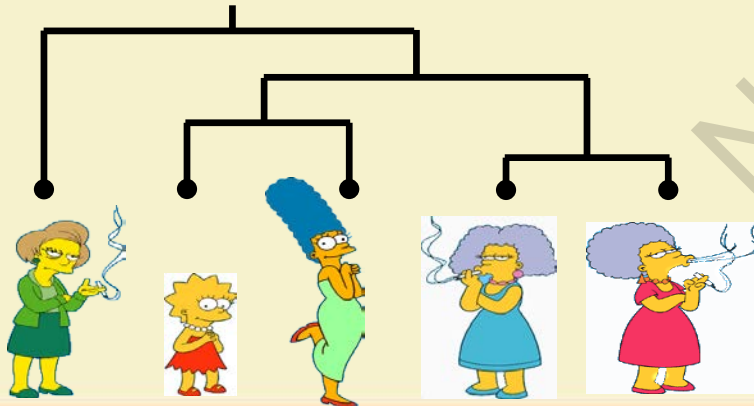


Peter Piotr

0.23

3

342.7

**Peter**  **Piotr**



```
"", ") = 0 d(s, ") = d(",
) = |s| --- i.e. length of
  d(s1+ch1, s2+ch2) =
in( d(s1, s2) + if
h1=ch2 then 0 else 1
, d(s1+ch1, s2) + 1
d(s1, s2+ch2) + 1 )
```

What properties should a distance measure have?

- $D$(A,B) = $D$(B,A)                    *Symmetry*
- $D$(A,B) = 0 iff A= B                  *Reflexive*
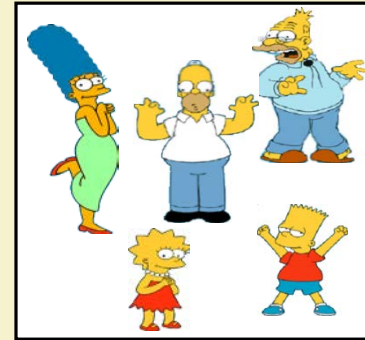- $D$(A,B) $\leq$ $D$(A,C) + $D$(B,C)  *Triangle Inequality*

# Two types of clustering

• **Partitional algorithms:** Construct various partitions and then evaluate them by some criterion

• **Hierarchical algorithms:** Create a hierarchical decomposition of the set of objects using some criterion
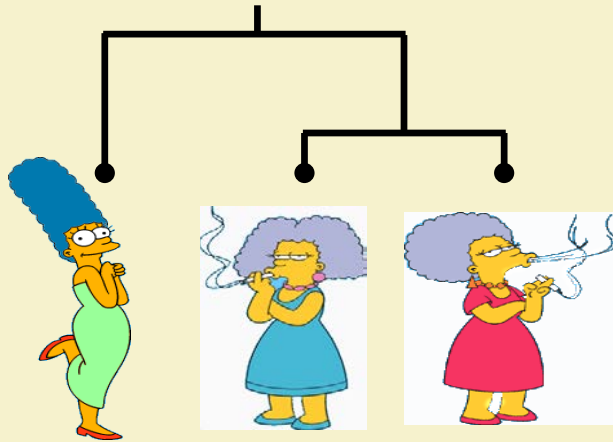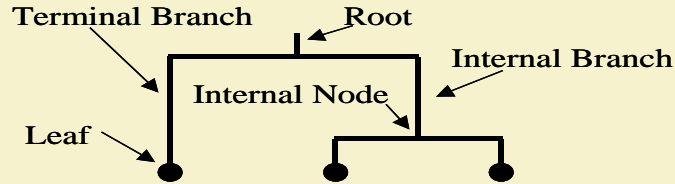
**Hierarchical**

**Partitional**
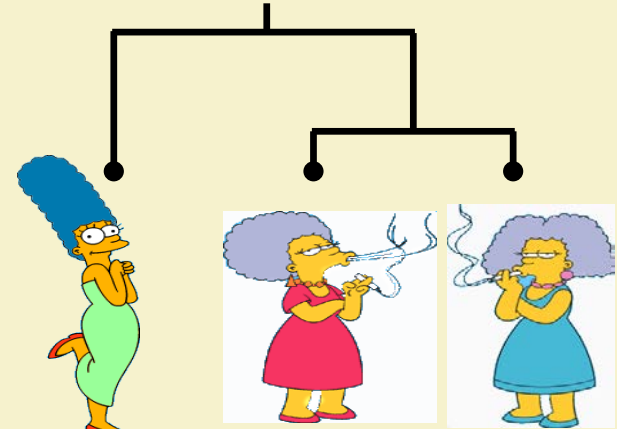
# Desirable Properties of clustering algorithm

- Scalability (in terms of both time and space)

- Ability to deal with different data types

- Minimal requirements for domain knowledge to determine input parameters

- Able to deal with noise and outliers

- Insensitive to order of input records

- Incorporation of user-specified constraints

- Interpretability and usability

# Summarizing similarity measurements

*Dendrogram.*



The similarity between two objects in a dendrogram is represented as the height of the lowest internal node they share.

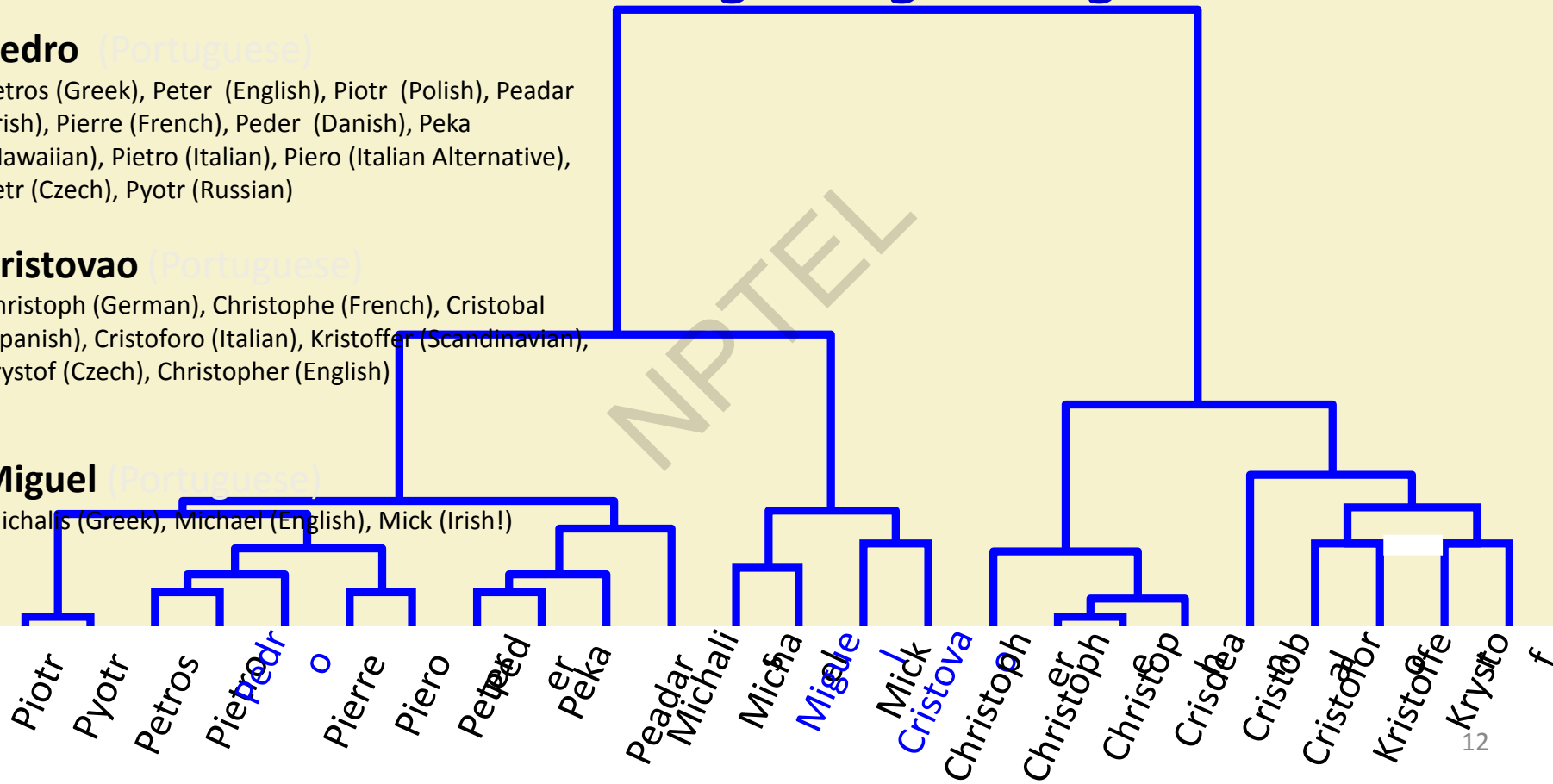# Hierarchical clustering using string edit distance

**Pedro** (Portuguese)

Petros (Greek), Peter (English), Piotr (Polish), Peadar (Irish), Pierre (French), Peder (Danish), Peka (Hawaiian), Pietro (Italian), Piero (Italian Alternative), Petr (Czech), Pyotr (Russian)

**Cristovao** (Portuguese)

Christoph (German), Christophe (French), Cristobal (Spanish), Cristoforo (Italian), Kristoffer (Scandinavian), Krystof (Czech), Christopher (English)
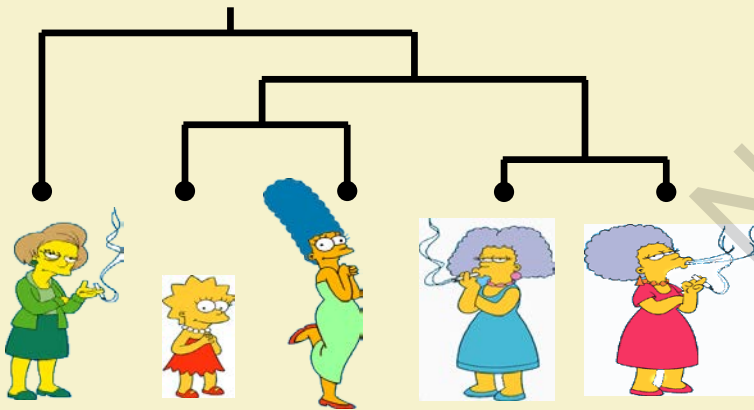
**Miguel** (Portuguese)

Michalis (Greek), Michael (English), Mick (Irish!)

# Hierarchical clustering

The number of dendrograms with $n$
leafs $= (2n - 3)! / [(2^{(n-2)})(n-2)!]$



**Bottom-Up (agglomerative):** Starting with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.

**Top-Down (divisive):** Starting with all the data in a single cluster, consider every possible way to divide the cluster into two. Choose the best division and recursively operate on both sides.

IIT KHARAGPUR

NPTEL ONLINE
CERTIFICATION COURSES

# Distance matrix

We begin with a distance matrix which contains the distances between every pair of objects in our database.

$$D(\text{🚬}, \text{👧}) = 8$$

$$D(\text{👩},\text{👩}) = 1$$

| | | | | |
|---|---|---|---|---|
| 0 | 8 | 8 | 7 | 7 |
| | 0 | 2 | 4 | 4 |
| | | 0 | 3 | 3 |
| | | | 0 | 1 |
| | | | | 0 |

# Bottom-Up (agglomerative)

Starting with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.

Consider all possible merges…

Choose the best

# Bottom-Up (agglomerative)

Starting with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.

Consider all possible merges…

Choose the best

Consider all possible merges…

Choose the best

# Bottom-Up (agglomerative)

Starting with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.



Consider all possible merges…    …    Choose the best

Consider all possible merges…    …    Choose the best

Consider all possible merges…    …    Choose the best

# Bottom-Up (agglomerative)

Starting with each item in its own cluster, find the best pair to merge into a new cluster. Repeat until all clusters are fused together.
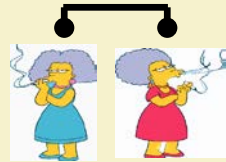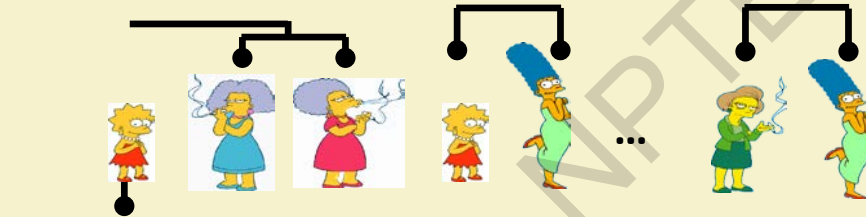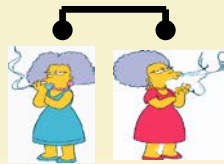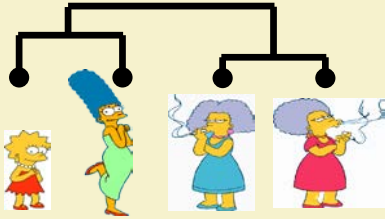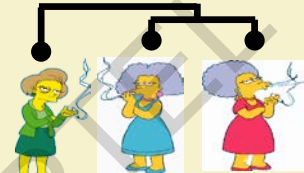
Consider all possible merges… … Choose the best

Consider all possible merges… … Choose the best

Consider all possible merges… … Choose the best

# Extending distance measure to clusters

We the distance between two objects, defining the distance between an object and a cluster, or defining the distance between two clusters:

- **Single linkage (nearest neighbor):** In this method the distance between two clusters is determined by the distance of the two closest objects (nearest neighbors) in the different clusters.
- **Complete linkage (farthest neighbor):** In this method, the distances between clusters are determined by the greatest distance between any two objects in the different clusters (i.e., by the "furthest neighbors").
- **Group average linkage:** In this method, the distance between two clusters is calculated as the average distance between all pairs of objects in the two different clusters.

# Minimal Spanning Tree – Single Linkage

- Build MST (Minimum Spanning Tree)
  - Start with a tree that consists of any point
  - In successive steps, look for the closest pair of points (p, q) such that one point (p) is in the current tree but the other (q) is not
  - Add q to the tree and put an edge between p and q

# MST: Divisive Hierarchical Clustering

- Use MST for constructing hierarchy of clusters

**Algorithm 7.5** MST Divisive Hierarchical Clustering Algorithm

1: Compute a minimum spanning tree for the proximity graph.
2: **repeat**
3:   Create a new cluster by breaking the link corresponding to the largest distance (smallest similarity).
4: **until** Only singleton clusters remain

# Summary of hierarchal clustering

- No need to specify the number of clusters in advance.
- Hierarchal nature maps nicely onto human intuition for some domains
- They do not scale well: time complexity of at least $O(n^2)$, where $n$ is the number of total objects.

# Partitional clustering

- Nonhierarchical, each instance is placed in exactly one of K nonoverlapping clusters.

- Since only one set of clusters is output, the user normally has to input the desired number of clusters K.

# *k-means*

1. Decide on a value for *k*.

2. Initialize the *k* cluster centers (randomly, if necessary).

3. Decide the class memberships of the *N* objects by assigning them to the nearest cluster center.

4. Re-estimate the *k* cluster centers, by assuming the memberships found above are correct.

5. If none of the *N* objects changed membership in the last iteration, exit. Otherwise goto 3.

# K-means clustering: step 1



Algorithm: k-means, Distance Metric: Euclidean Distance

# K-means clustering: step 2



Algorithm: k-means, Distance Metric: Euclidean Distance

# K-means clustering: step 3



Algorithm: k-means, Distance Metric: Euclidean Distance

# K-means clustering: step 4

Algorithm: k-means, Distance Metric: Euclidean Distance

# K-means clustering: step 5

Algorithm: k-means, Distance Metric: Euclidean Distance

# Evaluation of *K-means*

- Strength
  - *Relatively efficient*: $O(tkn)$, where $n$ is # objects, $k$ is # clusters, and $t$ is # iterations. Normally, $k$, $t$ << $n$.
  - Often terminates at a *local optimum*. The *global optimum* may be found using techniques such as: *deterministic annealing* and *genetic algorithms*
- Weakness
  - Applicable only when *mean* is defined, then what about categorical data?
  - Need to specify $k$, the *number* of clusters, in advance
  - Unable to handle noisy data and *outliers*
  - Not suitable for clusters with *non-convex shapes*

IIT KHARAGPUR

NPTEL ONLINE
CERTIFICATION COURSES

# DBSCAN

- DBSCAN is a density-based algorithm.
  - Density = number of points within a specified radius (Eps)

  - A point is a core point if it has more than a specified number of points (MinPts) within Eps
    - These are points that are at the interior of a cluster

  - A border point has fewer than MinPts within Eps, but is in the neighborhood of a core point

  - A noise point is any point that is not a core point or a border point.

# DBSCAN: Core, Border, and Noise Points

# DBSCAN Algorithm

- Eliminate noise points

- Perform clustering on the remaining points

$current\_cluster\_label \leftarrow 1$

**for** all core points **do**

    **if** the core point has no cluster label **then**

        $current\_cluster\_label \leftarrow current\_cluster\_label + 1$

        Label the current core point with cluster label $current\_cluster\_label$

    **end if**

    **for** all points in the $Eps$-neighborhood, except $i^{th}$ the point itself **do**

        **if** the point does not have a cluster label **then**

            Label the point with cluster label $current\_cluster\_label$

        **end if**

    **end for**

**end for**

# DBSCAN: Core, Border and Noise Points



Original Points

Point types: core, border and noise

Eps = 10, MinPts = 4

# When DBSCAN Works Well



Original Points

Clusters

- Resistant to Noise

- Can handle clusters of different shapes and sizes

# When DBSCAN Does NOT Work Well



Original Points



(MinPts=4, Eps=9.75).



(MinPts=4, Eps=9.92)

- Varying densities

- High-dimensional data

# DBSCAN: Determining EPS and MinPts

- Idea is that for points in a cluster, their k$^{th}$ nearest neighbors are at roughly the same distance

- Noise points have the k$^{th}$ nearest neighbor at farther distance

- So, plot sorted distance of every point to its k$^{th}$ nearest neighbor

# Summary of Clustering Algorithms

- K-Means – fast, works only for data where mean can be defined, generates spherical clusters, robust to noise

- Single linkage – produces non-convex clusters, slow for large data sets, sensitive to noise

- Complete linkage – produces non-convex clusters, very sensitive to noise, very slow for large data sets

- DBSCAN – produces arbitrary shaped clusters – works only for low dimensional data

# Cluster Validity

- For supervised classification we have a variety of measures to evaluate how good our model is
  - Accuracy, precision, recall

- For cluster analysis, the analogous question is how to evaluate the "goodness" of the resulting clusters?

- But "clusters are in the eye of the beholder"!

- Then why do we want to evaluate them?
  - To avoid finding patterns in noise
  - To compare clustering algorithms
  - To compare two sets of clusters
  - To compare two clusters

# Different Aspects of Cluster Validation

1.  Determining the clustering tendency of a set of data, i.e., distinguishing whether non-random structure actually exists in the data.

2.  Comparing the results of a cluster analysis to externally known results, e.g., to externally given class labels.

3.  Evaluating how well the results of a cluster analysis fit the data *without* reference to external information.

    - Use only the data

4.  Comparing the results of two different sets of cluster analyses to determine which is better.

5.  Determining the 'correct' number of clusters.

# Measures of Cluster Validity

- Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following three types.

    - **External Index:** Used to measure the extent to which cluster labels match externally supplied class labels.
        - Entropy

    - **Internal Index:** Used to measure the goodness of a clustering structure *without* respect to external information.
        - Sum of Squared Error (SSE)

    - **Relative Index:** Used to compare two different clusterings or clusters.
        - Often an external or internal index is used for this function, e.g., SSE or entropy

# Scatter Coefficient

- Cluster evaluation index

- Ratio of average intra-cluster distances to intra-cluster distances (Sum Squared Error)

# Internal Measures: Cohesion and Separation

- **Cluster Cohesion**: Measures how closely related are objects in a cluster
  - Example: SSE
- **Cluster Separation**: Measure how distinct or well-separated a cluster is from other clusters
- Example: Squared Error
  - Cohesion is measured by the within cluster sum of squares (SSE)
  $$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$
  - Separation is measured by the between cluster sum of squares
  $$BSS = \sum_i |C_i|(m - m_i)^2$$
    - Where $|C_i|$ is the size of cluster i

- Example: SSE
  - BSS + WSS = constant



K=1 cluster:

$$WSS = (1-3)^2 + (2-3)^2 + (4-3)^2 + (5-3)^2 = 10$$

$$BSS = 4 \times (3-3)^2 = 0$$

$$Total = 10 + 0 = 10$$

K=2 clusters:

$$WSS = (1-1.5)^2 + (2-1.5)^2 + (4-4.5)^2 + (5-4.5)^2 = 1$$

$$BSS = 2 \times (3-1.5)^2 + 2 \times (4.5-3)^2 = 9$$

$$Total = 1 + 9 = 10$$

# Internal Measures: Cohesion and Separation

- A proximity graph based approach can also be used for cohesion and separation.
  - Cluster cohesion is the sum of the weight of all links within a cluster.
  - Cluster separation is the sum of the weights between nodes in the cluster and nodes outside the cluster.

cohesion                                    separation

# Internal Measures: Silhouette Coefficient

- Silhouette Coefficient combine ideas of both cohesion and separation, but for individual points, as well as clusters and clusterings
- For an individual point, *i*
  - Calculate **a** = average distance of *i* to the points in its cluster
  - Calculate **b** = min (average distance of *i* to points in another cluster)
  - The silhouette coefficient for a point is then given by
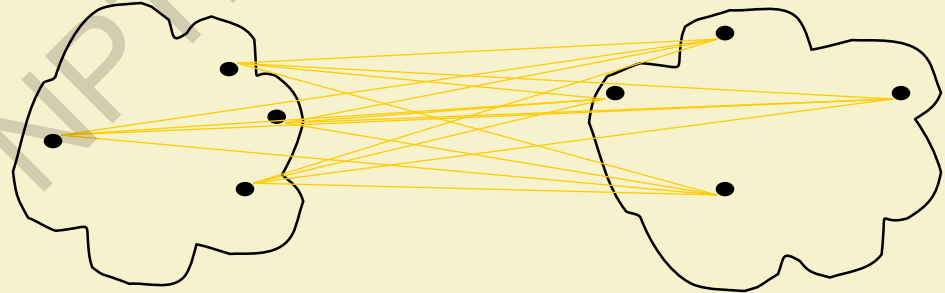
    s = 1 – a/b   if a < b,   (or s = b/a - 1    if a $\geq$ b, not the usual case)

  - Typically between 0 and 1.
  - The closer to 1 the better.

- Can calculate the Average Silhouette width for a cluster or a clustering

# External Measures of Cluster Validity: Entropy and Purity

**Table 5.9.** K-means Clustering Results for LA Document Data Set

| Cluster | Entertainment | Financial | Foreign | Metro | National | Sports | Entropy | Purity |
|---------|---------------|-----------|---------|-------|----------|--------|---------|--------|
| 1 | 3 | 5 | 40 | 506 | 96 | 27 | 1.2270 | 0.7474 |
| 2 | 4 | 7 | 280 | 29 | 39 | 2 | 1.1472 | 0.7756 |
| 3 | 1 | 1 | 1 | 7 | 4 | 671 | 0.1813 | 0.9796 |
| 4 | 10 | 162 | 3 | 119 | 73 | 2 | 1.7487 | 0.4390 |
| 5 | 331 | 22 | 5 | 70 | 13 | 23 | 1.3976 | 0.7134 |
| 6 | 5 | 358 | 12 | 212 | 48 | 13 | 1.5523 | 0.5525 |
| Total | 354 | 555 | 341 | 943 | 273 | 738 | 1.1450 | 0.7203 |

**entropy** For each cluster, the class distribution of the data is calculated first, i.e., for cluster $j$ we compute $p_{ij}$, the 'probability' that a member of cluster $j$ belongs to class $i$ as follows: $p_{ij} = m_{ij}/m_j$, where $m_j$ is the number of values in cluster $j$ and $m_{ij}$ is the number of values of class $i$ in cluster $j$. Then using this class distribution, the entropy of each cluster $j$ is calculated using the standard formula $e_j = \sum_{i=1}^{L} p_{ij} \log_2 p_{ij}$, where the $L$ is the number of classes. The total entropy for a set of clusters is calculated as the sum of the entropies of each cluster weighted by the size of each cluster, i.e., $e = \sum_{i=1}^{K} \frac{m_i}{m} e_j$, where $m_j$ is the size of cluster $j$, $K$ is the number of clusters, and $m$ is the total number of data points.

**purity** Using the terminology derived for entropy, the purity of cluster $j$, is given by $purity_j = \max p_{ij}$ and the overall purity of a clustering by $purity = \sum_{i=1}^{K} \frac{m_i}{m} purity_j$.

# Outliers Detection

- Important in many applications like anomaly detection

- Outliers are points not belonging to any cluster

- Many outlier detection algorithms available

# End of Clustering