



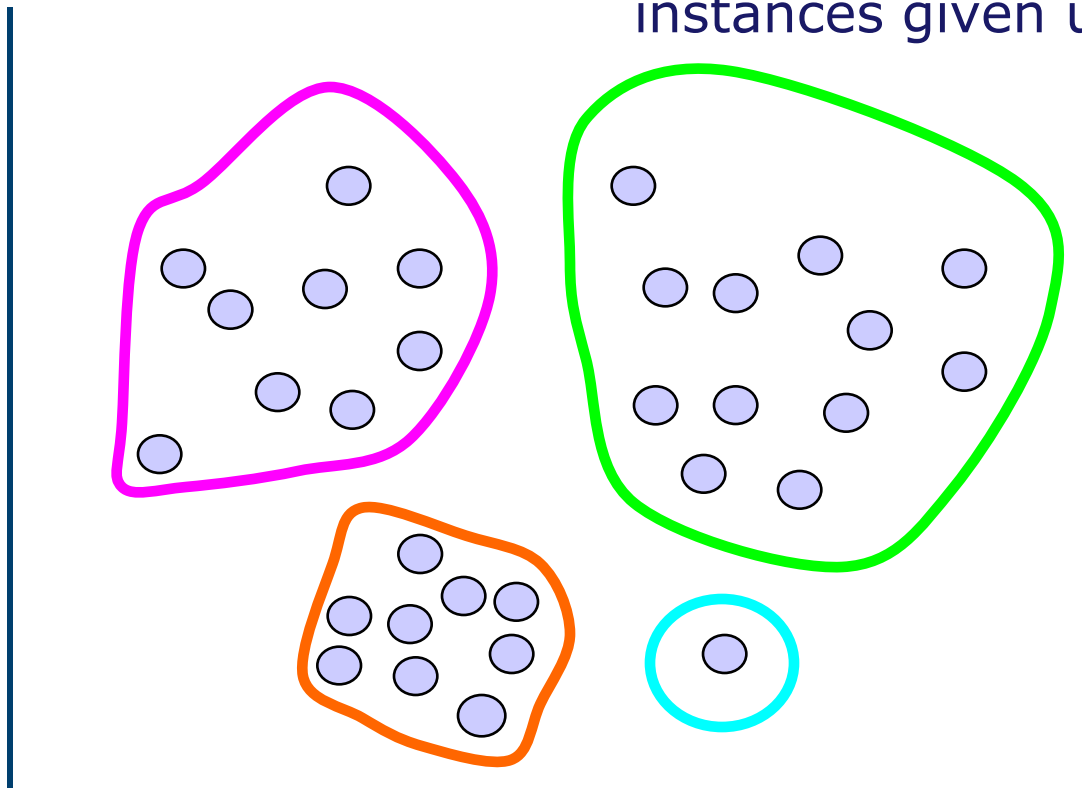
Clustering: Introduction

Data Mining and Text Mining (UIC 583 @ Politecnico di Milano)

- ❑ What is cluster analysis?
- ❑ Why clustering?
- ❑ What is good clustering?
- ❑ How to manage data types?
- ❑ What are the major clustering approaches?

- ❑ A **cluster** is a collection of data objects
 - ▶ Similar to one another within the same cluster
 - ▶ Dissimilar to the objects in other clusters
- ❑ Cluster analysis
 - ▶ Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters
- ❑ Unsupervised learning: no predefined classes
- ❑ Typical applications
 - ▶ As a stand-alone tool to get insight into data distribution
 - ▶ As a preprocessing step for other algorithms

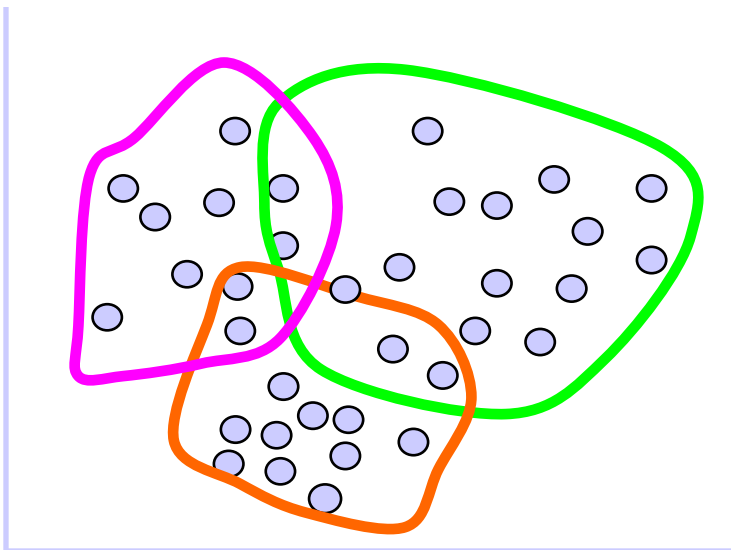
Clustering = Unsupervised learning:
Finds “natural” grouping of
instances given un-labeled data



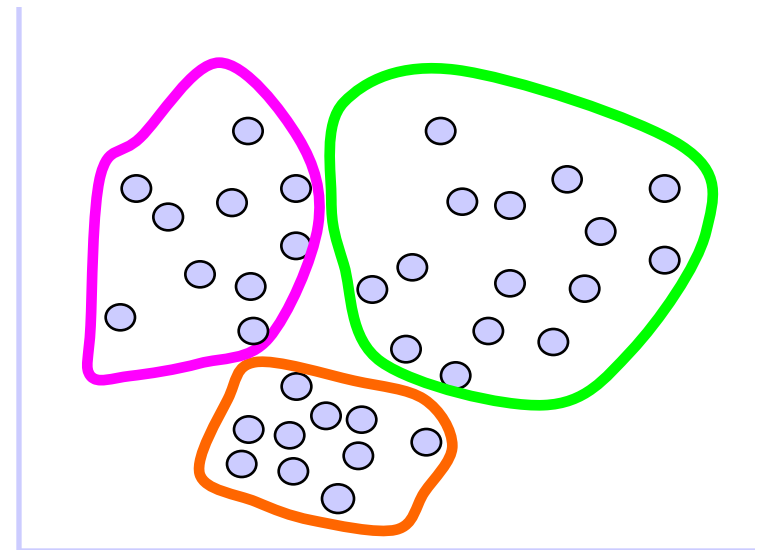
- ❑ Many different method and algorithms:
 - ▶ For numeric and/or symbolic data
 - ▶ Deterministic vs. probabilistic
 - ▶ Exclusive vs. overlapping
 - ▶ Hierarchical vs. flat
 - ▶ Top-down vs. bottom-up

Clusters: exclusive vs. overlapping

Overlapping



Non-overlapping



- ❑ Manual inspection
- ❑ Benchmarking on existing labels
- ❑ Cluster quality measures
 - ▶ distance measures
 - ▶ high similarity within a cluster,
low across clusters

- ❑ Pattern Recognition
- ❑ Spatial Data Analysis
 - ▶ Create thematic maps in GIS by clustering feature spaces
 - ▶ Detect spatial clusters or for other spatial mining tasks
- ❑ Image Processing
- ❑ Economic Science (especially market research)
- ❑ WWW
 - ▶ Document classification
 - ▶ Cluster Weblog data to discover groups of similar access patterns

- ❑ **Marketing**: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
- ❑ **Land use**: Identification of areas of similar land use in an earth observation database
- ❑ **Insurance**: Identifying groups of motor insurance policy holders with a high average claim cost
- ❑ **City-planning**: Identifying groups of houses according to their house type, value, and geographical location
- ❑ **Earth-quake studies**: Observed earth quake epicenters should be clustered along continent faults

- ❑ A good clustering consists of high quality clusters with
 - ▶ high intra-class similarity
 - ▶ low inter-class similarity
- ❑ The quality of a clustering result depends on both the similarity measure used by the method and its implementation
- ❑ The quality of a clustering method is also measured by its ability to discover some or all of the hidden patterns

- ❑ **Dissimilarity/Similarity metric**: Similarity is expressed in terms of a distance function, typically metric: $d(i, j)$
- ❑ There is a separate “quality” function that measures the “goodness” of a cluster.
- ❑ The definitions of **distance** functions are usually very different for interval-scaled, Boolean, categorical, ordinal ratio, and vector variables.
- ❑ Weights should be associated with different variables based on applications and data semantics.
- ❑ It is hard to define “similar enough” or “good enough”
 - ▶ the answer is typically highly subjective.

- ☐ Scalability
- ☐ Ability to deal with different types of attributes
- ☐ Ability to handle dynamic data
- ☐ Discovery of clusters with arbitrary shape
- ☐ Minimal requirements for domain knowledge to determine input parameters
- ☐ Able to deal with noise and outliers
- ☐ Insensitive to order of input records
- ☐ High dimensionality
- ☐ Incorporation of user-specified constraints
- ☐ Interpretability and usability

□ Data

| Outlook | Temperature | Humidity | Windy | Play |
|----------|-------------|----------|-------|------|
| Sunny | Hot | High | False | No |
| Sunny | Hot | High | True | No |
| Overcast | Hot | High | False | Yes |
| Rainy | Mild | Normal | False | Yes |
| ... | ... | ... | ... | ... |

□ Dissimilarity matrix

$$\begin{bmatrix}
 0 & & & & \\
 d(2,1) & 0 & & & \\
 d(3,1) & d(3,2) & 0 & & \\
 \vdots & \vdots & \vdots & & \\
 d(n,1) & d(n,2) & \dots & \dots & 0
 \end{bmatrix}$$

□ Data matrix

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}$$

□ Dissimilarity matrix

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

- ☐ Interval-scaled variables
- ☐ Binary variables
- ☐ Nominal, ordinal, and ratio variables
- ☐ Variables of mixed types

□ Standardize data

- ▶ Calculate the mean absolute deviation,

$$s_f = \frac{1}{n} (|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

- ▶ where $m_f = \frac{1}{n} (x_{1f} + x_{2f} + \dots + x_{nf})$.

- ▶ Calculate the standardized measurement (z-score)

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

- Using mean absolute deviation is more robust than using standard deviation

- ❑ Distances are normally used to measure the similarity or dissimilarity between two data objects
- ❑ Some popular ones include: Minkowski distance:

$$d(i, j) = \sqrt[q]{(|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{ip} - x_{jp}|^q)}$$

- ❑ where $i = (x_{i1}, x_{i2}, \dots, x_{ip})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jp})$ are two p -dimensional data objects, and q is a positive integer
- ❑ If $q = 1$, d is Manhattan distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

- If $q = 2$, d is Euclidean distance:

$$d(i,j) = \sqrt{(|x_{i_1} - x_{j_1}|^2 + |x_{i_2} - x_{j_2}|^2 + \dots + |x_{i_p} - x_{j_p}|^2)}$$

► Properties

- $d(i,j) \geq 0$
- $d(i,i) = 0$
- $d(i,j) = d(j,i)$
- $d(i,j) \leq d(i,k) + d(k,j)$

- Also, one can use weighted distance, parametric Pearson product moment correlation, or other dissimilarity measures

- Contingency table for binary data:

| | 1 | 0 | sum |
|-----|-----|-----|-----|
| 1 | a | b | a+b |
| 0 | c | d | c+d |
| sum | a+c | b+d | p |

- Distance measure for symmetric binary variables:
- Distance measure for asymmetric binary variables:
- Jaccard coefficient (similarity measure for asymmetric binary variables):

$$d(i, j) = \frac{b+c}{a+b+c+d}$$

$$d(i, j) = \frac{b+c}{a+b+c}$$

$$sim_{Jaccard}(i, j) = \frac{a}{a+b+c}$$

Example of Dissimilarity in Binary Variables

20

| Name | Gender | Fever | Cough | Test-1 | Test-2 | Test-3 | Test-4 |
|------|--------|-------|-------|--------|--------|--------|--------|
| Jack | M | Y | N | P | N | N | N |
| Mary | F | Y | N | P | N | P | N |
| Jim | M | Y | P | N | N | N | N |

- ❑ Gender is a symmetric attribute
- ❑ Remaining attributes are asymmetric binary
- ❑ Let the values Y and P be set to 1, and the value N be set to 0

$$d(jack, mary) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d(jack, jim) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d(jim, mary) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$

- ❑ A generalization of the binary variable in that it can take more than 2 states, e.g., red, yellow, blue, green
- ❑ Method 1: simple matching
 - ▶ m : # of matches, p : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

- ❑ Method 2: use a large number of binary variables
 - ▶ creating a new binary variable for each of the M nominal states

- ❑ An ordinal variable can be discrete or continuous
- ❑ Order is important, e.g., rank
- ❑ Can be treated like interval-scaled
 - ▶ replace x_{if} by their rank

$$r_{if} \in \{1, \dots, M_f\}$$

- ▶ map the range of each variable onto $[0, 1]$ by replacing i -th object in the f -th variable by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

- ▶ compute the dissimilarity using methods for interval-scaled variables

- ❑ Ratio-scaled variable: a positive measurement on a nonlinear scale, approximately at exponential scale, such as Ae^{Bt} or Ae^{-Bt}
- ❑ Methods:
 - ▶ treat them like interval-scaled variables—not a good choice! (why?—the scale can be distorted)
 - ▶ apply logarithmic transformation

$$y_{if} = \log(x_{if})$$

- ▶ treat them as continuous ordinal data treat their rank as interval-scaled

- ❑ A database may contain all the six types of variables
 - ▶ symmetric binary, asymmetric binary, nominal, ordinal, interval and ratio
- ❑ One may use a weighted formula to combine their effects

$$d(i, j) = \frac{\sum_{f=1}^p \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p \delta_{ij}^{(f)}}$$

- ▶ f is binary or nominal:
 $d_{ij}(f) = 0$ if $x_{if} = x_{jf}$, or $d_{ij}(f) = 1$ otherwise
- ▶ f is interval-based: use the normalized distance
- ▶ f is ordinal or ratio-scaled
 - compute ranks r_{if} and
 - and treat z_{if} as interval-scaled $z_{if} = \frac{r_{if} - 1}{M_f - 1}$

- ❑ Vector objects: keywords in documents, gene features in micro-arrays, etc.
- ❑ Broad applications: information retrieval, biologic taxonomy, etc.
- ❑ Cosine measure

$$s(\vec{X}, \vec{Y}) = \frac{\vec{X}^t \cdot \vec{Y}}{|\vec{X}| |\vec{Y}|},$$

- ❑ A variant: Tanimoto coefficient

\vec{X}^t is a transposition of vector \vec{X} , $|\vec{X}|$ is the Euclidean normal of vector \vec{X} ,

$$s(\vec{X}, \vec{Y}) = \frac{\vec{X}^t \cdot \vec{Y}}{\vec{X}^t \cdot \vec{X} + \vec{Y}^t \cdot \vec{Y} - \vec{X}^t \cdot \vec{Y}},$$

□ Partitioning approach:

- ▶ Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
- ▶ Typical methods: k-means, k-medoids, CLARANS

□ Hierarchical approach:

- ▶ Create a hierarchical decomposition of the set of data (or objects) using some criterion
- ▶ Typical methods: Diana, Agnes, BIRCH, ROCK, CAMELEON

□ Density-based approach:

- ▶ Based on connectivity and density functions
- ▶ Typical methods: DBSACN, OPTICS, DenClue

- ❑ Grid-based approach
 - ▶ based on a multiple-level granularity structure
 - ▶ Typical methods: STING, WaveCluster, CLIQUE
- ❑ Model-based
 - ▶ A model is hypothesized for each of the clusters and tries to find the best fit of that model to each other
 - ▶ Typical methods: EM, SOM, COBWEB
- ❑ Frequent pattern-based
 - ▶ Based on the analysis of frequent patterns
 - ▶ Typical methods: pCluster
- ❑ User-guided or constraint-based
 - ▶ Clustering by considering user-specified or application-specific constraints
 - ▶ Typical methods: COD (obstacles), constrained clustering

- ❑ Clusters are collection of data objects
- ❑ Objects should be
 - ▶ Similar to one another within the same cluster
 - ▶ Dissimilar to the objects in other clusters
- ❑ Cluster analysis searches for similarities between data according to the characteristics found in the data and groups similar data objects into clusters
- ❑ Similarity is defined in terms of distance between objects and between clusters
- ❑ Several approaches: partition-based, hierarchical, density-based, model-based, etc.