

Clustering

What is cluster analysis?

- What is a cluster?
 - A cluster is a collection of data objects which are
 - Similar (or related) to one another within the same group (i.e., cluster)
 - Dissimilar (or unrelated) to the objects in other groups (i.e., clusters)
- Cluster analysis (or *clustering*, *data segmentation*, ...)
 - Given a set of data points, partition them into a set of groups (i.e., clusters) which are as similar as possible
- Cluster analysis is **unsupervised learning** (i.e., no predefined classes)
 - This contrasts with *classification* (i.e., supervised learning)



What is cluster analysis?

- Typical ways to use/apply cluster analysis
 - As a stand-alone tool to get insight into data distribution, or
 - As a preprocessing (or intermediate) step for other algorithms

What Is Good Clustering?

- A good clustering method will produce high quality clusters which should have
 - High intra-class similarity: *Cohesive* within clusters
 - Low inter-class similarity: *Distinctive* between clusters
- Quality function
 - There is usually a separate “quality” function that measures the “goodness” of a cluster
 - It is hard to define “similar enough” or “good enough”
 - The answer is typically highly subjective
- There exist many similarity measures and/or functions for different applications
- Similarity measure is critical for cluster analysis



Cluster Analysis: Applications

- A key intermediate step for other data mining tasks
 - Generating a compact summary of data for classification, pattern discovery, hypothesis generation and testing, etc.
 - Outlier detection: Outliers—those “far away” from any cluster
- Data summarization, compression, and reduction
 - Ex. Image processing: Vector quantization
- Collaborative filtering, recommendation systems, or customer segmentation
 - Find like-minded users or similar products



Cluster Analysis: Applications

- *Dynamic trend detection*
 - *Clustering stream data and detecting trends and patterns*
- *Multimedia data analysis, biological data analysis and social network analysis*
 - *Ex. Clustering images or video/audio clips, gene/protein sequences, etc.*



Considerations for Cluster Analysis

- *Partitioning criteria (Single level vs. hierarchical partitioning)*
 - *Single level: All clusters are conceptually at the same level*
 - *Eg: partitioning customers into groups so that each group has its manager.*
 - *Hierarchical level: Clusters at different semantic levels.*
 - *Eg: general topics: “sports”, “politics” and subtopics in text mining.*
- *Separation of clusters*
 - *Exclusive (e.g., one customer belongs to only one region) vs. non-exclusive (e.g., one document may belong to more than one class)*



Considerations for Cluster Analysis

- *Similarity measure*
 - Distance-based (e.g., Euclidean, road network, vector) vs. similarity measure defined as connectivity-based (e.g., density or contiguity)
- *Clustering space*
 - Full space (often when low dimensional) vs. subspaces (often in high-dimensional clustering due to presence of irrelevant attributes)

Requirements and Challenges

- **Quality**
 - Ability to deal with different types of attributes: Numerical, categorical, text, multimedia, networks, and mixture of multiple types
 - Discovery of clusters with arbitrary shape
 - Ability to deal with noisy data
- **Scalability**
 - Clustering all the data instead of only samples
 - High dimensionality
 - Incremental or stream clustering and insensitivity to input order



Requirements and Challenges

- *Constraint-based clustering*
 - User-given preferences or constraints; domain knowledge; user queries
- *Interpretability and usability*
 - Clustering results should be interpretable, comprehensible and usable
 - Can able to tie with specific semantic interpretations and applications
- *Discovery of clusters with arbitrary shape*
 - Distance based clustering algorithm produces spherical clusters with similar size and density
 - Important to develop clusters of any shape.



Requirements and Challenges

- *Ability to deal with noisy data*
 - *Most data contains outliers, missing , unknown or erroneous data*
 - *Clustering algorithms are sensitive produce poor quality clusters*
 - *Need methods that are robust to noise*
- *Incremental clustering and insensitivity to input order:*
 - *Incremental updates requires recomputing from scratch and return different clusters depending of the order of the data given.*
- *High dimensionality: Need to handle high dimension data*



Type of data in clustering analysis

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

Interval-valued variables

- 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

Similarity and Dissimilarity Between Objects

- 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

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

Similarity and Dissimilarity Between Objects (Cont.)

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

$$d(i, j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

– 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

Binary Variables

- A contingency table for binary data

	1	0	sum
1	a	b	$a+b$
0	c	d	$c+d$
sum	$a+c$	$b+d$	p

- Simple matching coefficient (invariant, if the binary variable is symmetric): $d(i, j) = \frac{b + c}{a + b + c + d}$
- Jaccard coefficient (noninvariant if the binary variable is asymmetric): $d(i, j) = \frac{b + c}{a + b + c}$

Dissimilarity between Binary Variables

- Example

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
- the remaining attributes are asymmetric binary
- let the values Y and P be set to 1, and the value N be set to 0

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

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

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

Nominal Variables

- 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

Ordinal Variables

- 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 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

Variables of Mixed Types

- 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 $d_{ij}^{(f)} = 1$ o.w.

- 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

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