



# Cluster Analysis

——Types of Data in Cluster Analysis——

徐华

清华大学 计算机系 智能技术与系统国家重点实验室

xuhua@tsinghua.edu.cn

1

## Cluster Analysis



- ◉ What is Cluster Analysis?
- ◉ **Types of Data in Cluster Analysis**
- ◉ A Categorization of Major Clustering Methods
- ◉ Partitioning Methods
- ◉ Hierarchical Methods
- ◉ Density-Based Methods
- ◉ Grid-Based Methods
- ◉ Model-Based Clustering Methods
- ◉ Outlier Analysis
- 2 ◉ **Summary**



## Data Structures



- ◉ **Data matrix**
  - ◆ (two modes)

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

- ◉ **Dissimilarity matrix**
  - ◆ (one mode)

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

3



## Measure the Quality of Clustering



- ◉ **Dissimilarity/Similarity metric:** Similarity is expressed in terms of a distance function, which is 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 and ratio 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.

4



## Type of data in clustering analysis



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

5



## Interval-valued variables



- ◉ Standardize data
  - ◆ Calculate the mean absolute deviation(偏差):

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

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

6



## 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

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

7



## 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 (皮尔逊积差相关系数), or other dissimilarity measures

8



## Binary Variables



- A contingency table for binary data

		Object $j$		
		1	0	sum
Object $i$	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):

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

9



## 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 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

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

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

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

10



## 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
  - ◆  $m$ : the number of matches,  $p$ : the total number 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

11



## 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

12



## 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 and then treat their rank as interval-scaled

13



## 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:
 
$$d_{ij}^{(f)} = 0 \text{ if } x_{if} = x_{jf}, \text{ or } d_{ij}^{(f)} = 1$$
- ◆  $f$  is interval-based: use the normalized distance
- ◆  $f$  is ordinal or ratio-scaled
  - compute ranks  $r_{if}$  and
  - treat  $z_{if}$  as interval-scaled

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

14





**Thanks !**

