



Cluster Analysis

——What is 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



What is Cluster Analysis?



- ◉ **Cluster: a collection of data objects**
 - ◆ Similar to one another within the same cluster
 - ◆ Dissimilar to the objects in other clusters
- ◉ **Cluster analysis**
 - ◆ Grouping a set of data objects into clusters
- ◉ **Clustering is *unsupervised classification*: no predefined classes**
- ◉ **Typical applications**
 - ◆ As a *stand-alone tool* to get insight into data distribution
 - ◆ As a *preprocessing step* for other algorithms

3



General Applications of Clustering



- ◉ **Pattern Recognition**
- ◉ **Spatial Data Analysis**
 - ◆ create thematic maps in GIS by clustering feature spaces
 - ◆ detect spatial clusters and explain them in spatial data mining
- ◉ **Image Processing**
- ◉ **Economic Science (especially market research)**
- ◉ **WWW**
 - ◆ Document classification
 - ◆ Cluster Weblog data to discover groups of similar access patterns

4



Examples of Clustering Applications



- ◉ **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 (震中) should be clustered along continent faults

5



What Is Good Clustering?



- ◉ A **good clustering** method will produce 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.

6



Requirements of Clustering in Data Mining



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

7



Thanks !

8

