



Cluster Analysis

——Outlier——

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 Outlier Discovery?



- ◉ **What are outliers?**
 - ◆ The set of objects are considerably dissimilar from the remainder of the data
- ◉ **Problem**
 - ◆ Find top n outlier points
- ◉ **Applications:**
 - ◆ Credit card fraud detection
 - ◆ Telecom fraud detection
 - ◆ Customer segmentation
 - ◆ Medical analysis

3



Outlier Discovery: Statistical Approaches



- ◉ **Assume a model underlying distribution that generates data set (e.g. normal distribution)**
 - ◆ Use discordancy tests depending on
 - data distribution
 - distribution parameter (e.g., mean, variance)
 - number of expected outliers
 - ◆ Drawbacks
 - most tests are for single attribute
 - In many cases, data distribution may not be known

4



Outlier Discovery: Distance-Based Approach



- ◉ Introduced to counter the main limitations imposed by statistical methods
 - ◆ We need multi-dimensional analysis without knowing data distribution.
- ◉ Distance-based outlier: A $DB(p, D)$ -outlier is an object O in a dataset T such that at least a fraction p of the objects in T lies at a distance greater than D from O
- ◉ Algorithms for mining distance-based outliers
 - ◆ Index-based algorithm
 - ◆ Nested-loop algorithm
 - ◆ Cell-based algorithm

5



Outlier Discovery: Deviation-Based Approach



- ◉ Identifies outliers by examining the main characteristics of objects in a group
- ◉ Objects that “deviate” from this description are considered outliers
- ◉ Sequential exception technique
 - ◆ simulates the way in which humans can distinguish unusual objects from among a series of supposedly like objects
- ◉ OLAP data cube technique
 - ◆ uses data cubes to identify regions of anomalies (异常) in large multidimensional data

6



Summary



- ◉ **Considerable progress has been made in scalable clustering methods**
 - ◆ **Partitioning: k-means, k-medoids, CLARANS**
 - ◆ **Hierarchical: BIRCH, CURE**
 - ◆ **Density-based: DBSCAN, CLIQUE, OPTICS**
 - ◆ **Grid-based: STING, WaveCluster**
 - ◆ **Model-based: Autoclass, Denclue, Cobweb**
- ◉ **Current clustering techniques do not address all the requirements adequately**
- ◉ **Constraint-based clustering analysis: Constraints exist in data space (bridges and highways) or in user queries**

7



Thanks!

8

