# Cluster Analysis
——Grid-Based Methods——

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**
- ◉ **Summary**

2

## Grid-Based Clustering Method

- ⊙ **Using multi-resolution grid data structure**
- ⊙ **Several interesting methods**
  - ◆ **STING** (a STatistical INformation Grid approach) by Wang, Yang and Muntz (1997)
  - ◆ **WaveCluster** by Sheikholeslami, Chatterjee, and Zhang (VLDB'98)
    - · **A multi-resolution clustering approach using wavelet method**
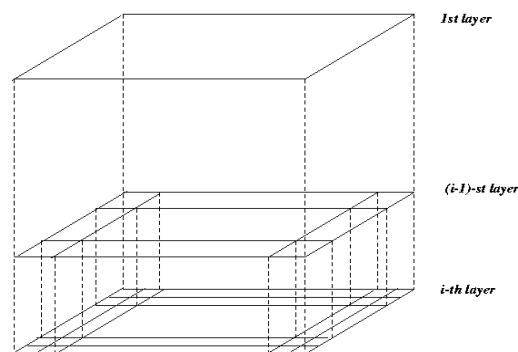  - ◆ **CLIQUE**: Agrawal, et al. (SIGMOD'98)

3

## STING: A Statistical Information Grid Approach

- ⊙ **Wang, Yang and Muntz (VLDB'97)**
- ⊙ **The spatial area is divided into rectangular cells**
- ⊙ **There are several levels of cells corresponding to different levels of resolution**



1st layer

(i-1)-st layer

i-th layer

4

## STING: A Statistical Information Grid Approach

- ◉ **Each cell at a high level is partitioned into a number of smaller cells in the next lower level**

- ◉ **Statistical info of each cell is calculated and stored beforehand and is used to answer queries**

- ◉ **Parameters of higher level cells can be easily calculated from parameters of lower level cell**
  - ◆ *count, mean, s, min, max*
  - ◆ type of distribution—normal, *uniform(均匀)*, etc.

- ◉ **Use a top-down approach to answer spatial data queries**

- ◉ **Start from a pre-selected layer—typically with a small number of cells**

- ◉ **For each cell in the current level compute the confidence interval**

5

---

## STING: A Statistical Information Grid Approach

- ◉ **Remove the irrelevant cells from further consideration**

- ◉ **When finish examining the current layer, proceed to the next lower level**

- ◉ **Repeat this process until the bottom layer is reached**

- ◉ **Advantages:**
  - ◆ **Query-independent, easy to parallelize, incremental update**
  - ◆ $O(K)$, **where** $K$ **is the number of grid cells at the lowest level**

- ◉ **Disadvantages:**
  - ◆ **All the cluster boundaries are either horizontal or vertical, and no diagonal（对角线的）boundary is detected**

6

# Thanks!

7