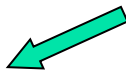# Unit IV : Clustering

Cluster Analysis – Partitioning Methods – Hierarchical Methods – Density Based Methods – Grid Based Methods – Outlier Analysis
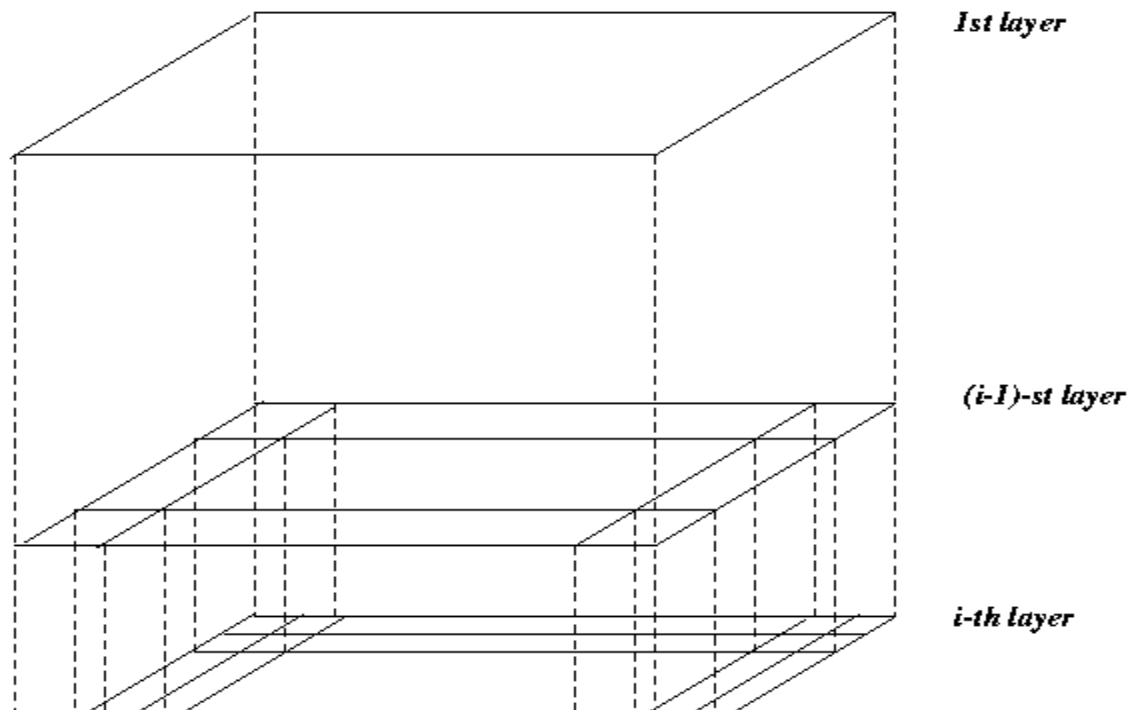
# Cluster Analysis

1.  What is Cluster Analysis?

2.  Types of Data in Cluster Analysis

3.  A Categorization of Major Clustering Methods

4.  Partitioning Methods

5.  Hierarchical Methods

6.  Density-Based Methods

7.  Grid-Based Methods

8.  Outlier Analysis

9.  Summary

# 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)

        - On high-dimensional data (thus put in the section of clustering high-dimensional data

# STING: A Statistical Information Grid Approach

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

# The STING Clustering Method

- 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*, 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
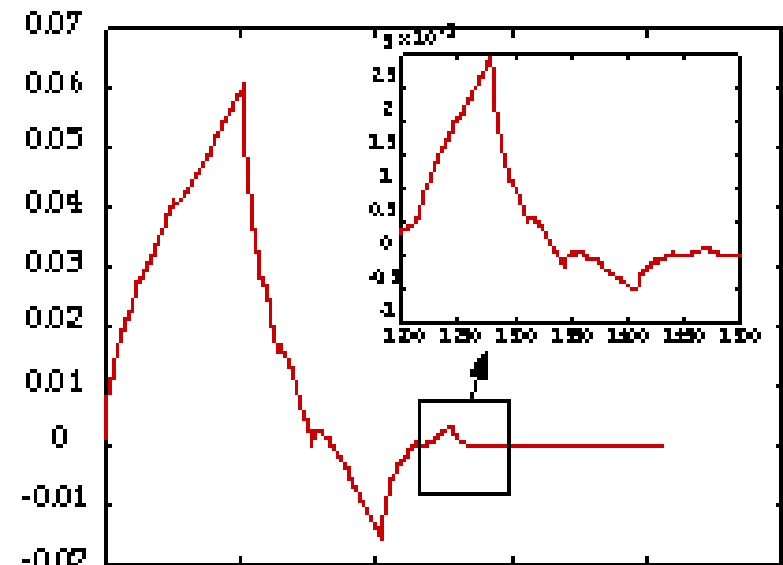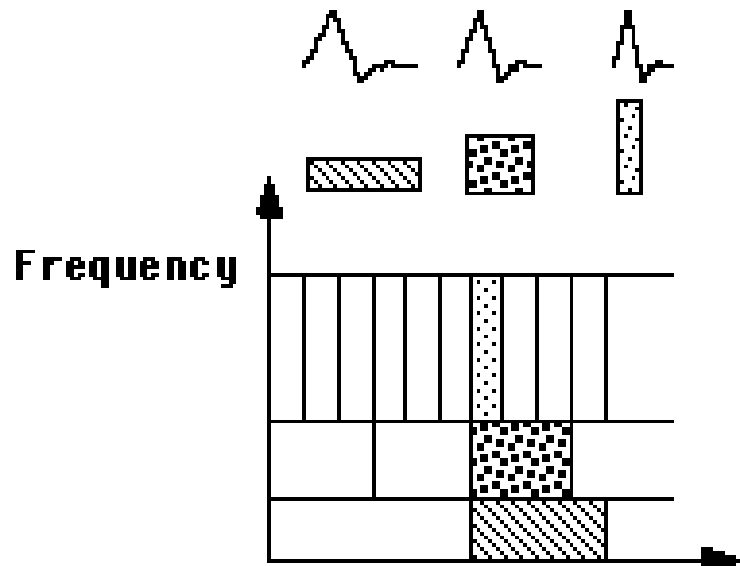
# Comments on STING

- 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 boundary is detected

# WaveCluster: Clustering by Wavelet Analysis (1998)

- Sheikholeslami, Chatterjee, and Zhang (VLDB'98)

- A multi-resolution clustering approach which applies wavelet transform to the feature space

- How to apply wavelet transform to find clusters

  - Summarizes the data by imposing a multidimensional grid structure onto data space

  - These multidimensional spatial data objects are represented in a n-dimensional feature space

  - Apply wavelet transform on feature space to find the dense regions in the feature space

  - Apply wavelet transform multiple times which result in clusters at different scales from fine to coarse

# Wavelet Transform

- Wavelet transform: A signal processing technique that decomposes a signal into different frequency sub-band (can be applied to n-dimensional signals)

- Data are transformed to preserve relative distance between objects at different levels of resolution

- Allows natural clusters to become more distinguishable

# The WaveCluster Algorithm

- Input parameters
  - # of grid cells for each dimension
  - the wavelet, and the # of applications of wavelet transform
- Why is wavelet transformation useful for clustering?
  - Use hat-shape filters to emphasize region where points cluster, but simultaneously suppress weaker information in their boundary
  - Effective removal of outliers, multi-resolution, cost effective
- Major features:
  - Complexity O(N)
  - Detect arbitrary shaped clusters at different scales
  - Not sensitive to noise, not sensitive to input order
  - Only applicable to low dimensional data
- Both grid-based and density-based

# Quantization & Transformation

- First, quantize data into m-D grid structure, then wavelet transform
    - a) scale 1: high resolution
    - b) scale 2: medium resolution
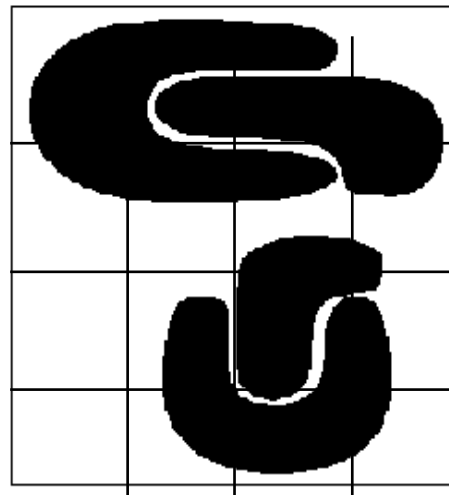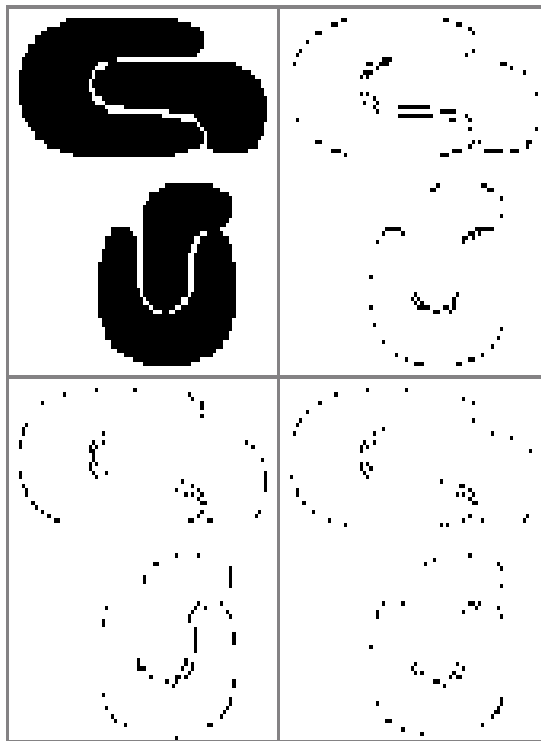    - c) scale 3: low resolution
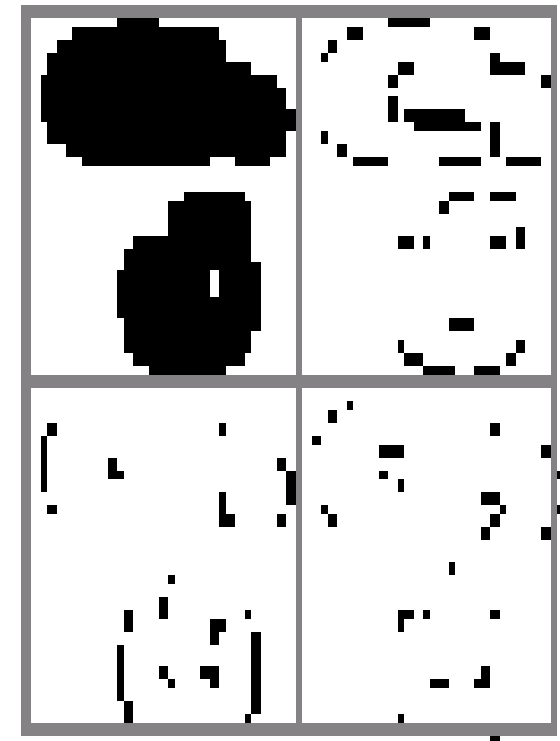
Figure 1: A sample 2-dimensional feature space.

# Quantization & Transformation



a)          b)          c)