

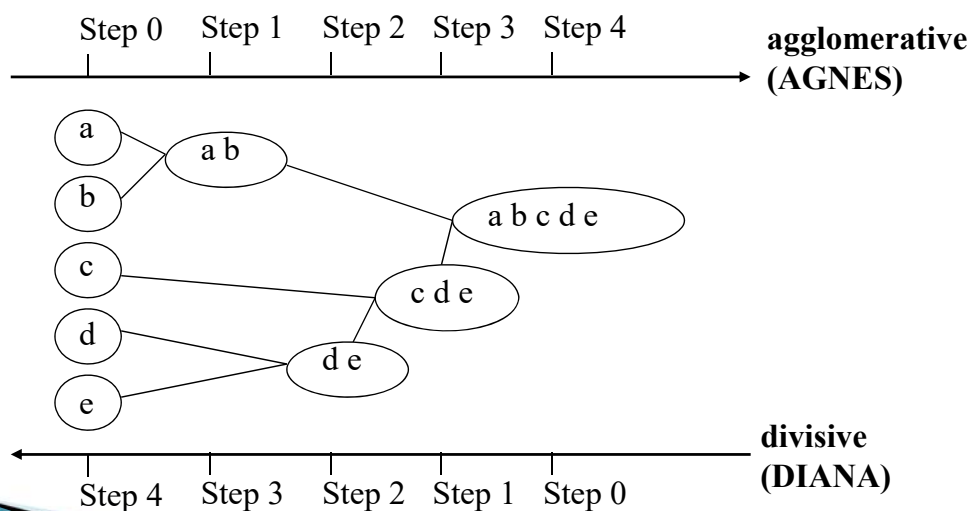
# Data Mining -- Clustering

Instructor: Jen-Wei Huang

Office: 92528 in the EE building  
jwhuang@mail.ncku

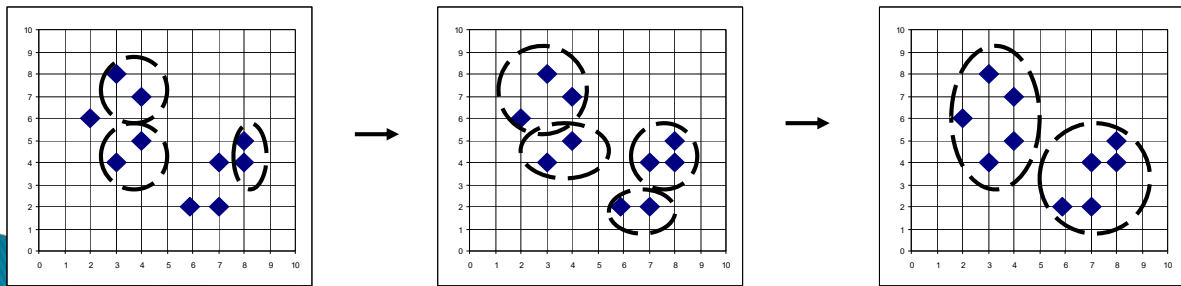
## Hierarchical Clustering

- ▶ Use distance matrix as clustering criteria. This method does not require the number of clusters  $k$  as an input, but needs a termination condition



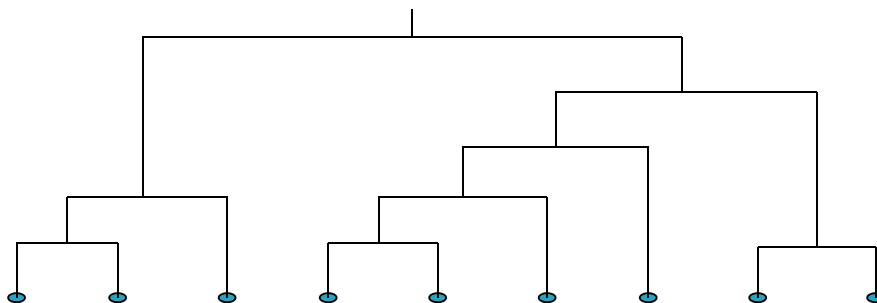
# AGNES (Agglomerative Nesting)

- ▶ Introduced in Kaufmann and Rousseeuw (1990)
- ▶ Implemented in statistical analysis packages, e.g., Splus
- ▶ Use the Single-Link method and the dissimilarity matrix.
- ▶ Merge nodes that have the least dissimilarity
- ▶ Go on in a non-descending fashion
- ▶ Eventually all nodes belong to the same cluster



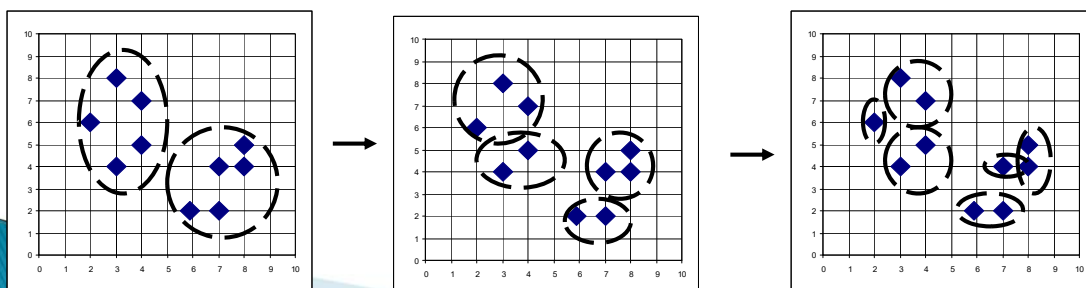
## Dendrogram

- ▶ Shows how the clusters are merged
- ▶ Decompose data objects into a several levels of nested partitioning (tree of clusters), called a dendrogram.
- ▶ A clustering of the data objects is obtained by cutting the dendrogram at the desired level, then each connected component forms a cluster.



# DIANA (Divisive Analysis)

- ▶ Introduced in Kaufmann and Rousseeuw (1990)
- ▶ Implemented in statistical analysis packages, e.g., Splus
- ▶ Inverse order of AGNES
- ▶ Eventually each node forms a cluster on its own



Data Mining & Social Network Analysis 2021/03/24

5

## Comments

- ▶ Major weakness of agglomerative clustering methods
  - do not scale well: time complexity of at least  $O(n^2)$ , where  $n$  is the number of total objects
  - can never undo what was done previously
- ▶ Integration of hierarchical with distance-based clustering
  - BIRCH (1996): uses CF-tree and incrementally adjusts the quality of sub-clusters
  - ROCK (1999): clustering categorical data by neighbor and link analysis
  - CHAMELEON (1999): hierarchical clustering using dynamic modeling

Data Mining & Social Network Analysis 2021/03/24

6

# BIRCH

- ▶ Birch: Balanced Iterative Reducing and Clustering using Hierarchies (Zhang, Ramakrishnan & Livny, SIGMOD'96)
- ▶ Incrementally construct a CF (Clustering Feature) tree, a hierarchical data structure for multiphase clustering
  - Phase 1: scan DB to build an initial in-memory CF tree (a multi-level compression of the data that tries to preserve the inherent clustering structure of the data)
  - Phase 2: use an arbitrary clustering algorithm to cluster the leaf nodes of the CF-tree

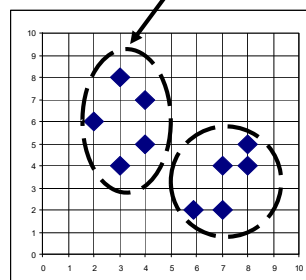
## Clustering Feature Vector

- ▶ Summary of the statistics for a given subcluster
- ▶ Registers crucial measurements for computing cluster and utilizes storage efficiently
- ▶ in BIRCH, **Clustering Feature:  $CF = (N, LS, SS)$**

$N$ : Number of data points

$$LS: \sum_{i=1}^N X_i$$

$$SS: \sum_{i=1}^N X_i^2$$



$$CF = (5, (16,30),(54,190))$$

(3,4)

(2,6)

(4,5)

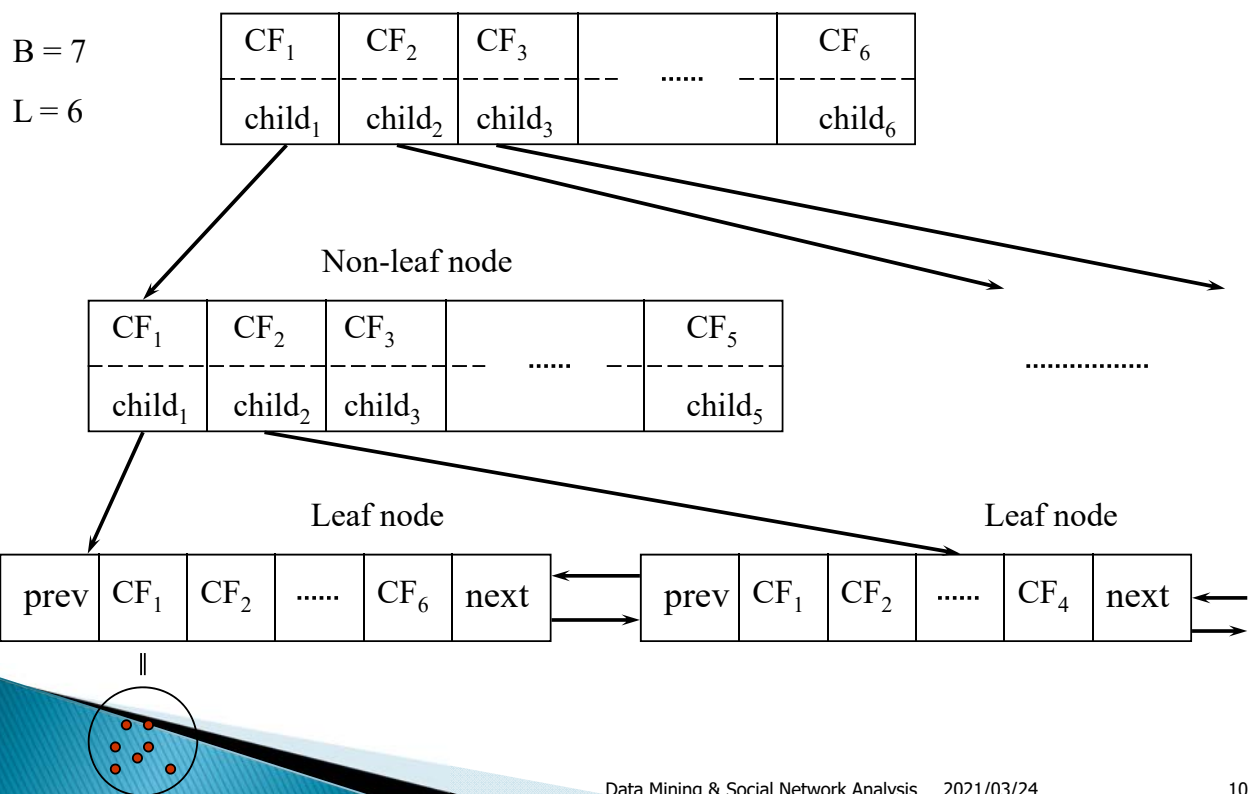
(4,7)

(3,8)

# CF-Tree

- ▶ A CF tree is a height-balanced tree that stores the clustering features for a hierarchical clustering
  - A nonleaf node in a tree has descendants or “children”
  - The nonleaf nodes store sums of the CFs of their children
- ▶ A CF tree has two parameters
  - Branching factor: specify the maximum number of children.
  - threshold: max diameter of sub-clusters stored at the leaf nodes

## CF-Tree



# Comments of BIRCH

- ▶ *Scales linearly*: finds a good clustering with a single scan and improves the quality with a few additional scans
- ▶ *Weakness*: handles only numeric data, and sensitive to the order of the data record.

# Clustering Categorical Data

- ▶ ROCK algorithm: RObust Clustering using linKs
  - S. Guha, R. Rastogi & K. Shim, ICDE'99
- ▶ Major ideas
  - Use links to measure similarity/proximity
  - Not distance-based
  - Computational complexity:
- ▶ Algorithm: sampling-based clustering
  - Draw random sample
  - Cluster with links
  - Label data in disk
- ▶ Experiments
  - Congressional voting, mushroom data

# Similarity Measure

- ▶ Traditional measures for categorical data may not work well, e.g., Jaccard coefficient
- ▶ Example: Two groups (clusters) of transactions
  - $C_1$ .  $\langle a, b, c, d, e \rangle$ :  $\{a, b, c\}, \{a, b, d\}, \{a, b, e\}, \{a, c, d\}, \{a, c, e\}, \{a, d, e\}, \{b, c, d\}, \{b, c, e\}, \{b, d, e\}, \{c, d, e\}$
  - $C_2$ .  $\langle a, b, f, g \rangle$ :  $\{a, b, f\}, \{a, b, g\}, \{a, f, g\}, \{b, f, g\}$
- ▶ Jaccard coefficient may lead to wrong clustering result
  - $C_1$ : 0.2 ( $\{a, b, c\}, \{b, d, e\}$ ) to 0.5 ( $\{a, b, c\}, \{a, b, d\}$ )
  - $C_1$  &  $C_2$ : could be as high as 0.5 ( $\{a, b, c\}, \{a, b, f\}$ )
- ▶ Jaccard coefficient-based similarity function:

$$Sim(T_1, T_2) = \frac{|T_1 \cap T_2|}{|T_1 \cup T_2|}$$

- Ex. Let  $T_1 = \{a, b, c\}$ ,  $T_2 = \{c, d, e\}$

$$Sim(T_1, T_2) = \frac{|\{c\}|}{|\{a, b, c, d, e\}|} = \frac{1}{5} = 0.2$$

# Link Measure

- ▶ Links: # of common neighbors
  - $C_1$   $\langle a, b, c, d, e \rangle$ :  $\{a, b, c\}, \{a, b, d\}, \{a, b, e\}, \{a, c, d\}, \{a, c, e\}, \{a, d, e\}, \{b, c, d\}, \{b, c, e\}, \{b, d, e\}, \{c, d, e\}$
  - $C_2$   $\langle a, b, f, g \rangle$ :  $\{a, b, f\}, \{a, b, g\}, \{a, f, g\}, \{b, f, g\}$
- ▶ Let  $T_1 = \{a, b, c\}$ ,  $T_2 = \{c, d, e\}$ ,  $T_3 = \{a, b, f\}$ 
  - $link(T_1, T_2) = 4$ , since they have 4 common neighbors
    - $\{a, c, d\}, \{a, c, e\}, \{b, c, d\}, \{b, c, e\}$
  - $link(T_1, T_3) = 3$ , since they have 3 common neighbors
    - $\{a, b, d\}, \{a, b, e\}, \{a, b, g\}$
- ▶ Thus link is a better measure than Jaccard coefficient



# CHAMELEON

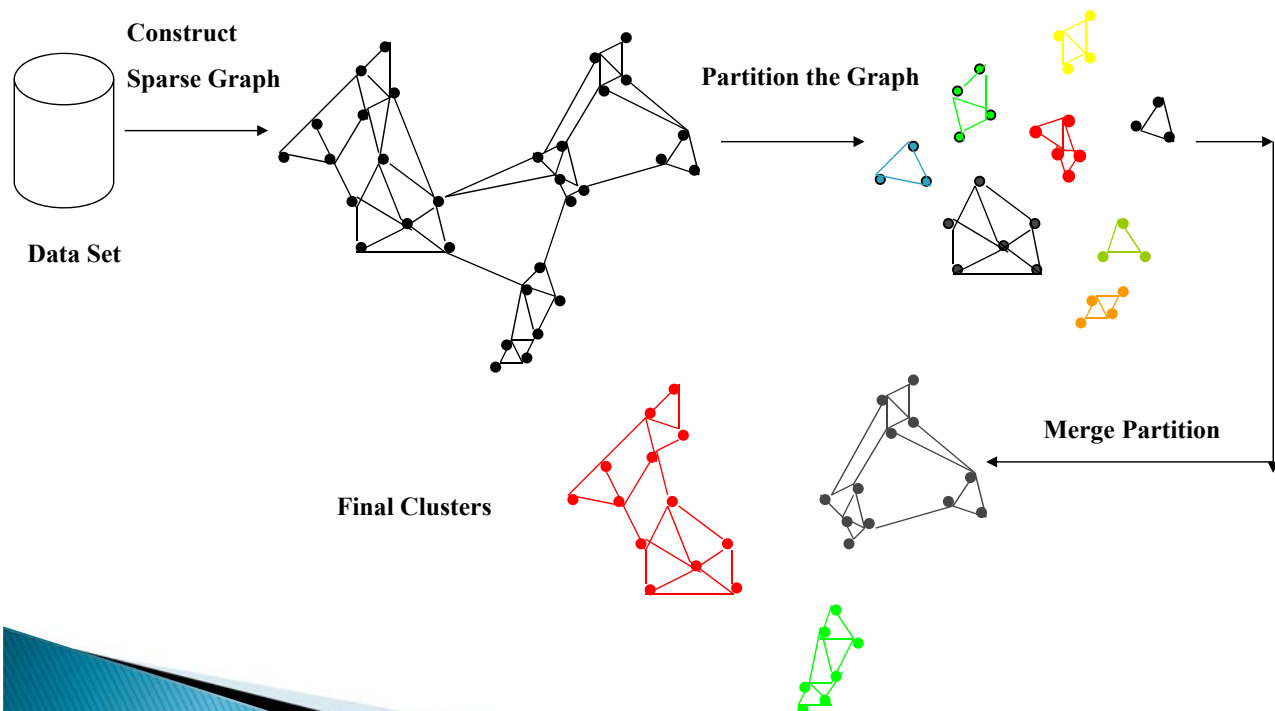
- ▶ Hierarchical Clustering Using Dynamic Modeling by G. Karypis, E.H. Han, and V. Kumar'99
- ▶ Measures the similarity based on a dynamic model
  - Two clusters are merged only if the *interconnectivity* and *closeness (proximity)* between two clusters are high *relative to* the internal interconnectivity of the clusters and closeness of items within the clusters
  - Cure ignores information about *interconnectivity* of the objects, Rock ignores information about the *closeness* of two clusters

# CHAMELEON

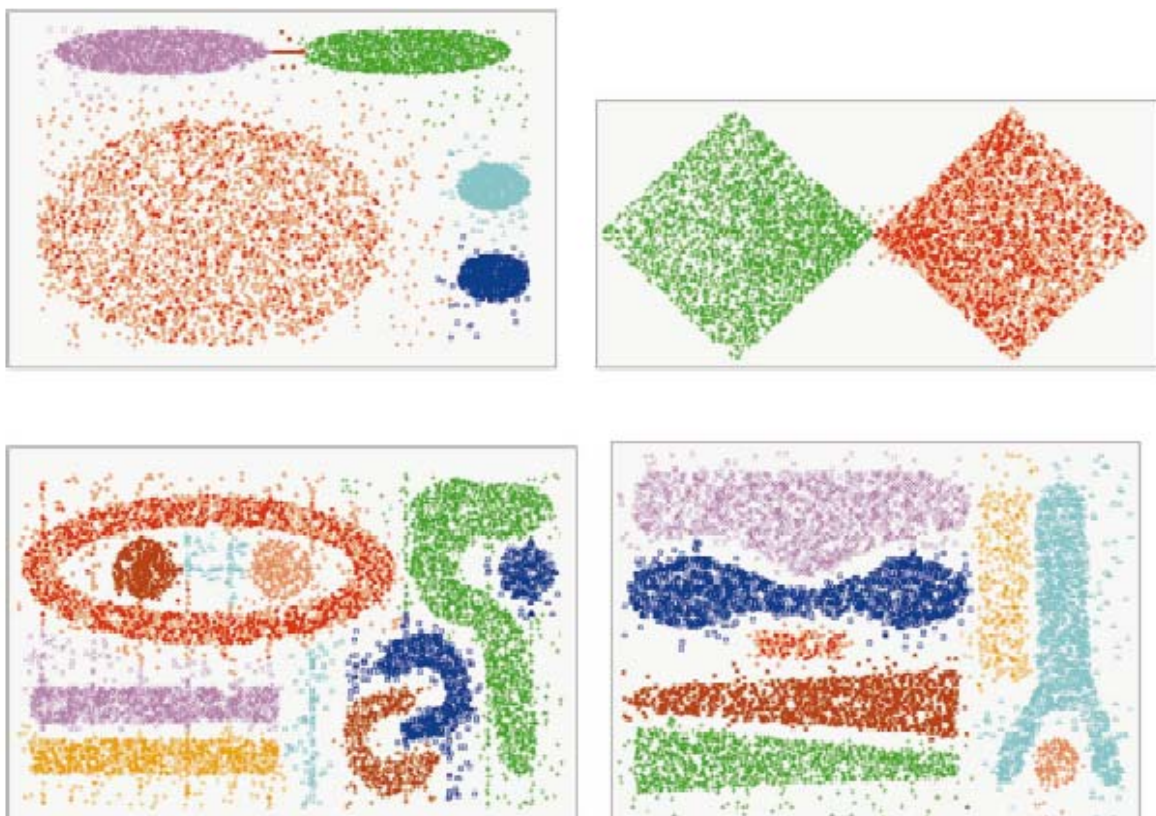
- ▶ A two-phase algorithm
  1. Use a graph partitioning algorithm: cluster objects into a large number of relatively small sub-clusters
  2. Use an agglomerative hierarchical clustering algorithm: find the genuine clusters by repeatedly combining these sub-clusters



# Framework of CHAMELEON



## Clustering Complex Objects



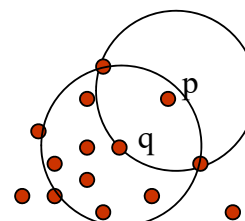
# Density-Based Clustering

- ▶ Major features:
  - Discover clusters of arbitrary shape
  - Handle noise
  - One scan
  - Need density parameters as termination condition
- ▶ Several interesting studies:
  - DBSCAN: Ester, et al. (KDD'96)
  - OPTICS: Ankerst, et al (SIGMOD'99).
  - DENCLUE: Hinneburg & D. Keim (KDD'98)
  - CLIQUE: Agrawal, et al. (SIGMOD'98) (more grid-based)

## Basic Concepts

- ▶ Two parameters:
  - *Eps*: Maximum radius of the neighbourhood
  - *MinPts*: Minimum number of points in an *Eps*-neighbourhood of that point
- ▶  $N_{Eps}(p)$ :  $\{q \text{ belongs to } D \mid dist(p,q) \leq Eps\}$
- ▶ **Directly density-reachable**: A point  $p$  is directly density-reachable from a point  $q$  w.r.t.  $Eps$ ,  $MinPts$  if
  - $p$  belongs to  $N_{Eps}(q)$
  - core point condition:

$$|N_{Eps}(q)| \geq MinPts$$

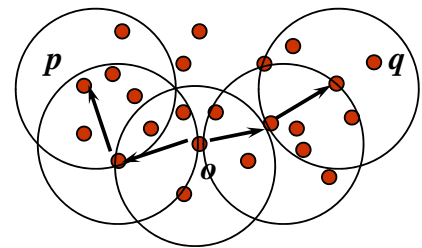
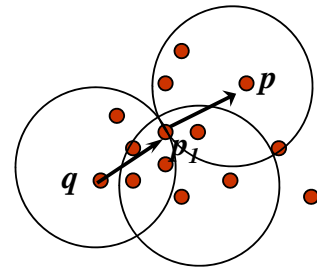


MinPts = 5

Eps = 1 cm

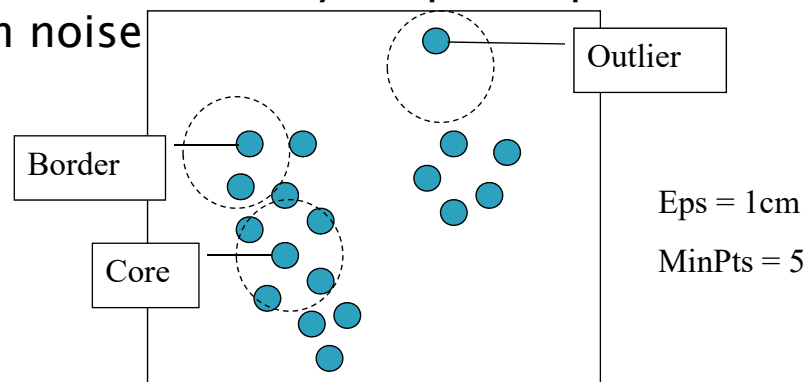
# Density-Reachable and Density-Connected

- ▶ Density-reachable:
  - A point  $p$  is **density-reachable** from a point  $q$  w.r.t.  $Eps$ ,  $MinPts$  if there is a chain of points  $p_1, \dots, p_n$ ,  $p_1 = q$ ,  $p_n = p$  such that  $p_{i+1}$  is directly density-reachable from  $p_i$
- ▶ Density-connected:
  - A point  $p$  is **density-connected** to a point  $q$  w.r.t.  $Eps$ ,  $MinPts$  if there is a point  $o$  such that both,  $p$  and  $q$  are density-reachable from  $o$  w.r.t.  $Eps$  and  $MinPts$



## DBSCAN

- ▶ Density Based Spatial Clustering of Applications with Noise
- ▶ Relies on a *density-based* notion of cluster: A *cluster* is defined as a maximal set of density-connected points
- ▶ Discovers clusters of arbitrary shape in spatial databases with noise



# DBSCAN

- ▶ Arbitrary select a point  $p$
- ▶ Retrieve all points density-reachable from  $p$  w.r.t.  $Eps$  and  $MinPts$ .
- ▶ If  $p$  is a core point, a cluster is formed.
- ▶ If  $p$  is a border point, no points are density-reachable from  $p$  and DBSCAN visits the next point of the database.
- ▶ Continue the process until all of the points have been processed.

## Sensitive to Parameters

Figure 8. DBScan results for DS1 with  $MinPts$  at 4 and  $Eps$  at (a) 0.5 and (b) 0.4.

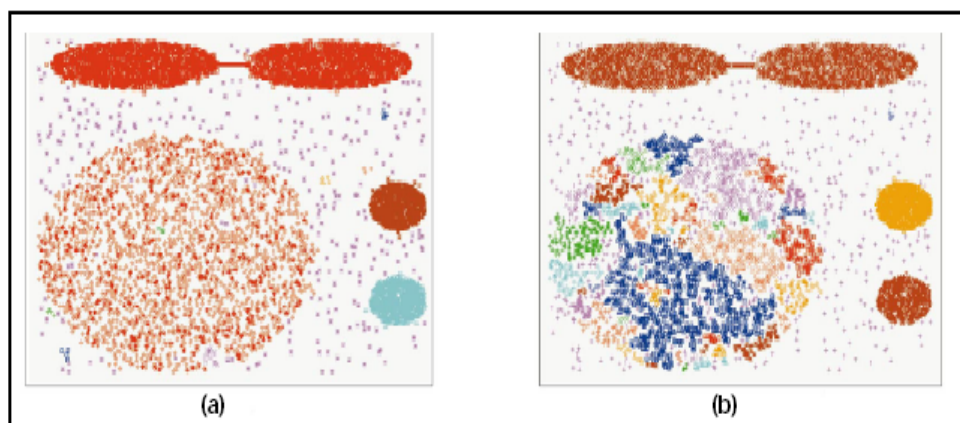
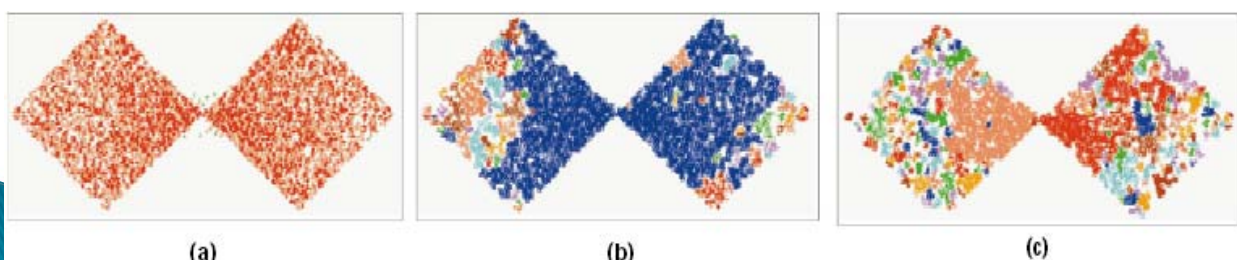


Figure 9. DBScan results for DS2 with  $MinPts$  at 4 and  $Eps$  at (a) 5.0, (b) 3.5, and (c) 3.0.





# OPTICS

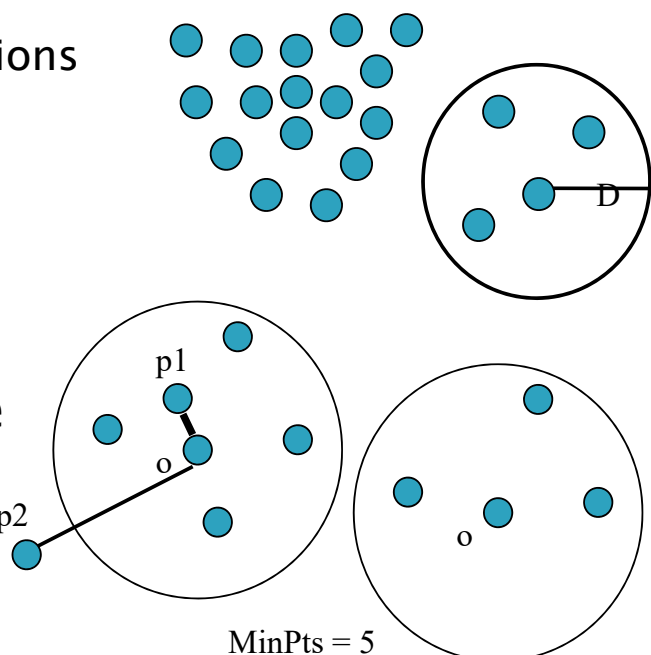
- ▶ Ordering Points To Identify the Clustering Structure, Ankerst, Breunig, Kriegel, and Sander (SIGMOD'99)
- ▶ Produces a special order of the database wrt its density-based clustering structure
- ▶ This cluster-ordering contains info equiv to the density-based clusterings corresponding to a broad range of parameter settings
- ▶ Good for both automatic and interactive cluster analysis, including finding intrinsic clustering structure
- ▶ Can be represented graphically or using visualization techniques

## Some Extension from DBSCAN

- ▶ Index-based:
  - $k$  = number of dimensions
  - $N = 20$
  - $p = 75\%$
  - $M = N(1-p) = 5$
  - Complexity:  $O(kN^2)$
- ▶ Core Distance
- ▶ Reachability Distance

$\text{Max}(\text{core-distance}(o), d(o, p))$

$r(p1, o) = 2.8\text{cm}$ .  $r(p2, o) = 4\text{cm}$



$\epsilon = 3\text{ cm}$

# DENCLUE

- ▶ DENSity-based CLUstEring by Hinneburg & Keim (KDD'98)
- ▶ Using statistical density functions:

$$f_{\text{Gaussian}}(x, y) = e^{-\frac{d(x, y)^2}{2\sigma^2}}$$

$$f_{\text{Gaussian}}^D(x) = \sum_{i=1}^N e^{-\frac{d(x, x_i)^2}{2\sigma^2}}$$

- ▶ Major features

$$\nabla f_{\text{Gaussian}}^D(x, x_i) = \sum_{i=1}^N (x_i - x) \cdot e^{-\frac{d(x, x_i)^2}{2\sigma^2}}$$

- Solid mathematical foundation
- Good for data sets with large amounts of noise
- Allows a compact mathematical description of arbitrarily shaped clusters in high-dimensional data sets
- Significant faster than existing algorithm (e.g., DBSCAN)
- But needs a large number of parameters

# DENCLUE

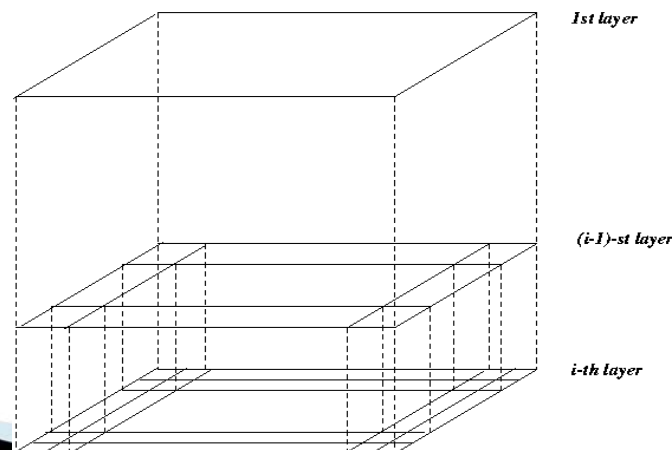
- ▶ Uses grid cells but only keeps information about grid cells that do actually contain data points and manages these cells in a tree-based access structure
- ▶ Influence function: describes the impact of a data point within its neighborhood
- ▶ Overall density of the data space can be calculated as the sum of the influence function of all data points
- ▶ Clusters can be determined mathematically by identifying density attractors
- ▶ Density attractors are local maximal of the overall density function

# Grid-Based Clustering

- ▶ 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 is divided into rectangular cells
- ▶ There are several levels of cells corresponding to different levels of resolution





# STING

- ▶ 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

# STING

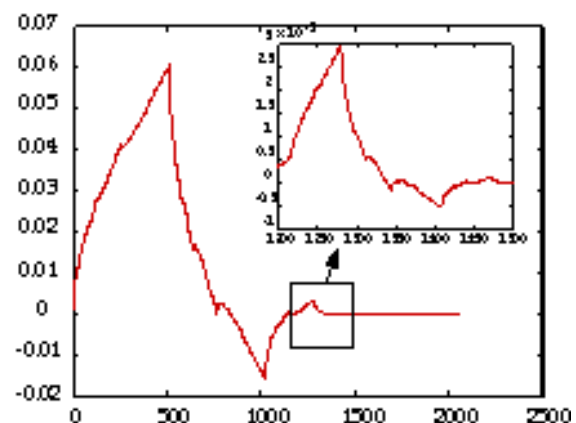
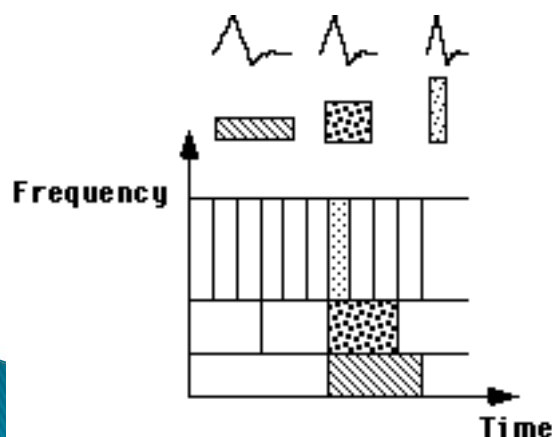
- ▶ Remove the irrelevant cells from further consideration
- ▶ When finish examining the current layer, proceed to the next lower level
- ▶ Repeat this process until the bottom layer is reached
- ▶ Advantages:
  - Query-independent, easy to parallelize, incremental update
  - $O(K)$ , where  $K$  is 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, 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



# WaveCluster

- ▶ 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

## Model-Based Clustering

- ▶ What is model-based clustering?
  - Attempt to optimize the fit between the given data and some mathematical model
  - Based on the assumption: Data are generated by a mixture of underlying probability distribution
- ▶ Typical methods
  - Statistical approach
    - EM (Expectation maximization), AutoClass
  - Machine learning approach
    - COBWEB, CLASSIT
  - Neural network approach
    - SOM (Self-Organizing Feature Map)

# EM

- ▶ Expectation Maximization (A popular iterative refinement algorithm)
- ▶ General idea
  - Starts with an initial estimate of the parameter vector
  - Iteratively rescores the patterns against the mixture density produced by the parameter vector
  - The rescored patterns are used to update the parameter updates
  - Patterns belonging to the same cluster, if they are placed by their scores in a particular component
- ▶ Algorithm converges fast but may not be in global optima

# EM

- ▶ Initially, randomly assign k cluster centers
- ▶ Iteratively refine the clusters based on two steps
  - Expectation step: assign each data point  $X_i$  to cluster  $C_i$  with the following probability

$$P(X_i \in C_k) = p(C_k|X_i) = \frac{p(C_k)p(X_i|C_k)}{p(X_i)},$$

- Maximization step:
  - Estimation of model parameters

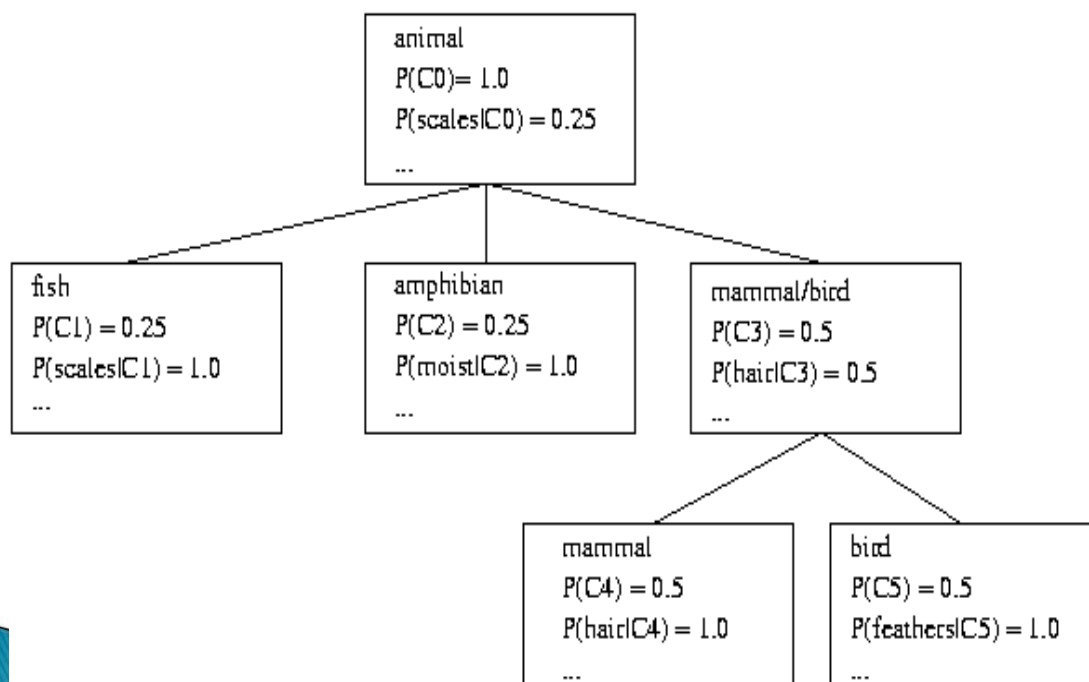
$$m_k = \frac{1}{N} \sum_{i=1}^N \frac{X_i P(X_i \in C_k)}{\sum_j P(X_i \in C_j)}.$$

# Conceptual Clustering

- ▶ Conceptual clustering
  - A form of clustering in machine learning
  - Produces a classification scheme for a set of unlabeled objects
  - Finds characteristic description for each concept (class)
- ▶ COBWEB (Fisher'87)
  - A popular and simple method of incremental conceptual learning
  - Creates a hierarchical clustering in the form of a **classification tree**
  - Each node refers to a concept and contains a probabilistic description of that concept

## COBWEB

- ▶ **A classification tree**



# Limitations of COBWEB

- ▶ The assumption that the attributes are independent of each other is often too strong because correlation may exist
- ▶ Not suitable for clustering large database data – skewed tree and expensive probability distributions
- ▶ CLASSIT
  - an extension of COBWEB for incremental clustering of continuous data
  - suffers similar problems as COBWEB
- ▶ AutoClass (Cheeseman and Stutz, 1996)
  - Uses Bayesian statistical analysis to estimate the number of clusters
  - Popular in industry

# Neural Network

- ▶ Neural network approaches
  - Represent each cluster as an exemplar, acting as a “prototype” of the cluster
  - New objects are distributed to the cluster whose exemplar is the most similar according to some distance measure
- ▶ Typical methods
  - SOM (Soft-Organizing feature Map)
  - Competitive learning
    - Involves a hierarchical architecture of several units (neurons)
    - Neurons compete in a “winner-takes-all” fashion for the object currently being presented

# SOM

- ▶ Self-Organizing Feature Map, also called topological ordered maps, or Kohonen Self-Organizing Feature Map (KSOMs)
- ▶ It maps all the points in a high-dimensional source space into a 2 to 3-d target space, s.t., the distance and proximity relationship (i.e., topology) are preserved as much as possible
- ▶ Similar to k-means: cluster centers tend to lie in a low-dimensional manifold in the feature space

# SOM

- ▶ Clustering is performed by having several units competing for the current object
  - The unit whose weight vector is closest to the current object wins
  - The winner and its neighbors learn by having their weights adjusted
- ▶ SOMs are believed to resemble processing that can occur in the brain
- ▶ Useful for visualizing high-dimensional data in 2- or 3-D space



# References

- ▶ Slides from Prof. J.-W. Han, UIUC
- ▶ Slides from Prof. M.-S. Chen, NTU
- ▶ Slides from Prof. W.-Z. Peng, NCTU