Università degli Studi di Milano

Master Degree in Computer Science

# Information Management course

Teacher: Alberto Ceselli

Lecture 20: 09/12/2014

# Data Mining:
## Concepts and Techniques
### (3rd ed.)

### — Chapter 10 —

Jiawei Han, Micheline Kamber, and Jian Pei

University of Illinois at Urbana-Champaign &

Simon Fraser University

# Cluster Analysis: Concepts and Methods

- Cluster Analysis: Basic Concepts

- Partitioning Methods

- Hierarchical Methods

- Density-Based Methods

- Grid-Based Methods

- Evaluation of Clustering

- Summary

# What is Cluster Analysis?

- Cluster: A collection of data objects
  - similar (or related) to one another within the same group
  - dissimilar (or unrelated) to the objects in other groups
- Cluster analysis (or *clustering*, *data segmentation, …*)
  - Finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters
- Unsupervised learning: no predefined classes (i.e., *try to learn by extracting regularities in data*)
- Typical applications
  - As a stand-alone tool to get insight into data distribution
  - As a preprocessing step for other algorithms

# Applications of Cluster Analysis

- Data reduction
  - Summarization: Preprocessing for regression, PCA, classification, and association analysis
  - Compression: Image processing: vector quantization
- Hypothesis generation and testing
- Prediction based on groups
  - Cluster & find characteristics/patterns for each group
- Finding K-nearest Neighbors
  - Localizing search to one or a small number of clusters
- Outlier detection: Outliers are often viewed as those "far away" from any cluster

# Clustering: Application Examples

- Biology: taxonomy of living things: kingdom, phylum, class, order, family, genus and species
- Information retrieval: document clustering
- Land use: Identification of areas of similar land use in an earth observation database
- Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
- City-planning: Identifying groups of houses according to their house type, value, and geographical location
- Earth-quake studies: Observed earth quake epicenters should be clustered along continent faults
- Climate: understanding earth climate, find patterns of atmospheric and ocean
- Economic Science: market research

# Basic Steps to Develop a Clustering Task

- Feature selection
  - Select info concerning the task of interest
  - Minimal information redundancy
- Proximity measure
  - Similarity of two feature vectors
- Clustering criterion
  - Expressed via a cost function or some rules
- Clustering algorithms
  - Choice of algorithms
- Validation of the results
  - Validation test (also, *clustering tendency* test)
- Interpretation of the results
  - Integration with applications

# Quality: What Is Good Clustering?

- A <u>good clustering</u> method will produce high quality clusters
  - high <u>intra-class</u> similarity: <span style="color:red">cohesive</span> within clusters
  - low <u>inter-class</u> similarity: <span style="color:red">distinctive</span> between clusters
- The <u>quality</u> of a clustering method depends on
  - the similarity measure used by the method
  - its implementation (optimality guarantees + computational effectiveness), and
  - Its ability to discover some or all of the <u>hidden</u> patterns (practical behavior)

# Measure the Quality of Clustering

- Dissimilarity/Similarity metric
  - Similarity is expressed in terms of a (typically metric) <u>pairwise</u> distance function $d(i, j)$
  - The definitions of distance functions are usually rather different for interval-scaled, boolean, categorical, ordinal ratio, and vector variables
  - Weights should be associated with different variables based on applications and data semantics
- Quality of clustering:
  - There is usually a separate <u>global quality function</u> that measures the "goodness" of a cluster.
  - It is hard to define "similar enough" or "good enough" (need to stick to the application!)
    - The answer is typically highly subjective (i.e. don't blame the *algorithm* for *modeling* errors)

# Major Clustering Approaches

- Partitioning approach:
  - Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors
  - Typical methods: k-means, k-medoids, CLARANS
- Hierarchical approach:
  - Create a hierarchical decomposition of the set of data (or objects) using some criterion
  - Agglomerative (bottom-up) or divisive (top-down)
- Density-based approach:
  - Based on connectivity and density functions (keep growing as points are still in the neighborhood of cluster elements)
  - Find arbitrarily shaped clusters
- Grid-based approach:
  - Quantize object space in a grid structure
  - build a multiple-level granularity structure

# Cluster Analysis: Basic Concepts and Methods

- Cluster Analysis: Basic Concepts

- Partitioning Methods

- Hierarchical Methods

- Density-Based Methods

- Grid-Based Methods

- Evaluation of Clustering

- Summary

# Partitioning Algorithms: Basic Concept

Partitioning method:

- given the number of clusters k
- given a dissimilarity measure (partitioning criterion)
- given a database $D$ of $n$ objects
- partition it into a set of $k$ clusters
- such that the sum of dissimilarities with respect to a cluster representative $c_i$ is minimized (e.g. squared distances to the centroid or medoid of cluster $C_i$)

$$E = \sum_{i=1}^{k} \sum_{p \in C_i} \left( d\left( p, c_i \right) \right)^2$$

# Partitioning Algorithms: Basic Concept

Partitioning method:

.

- Exact (globally opt) methods: mixed integer programming
- Heuristic methods: *k-means* and *k-medoids* algorithms
- *k-means* (MacQueen'67, Lloyd'57/'82): Each cluster is represented by the center (attibute-wise means) of the cluster
- *k-medoids* or k-medians or PAM (Partition Around Medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

# The *K-Means* Clustering Method

- Given *k*, the *k-means* algorithm is implemented in four steps:

  1) Partition objects into *k* nonempty subsets

  2) Compute seed points as the centroids of the clusters of the current partitioning (the centroid is the center, i.e., *mean point*, of the cluster)

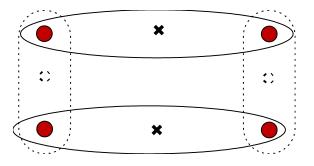  3) Assign each object to the cluster with the nearest seed point

  4) Go back to Step 2, stop when the assignment does not change

# Comments on the *K-Means* Method

- <u>Strength:</u> *Efficient*: $O(tkn)$, where
  - *n* is # objects, *k* is # clusters, and *t* is # iterations.
  - Normally, *k*, *t* << *n*.
  - Comparing: PAM: $O(k(n-k)^2)$, CLARA: $O(ks^2 + k(n-k))$
- <u>Weakness</u>
  - Heuristc; often terminates at a *local optimal*
  - Applicable only to objects in a continuous n-dimensional space
    - Using the k-modes for categorical data
    - Using the k-medoids for a wider range of data
  - Need to give *k,* the *number* of clusters, as input (there are ways to guess meaningful k, see Hastie et al. 2009)
  - Sensitive to noisy data and *outliers*
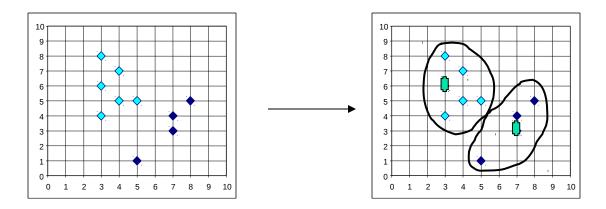  - Not suitable to discover clusters with *non-convex shapes*

# Variations of the *K-Means* Method

- Most of the variants of the *k-means* which differ in

    - Selection of the initial *k* means

    - Dissimilarity calculations

    - Strategies to calculate cluster means

- Handling categorical data: *k-modes*

    - Replacing means of clusters with <u>modes</u>

    - Using new dissimilarity measures to deal with categorical objects

    - Using a <u>frequency</u>-based method to update modes of clusters

    - A mixture of categorical and numerical data: *k-prototype* method
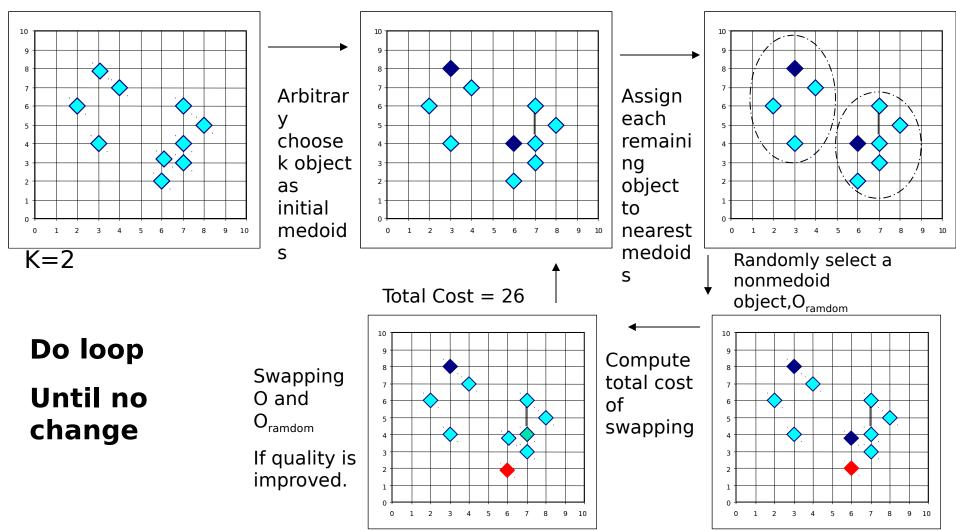
# What Is the Problem of the K-Means Method?

- The k-means algorithm is sensitive to outliers !

  - Since an object with an extremely large value may substantially distort the distribution of the data

- K-Medoids:  Instead of taking the **mean** value of the object in a cluster as a reference point, **medoids** can be used, which is the **most centrally located** object in a cluster

# PAM: A Typical K-Medoids Algorithm

Total Cost = 20



Arbitrary choose k object as initial medoids

K=2

Assign each remaining object to nearest medoids

Randomly select a nonmedoid object,$O_{ramdom}$

**Do loop**

**Until no change**

Total Cost = 26

Swapping O and $O_{ramdom}$

If quality is improved.

Compute total cost of swapping

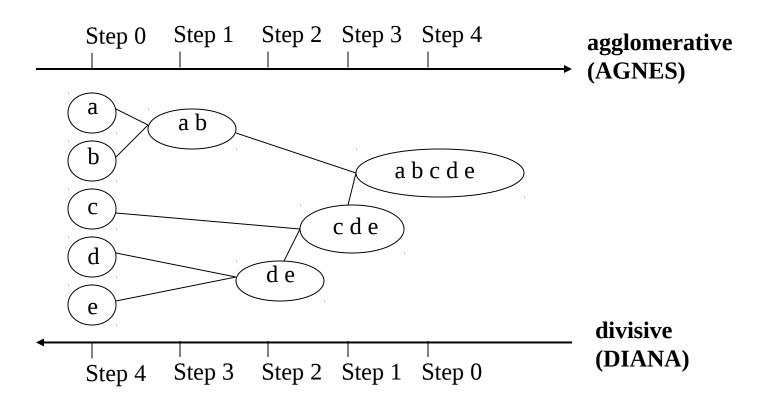# The K-Medoid Clustering Method

- *K-Medoids* Clustering: Find *representative* objects (<u>medoids</u>) in clusters

  - *PAM* (Partitioning Around Medoids, Kaufmann & Rousseeuw 1987)

    - Starts from an initial set of medoids and iteratively replaces one of the medoids by one of the non-medoids if it improves the total distance of the resulting clustering

    - *PAM* works effectively for small data sets, but does not scale well for large data sets (due to the computational complexity)

- Efficiency improvement on PAM

  - *CLARA* (Kaufmann & Rousseeuw, 1990): PAM on samples

  - *CLARANS* (Ng & Han, 1994): Randomized re-sampling

# Chapter 10. Cluster Analysis: Basic Concepts and Methods

- Cluster Analysis: Basic Concepts

- Partitioning Methods

- Hierarchical Methods ⬅

- Density-Based Methods

- Grid-Based Methods
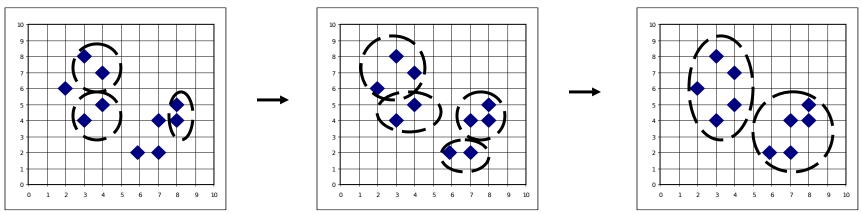
- Evaluation of Clustering

- Summary

# 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



Step 0  Step 1  Step 2  Step 3  Step 4    **agglomerative (AGNES)**

a
b
a b
a b c d e
c
d
e
c d e
d e

Step 4  Step 3  Step 2  Step 1  Step 0    **divisive (DIANA)**
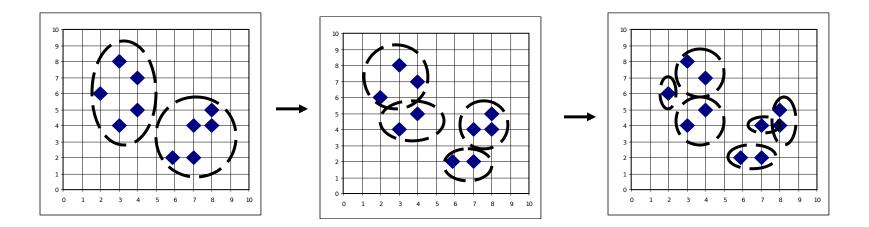
25

# AGNES (AGglomerative NESting)

- Introduced in Kaufmann and Rousseeuw (1990)
- Implemented in statistical 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
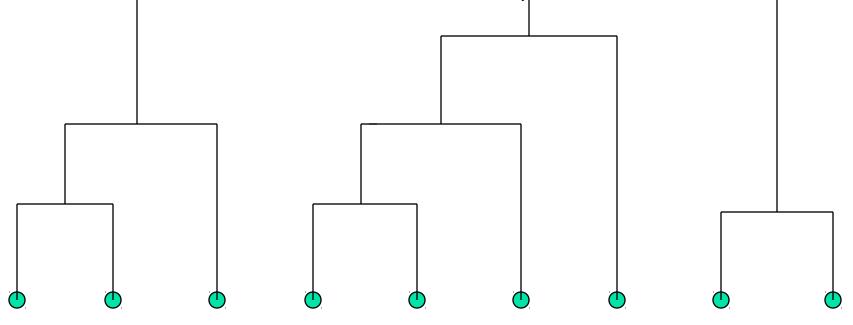
# DIANA (DIvisive ANAlysis)

- 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



27

# *Dendrogram:* Shows How Clusters are Merged

Decompose data objects into a several levels of nested partitioning (<u>tree</u> of clusters), called a <u>dendrogram</u>

A <u>clustering</u> of the data objects is obtained by <u>cutting</u> the dendrogram at the desired level, then each <u>connected component</u> forms a cluster
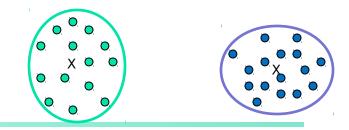
# Distance between Clusters

- Single link:  smallest distance between an element in one cluster and an element in the other, i.e.,  $\text{dist}(K_i, K_j) = \min_{p \text{ in } Ki, q \text{ in } Kj} d(p,q)$

- Complete link: largest distance between an element in one cluster and an element in the other, i.e., $\text{dist}(K_i, K_j) = \max_{p \text{ in } Ki, q \text{ in } Kj} d(p,q)$

- Average: avg distance between an element in one cluster and an element in the other, i.e., $\text{dist}(K_i, K_j) = \text{sum}_{p \text{ in } Ki, q \text{ in } Kj} d(p,q) / (|Ki||Kj|)$

- Centroid: distance between the centroids of two clusters, e.g., $p = \text{mean}(Ki), q = \text{mean}(Kj), \text{dist}(K_i, K_j) = d(p,q)$

- Medoid: distance between the medoids of two clusters, i.e., $p = \text{median}(Ki), q = \text{median}(Kj), \text{dist}(K_i, K_j) = d(p,q)$

# Distance between Clusters

- Algorithms using minimum distance are also called nearest-neighbor clustering algorithms
  - they build minimum spanning trees
  - if clustering is terminated when the minimum inter-cluster distance exceeds a given threshold they are called single-linkage
- Algorithms using maximum distance are also called farthest-neighbor clustering algorithms
  - If clustering is terminated when the maximum inter-cluster distance between nearest clusters exceeds a given threshold they are called complete-linkage

# Centroid, Radius and Diameter of a Cluster (for numerical data sets)

- Centroid: the "middle" of a cluster

$$C_m = \frac{\sum_{i=1}^{N}(t_{ip})}{N}$$

- Radius: square root of average distance from any point of the cluster to its centroid

$$R_m = \sqrt{\frac{\sum_{i=1}^{N}(t_{ip} - c_m)^2}{N}}$$

- Diameter: square root of average mean squared distance between all pairs of points in the cluster

$$D_m = \sqrt{\frac{\sum_{i=1}^{N}\sum_{i=1}^{N}(t_{ip} - t_{iq})^2}{N(N-1)}}$$

# Extensions to Hierarchical Clustering

- Major weakness of agglomerative clustering methods

  - <u>Can never undo what was done previously</u>

  - <u>Do not scale</u> well: time complexity of at least $O(n^2)$, where $n$ is the number of total objects

- Integration of hierarchical & distance-based clustering

  - <u>BIRCH (1996)</u>: uses CF-tree and incrementally adjusts the quality of sub-clusters

  - <u>CHAMELEON (1999)</u>: hierarchical clustering using dynamic modeling

# BIRCH (Balanced Iterative Reducing and Clustering Using Hierarchies)

- Zhang, Ramakrishnan & Livny, SIGMOD'96

- Clustering Feature (CF): <n, LS, SS>
  - n: number of points, LS: their sum, SS: their sum of squares
- Easy to compute centroid, radius and diameter from CF
- CFs are <u>additive</u>

- Incrementally construct a CF 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 its inherent clustering structure)

  - Phase 2: use an arbitrary clustering algorithm to cluster the leaf nodes of the CF-tree

# BIRCH (Balanced Iterative Reducing and Clustering Using Hierarchies)

- *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 Feature Vector in BIRCH
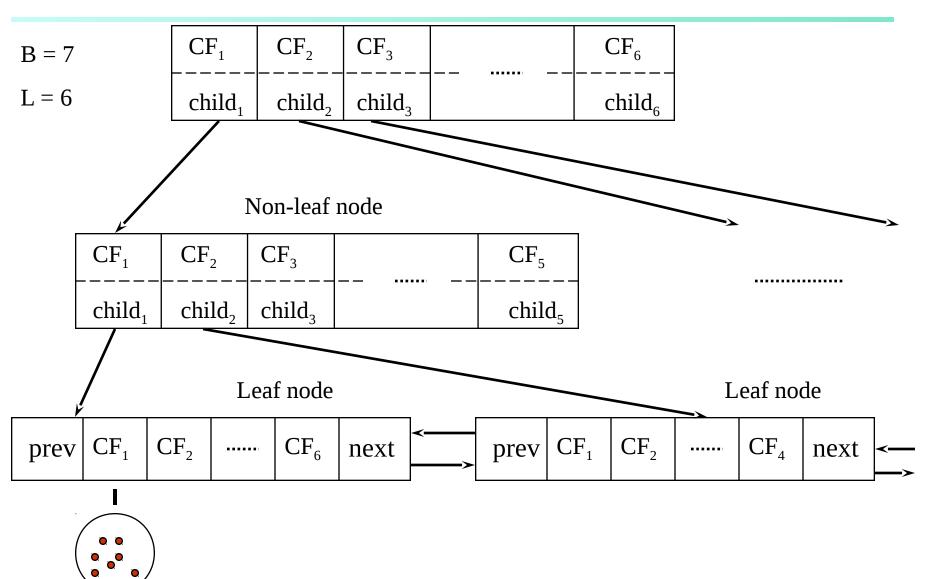
**Clustering Feature (CF):** *CF = (N, LS, SS)*

*N*: **Number of data points**

*LS: linear sum of N points:*

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

*SS: square sum of N points*

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

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



(3,4)
(2,6)
(4,5)
(4,7)
(3,8)

35

# CF-Tree in BIRCH

- Clustering feature:

  - Summary of the statistics for a given subcluster: the 0-th, 1st, and 2nd moments of the subcluster from the statistical point of view

  - Registers crucial measurements for computing cluster and utilizes storage efficiently

- 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: max # of children

  - Threshold: max diameter of sub-clusters stored at the leaf nodes

# The CF Tree Structure

Root

$B = 7$

$L = 6$

| $CF_1$ | $CF_2$ | $CF_3$ | | ...... | $CF_6$ |
|--------|--------|--------|---|---|--------|
| child$_1$ | child$_2$ | child$_3$ | | | child$_6$ |

Non-leaf node

| $CF_1$ | $CF_2$ | $CF_3$ | | ...... | $CF_5$ |
|--------|--------|--------|---|---|--------|
| child$_1$ | child$_2$ | child$_3$ | | | child$_5$ |

.................

Leaf node

Leaf node

| prev | $CF_1$ | $CF_2$ | ...... | $CF_6$ | next |
|------|--------|--------|--------|--------|------|

| prev | $CF_1$ | $CF_2$ | ...... | $CF_4$ | next |
|------|--------|--------|--------|--------|------|

# The Birch Algorithm

- Cluster Diameter

$$\sqrt{\frac{1}{n(n-1)}\sum(x_i-x_j)^2}$$

- For each point in the input
  - Find closest leaf entry
  - Add point to leaf entry and update CF
  - If entry diameter > max_diameter, then split leaf, and possibly parents
- Algorithm is O(n)
- Concerns
  - Sensitive to insertion order of data points
  - Since we fix the size of leaf nodes, so clusters may not be so natural
  - Clusters tend to be spherical given the radius and diameter measures

# CHAMELEON: Hierarchical Clustering Using Dynamic Modeling (1999)

- CHAMELEON: G. Karypis, E. H. Han, and V. Kumar, 1999
- 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
- Graph-based, and 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

# KNN Graphs & Interconnectivity

- k-nearest graphs from an original data in 2D:



(a) Original Data in 2D    (b) 1-nearest neighbor graph    (c) 2-nearest neighbor graph    (d) 3-nearest neighbor graph

- $EC_{\{Ci,Cj\}}$ : The absolute inter-connectivity between $C_i$ and $C_j$: *the sum* of the weight of the edges that connect vertices in $C_i$ *to vertices in* $C_j$

- Internal inter-connectivity of a cluster $C_i$ : *the size of its min-cut bisector $EC_{Ci}$ (i.e., the weighted sum* of edges that partition the graph into two roughly equal parts)

- Relative Inter-connectivity (RI):

$$RI(C_i, C_j) = \frac{|EC_{\{C_i,C_j\}}|}{\frac{|EC_{C_i}|+|EC_{C_j}|}{2}}$$

40

# Relative Closeness & Merge of Sub-Clusters

- **Relative closeness** between a pair of clusters $C_i$ and $C_j$ : the absolute closeness between $C_i$ and $C_j$ normalized w.r.t. the internal closeness of the two clusters $C_i$ and $C_j$
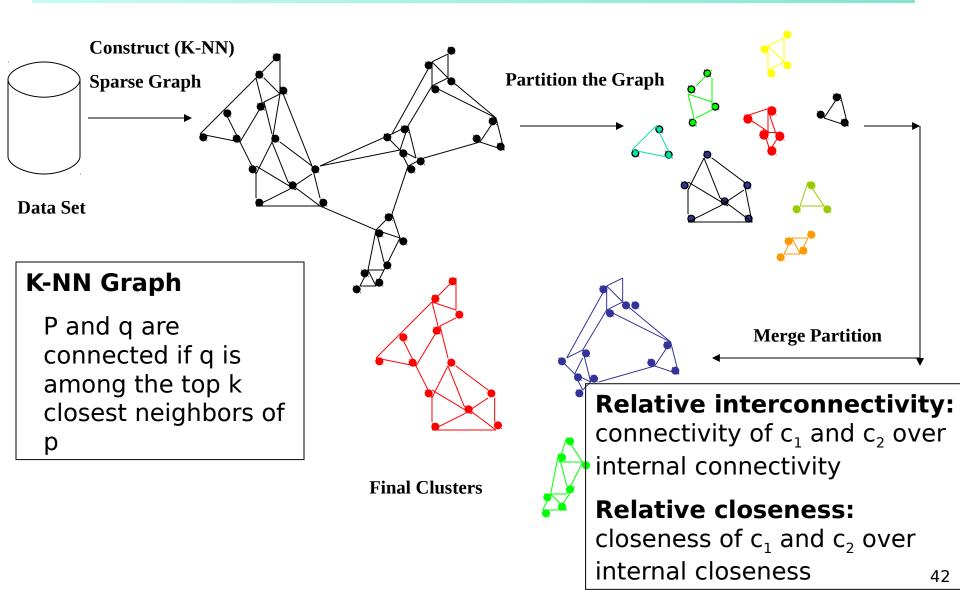
$$RC(C_i, C_j) = \frac{\overline{S}_{EC_{\{C_i,C_j\}}}}{\frac{|C_i|}{|C_i|+|C_j|}\overline{S}_{EC_{C_i}} + \frac{|C_j|}{|C_i|+|C_j|}\overline{S}_{EC_{C_j}}}$$

  - $\overline{S}_{EC_{C_i}}$ and $\overline{S}_{EC_{C_j}}$ are the average weights of the edges that belong in the min-cut bisector of clusters $C_i$ and $C_j$ , respectively, and $\overline{S}_{EC_{\{C_i,C_j\}}}$ is the average weight of the edges that connect vertices in $C_i$ to vertices in $C_j$
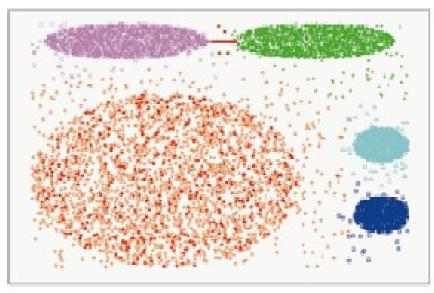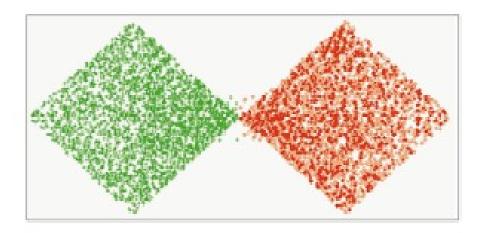
- **Merge Sub-Clusters:**
  - Merges only those pairs of clusters whose RI and RC are both above some user-specified thresholds
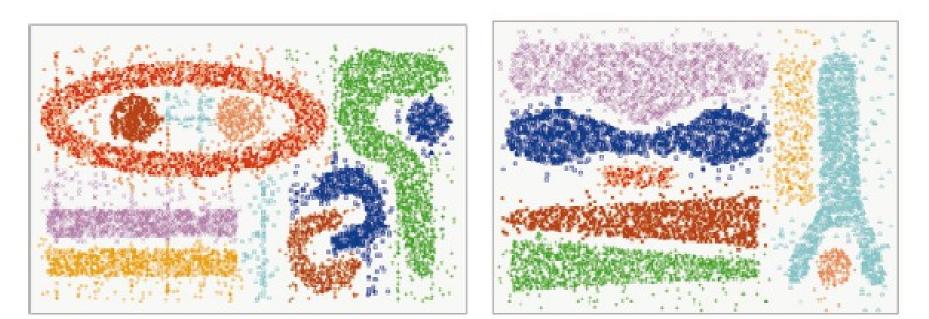  - Merge those maximizing a function combining RI & RC

# Overall Framework of CHAMELEON

**Construct (K-NN)**

**Sparse Graph**

**Data Set**

**Partition the Graph**

### K-NN Graph

P and q are connected if q is among the top k closest neighbors of p

**Final Clusters**

**Merge Partition**

**Relative interconnectivity:** connectivity of $c_1$ and $c_2$ over internal connectivity

**Relative closeness:** closeness of $c_1$ and $c_2$ over internal closeness

42

# CHAMELEON (Clustering Complex Objects)

# Hierarchical Clustering Summary

- Hierarchical clustering strengths
  - Produce at once clustering solutions for different k values
  - Link them, highlighting regularities
- Hierarchical clustering weaknesses
  - Nontrivial to choose a good distance measure
  - Hard to handle missing attribute values
  - Algorithmically (besides theoretically) hard: mainly heuristics in practical settings

# Probabilistic Hierarchical Clustering

- Hierarchical (distance-based) clustering strengths …

- Hierarchical (distance-based) clustering weaknesses …

- **Probabilistic** ("fitting") hierarchical clustering

  - Use probabilistic models to measure distances between clusters

  - Generative model: Regard the set of data objects to be clustered as a sample of the underlying data generation mechanism to be analyzed

  - Easy to understand, same efficiency as algorithmic agglomerative clustering method, can handle partially observed data

- In practice, assume the generative models adopt common distributions functions, e.g., Gaussian distribution or Bernoulli distribution, governed by parameters

# Generative Model

- Given a set of 1-D points $X = \{x_1, ..., x_n\}$ for clustering analysis & assuming they are generated by a Gaussian distribution:

$$\mathcal{N}(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- The probability that a point $x_i \in X$ is generated by the model

$$P(x_i|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

- The likelihood that $X$ is generated by the model:

$$L(\mathcal{N}(\mu, \sigma^2) : X) = P(X|\mu, \sigma^2) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

- The task of learning the generative model: find the parameters $\mu$ and $\sigma^2$ such that

$$\mathcal{N}(\mu_0, \sigma_0^2) = \arg\max\{L(\mathcal{N}(\mu, \sigma^2) : X)\}$$

the maximum likelihood

46

# A Probabilistic Hierarchical Clustering Algorithm

- For a set of objects partitioned into $m$ clusters $C_1, \ldots, C_m$, the quality can be measured by,

$$Q(\{C_1, \ldots, C_m\}) = \prod_{i=1}^{m} P(C_i)$$

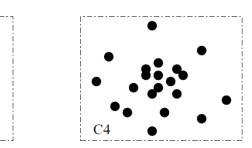  where $P()$ is the maximum likelihood

- If we merge two clusters $C_{j1}$ and $C_{j2}$ into a cluster $C_{j1} \cup C_{j2}$, then, the change in quality of the overall clustering is

$$Q((\{C_1, \ldots, C_m\} - \{C_{j_1}, C_{j_2}\}) \cup \{C_{j_1} \cup C_{j_2}\}) - Q(\{C_1, \ldots, C_m\})$$

$$= \frac{\prod_{i=1}^{m} P(C_i) \cdot P(C_{j_1} \cup C_{j_2})}{P(C_{j_1}) P(C_{j_2})} - \prod_{i=1}^{m} P(C_i)$$

$$= \prod_{i=1}^{m} P(C_i) \left( \frac{P(C_{j_1} \cup C_{j_2})}{P(C_{j_1}) P(C_{j_2})} - 1 \right)$$

- Distance between clusters $C_1$ and $C_2$:

$$dist(C_i, C_j) = -\log \frac{P(C_1 \cup C_2)}{P(C_1) P(C_2)}$$

# Cluster Analysis: Basic Concepts and Methods

- Cluster Analysis: Basic Concepts

- Partitioning Methods

- Hierarchical Methods

- Density-Based Methods ⬅

- Grid-Based Methods
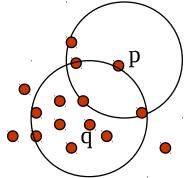
- Evaluation of Clustering

- Summary

# Density-Based Clustering Methods

- Clustering based on density (local cluster criterion), such as density-connected points
- 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)

# Density-Based Clustering: Basic Concepts

- Two parameters:

  - *Eps*: Maximum radius of the neighbourhood

  - *MinPts*: Minimum number of points in an Eps-neighbourhood of that point

- $N_{Eps}(q)$: {p belongs to D | dist(p,q) ≤ Eps}

- Directly density-reachable: A point *p* is directly density-reachable from a point *q* w.r.t. (*Eps*, *MinPts)* if

  - *p* belongs to $N_{Eps}(q)$
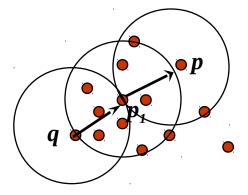
  - core point condition:

    $$|N_{Eps}(q)| \geq MinPts$$
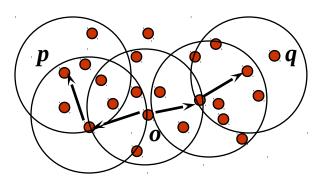
MinPts = 5

Eps = 1 cm

# Density-Reachable and Density-Connected

- Density-reachable:

  - A point *p* is <span style="color:red">density-reachable</span> from a point *q* w.r.t. (*Eps*, *MinPts)* if there is a chain of points $p_1$, ..., $p_n$, $p_1 = q$, $p_n = p$ such that $p_{i+1}$ is <u>directly</u> density-reachable from $p_i$
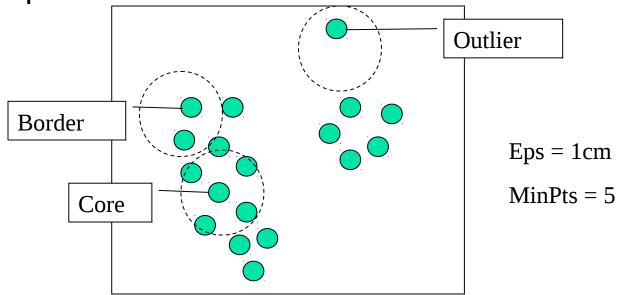
- Density-connected

  - A point *p* is <span style="color:red">density-connected</span> 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

- Experimentally, discovers clusters of arbitrary shape in spatial databases with noise

Outlier

Border

Core

Eps = 1cm

MinPts = 5

# DBSCAN: The Algorithm

- 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

- If a spatial index is used, the computational complexity of DBSCAN is O(nlogn), where n is the number of database objects. Otherwise, the complexity is O(n$^2$)
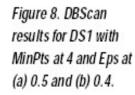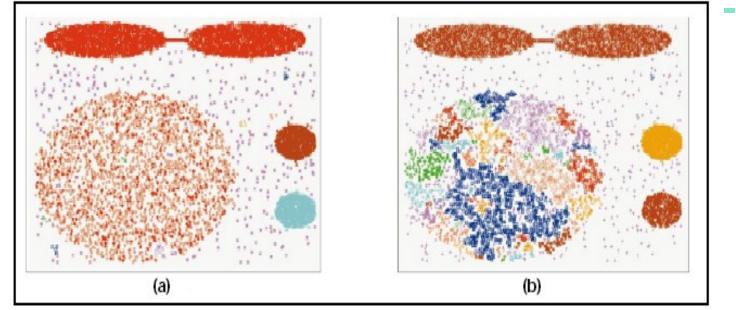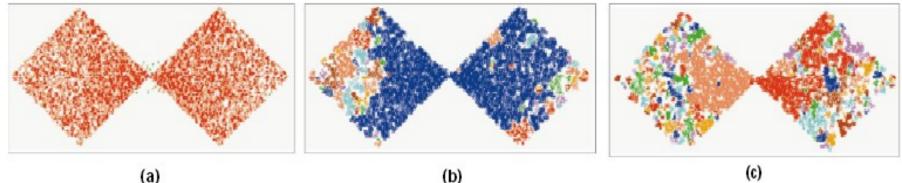
# DBSCAN: Sensitive to Parameters



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.

(a)   (b)

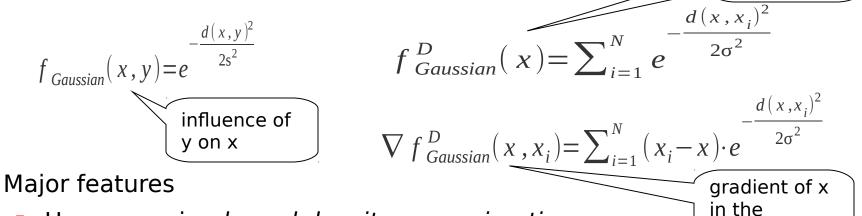(a)   (b)   (c)

**DBSCAN online Demo:**
http://webdocs.cs.ualberta.ca/~yaling/Cluster/Applet/Code/Cluster.html

# OPTICS: A Cluster-Ordering Method (1999)

- 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

# DENCLUE: Using Statistical Density Functions

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

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

influence of y on x

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

total influence on x

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

gradient of x in the direction of $x_i$

- Major features
  - Uses gaussian *kernel density approximation*

$$f(x) = \frac{1}{ns} \sum_{i=1}^{n} K\left(\frac{x - x_i}{s}\right)$$

  - Clusters can be determined mathematically by identifying density attractors (local maxima of the overall density function)
  - Center defined clusters: assign to each density attractor the points density attracted to it (pick each point and follow the gradient)
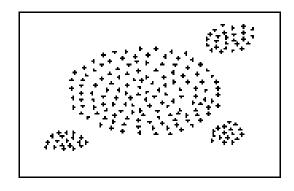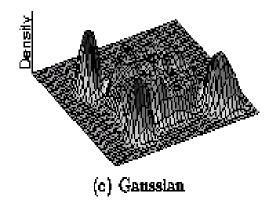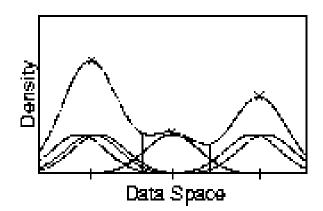
# Denclue: Technical Essence

- Arbitrary shaped cluster: merge density attractors that are connected through paths of high density (> threshold)
- 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)
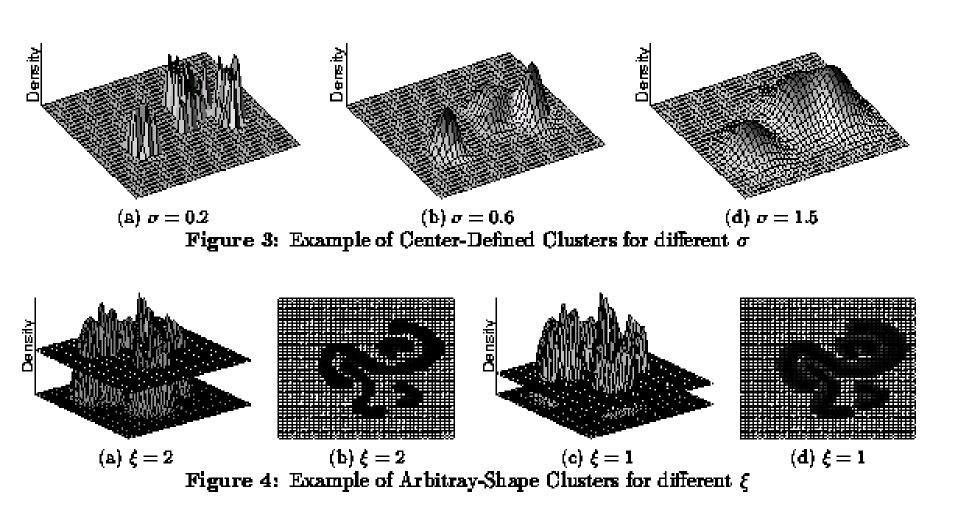- But needs a large number of parameters

# Density Attractor



(a) Data Set

(c) Gaussian

Density

Data Space

# Center-Defined and Arbitrary



(a) $\sigma = 0.2$     (b) $\sigma = 0.6$     (d) $\sigma = 1.5$

Figure 3: Example of Center-Defined Clusters for different $\sigma$



(a) $\xi = 2$     (b) $\xi = 2$     (c) $\xi = 1$     (d) $\xi = 1$

Figure 4: Example of Arbitray-Shape Clusters for different $\xi$

# Chapter 10. Cluster Analysis: Basic Concepts and Methods

- Cluster Analysis: Basic Concepts

- Partitioning Methods

- Hierarchical Methods

- Density-Based Methods

- Grid-Based Methods

- Evaluation of Clustering

- 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)
    - Both grid-based and subspace clustering

# 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

# 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
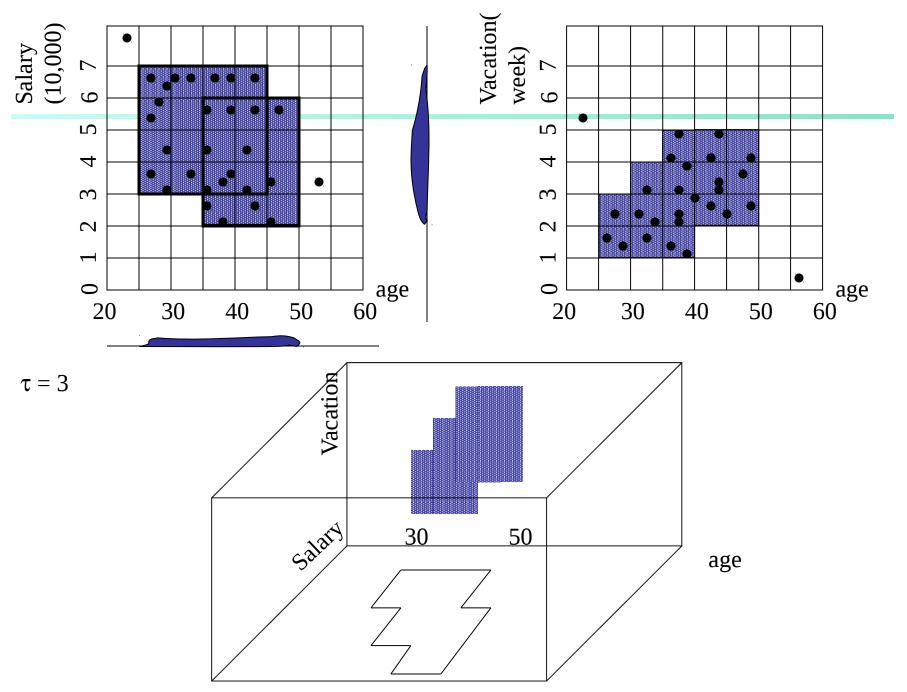
# STING Algorithm and Its Analysis

- 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

# CLIQUE (Clustering In QUEst)

- Agrawal, Gehrke, Gunopulos, Raghavan (SIGMOD'98)

- Automatically identifying subspaces of a high dimensional data space that allow better clustering than original space

- CLIQUE can be considered as both density-based and grid-based

  - It partitions each dimension into the same number of equal length interval

  - It partitions an m-dimensional data space into non-overlapping rectangular units

  - A unit is dense if the fraction of total data points contained in the unit exceeds the input model parameter

  - A cluster is a maximal set of connected dense units within a subspace

# CLIQUE: The Major Steps

- Partition the data space and find the number of points that lie inside each cell of the partition.

- Identify the subspaces that contain clusters using the Apriori principle

- Identify clusters

  - Determine dense units in all subspaces of interests

  - Determine connected dense units in all subspaces of interests.

- Generate minimal description for the clusters

  - Determine maximal regions that cover a cluster of connected dense units for each cluster

  - Determination of minimal cover for each cluster

$\tau = 3$

# Strength and Weakness of *CLIQUE*

- Strength
  - *automatically* finds subspaces of the highest dimensionality such that high density clusters exist in those subspaces
  - *insensitive* to the order of records in input and does not presume some canonical data distribution
  - scales *linearly* with the size of input and has good scalability as the number of dimensions in the data increases
- Weakness
  - The accuracy of the clustering result may be degraded at the expense of simplicity of the method