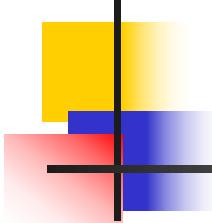


Data Mining: Concepts and Techniques

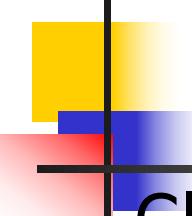
— Slides for Textbook —
— Chapter 8 —

©Jiawei Han and Micheline Kamber
Intelligent Database Systems Research Lab
School of Computing Science
Simon Fraser University, Canada
<http://www.cs.sfu.ca>



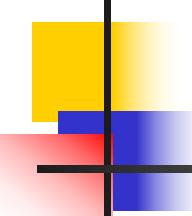
Chapter 8. Cluster Analysis

- n What is Cluster Analysis?
- n Types of Data in Cluster Analysis
- n A Categorization of Major Clustering Methods
- n Partitioning Methods
- n Hierarchical Methods
- n Density-Based Methods
- n Grid-Based Methods
- n Model-Based Clustering Methods
- n Outlier Analysis
- n Summary



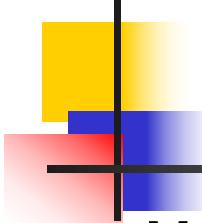
What is Cluster Analysis?

- n Cluster: a collection of data objects
 - n Similar to one another within the same cluster
 - n Dissimilar to the objects in other clusters
- n Cluster analysis
 - n Grouping a set of data objects into clusters
- n Clustering is **unsupervised classification**: no predefined classes
- n Typical applications
 - n As a **stand-alone tool** to get insight into data distribution
 - n As a **preprocessing step** for other algorithms



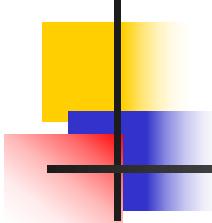
General Applications of Clustering

- n Pattern Recognition
- n Spatial Data Analysis
 - n create thematic maps in GIS by clustering feature spaces
 - n detect spatial clusters and explain them in spatial data mining
- n Image Processing
- n Economic Science (especially market research)
- n WWW
 - n Document classification
 - n Cluster Weblog data to discover groups of similar access patterns



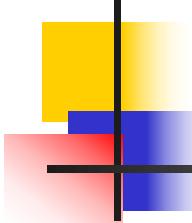
Examples of Clustering Applications

- n Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs
- n Land use: Identification of areas of similar land use in an earth observation database
- n Insurance: Identifying groups of motor insurance policy holders with a high average claim cost
- n City-planning: Identifying groups of houses according to their house type, value, and geographical location
- n Earthquake studies: Observed earth quake epicenters should be clustered along continent faults



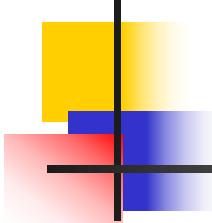
What Is Good Clustering?

- A good clustering method will produce high quality clusters with
 - high intra-class similarity
 - low inter-class similarity
- The quality of a clustering result depends on both the similarity measure used by the method and its implementation.
- The quality of a clustering method is also measured by its ability to discover some or all of the hidden patterns.



Requirements of Clustering in Data Mining

- n Scalability
- n Ability to deal with different types of attributes
- n Discovery of clusters with arbitrary shape
- n Minimal requirements for domain knowledge to determine input parameters
- n Able to deal with noise and outliers
- n Insensitive to order of input records
- n High dimensionality
- n Incorporation of user-specified constraints
- n Interpretability and usability



Chapter 8. Cluster Analysis

- n What is Cluster Analysis?
- n **Types of Data in Cluster Analysis**
- n A Categorization of Major Clustering Methods
- n Partitioning Methods
- n Hierarchical Methods
- n Density-Based Methods
- n Grid-Based Methods
- n Model-Based Clustering Methods
- n Outlier Analysis
- n Summary

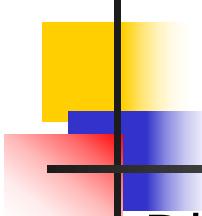
Data Structures

- n Data matrix
 - n (two modes)

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}$$

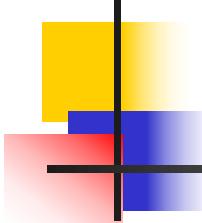
- n Dissimilarity matrix
 - n (one mode)

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$



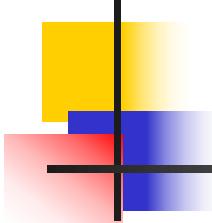
Measure the Quality of Clustering

- n Dissimilarity/Similarity metric: Similarity is expressed in terms of a distance function, which is typically metric:
 $d(i, j)$
- n There is a separate “quality” function that measures the “goodness” of a cluster.
- n The definitions of distance functions are usually very different for interval-scaled, boolean, categorical, ordinal and ratio variables.
- n Weights should be associated with different variables based on applications and data semantics.
- n It is hard to define “similar enough” or “good enough”
 - n the answer is typically highly subjective.



Type of data in clustering analysis

- n Interval-scaled variables:
- n Binary variables:
- n Nominal, ordinal, and ratio variables:
- n Variables of mixed types:



Interval-valued variables

n Standardize data

n Calculate the mean absolute deviation:

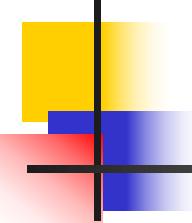
$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

where $m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf})$.

n Calculate the standardized measurement (*z-score*)

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

n Using mean absolute deviation is more robust than using standard deviation



Similarity and Dissimilarity Between Objects

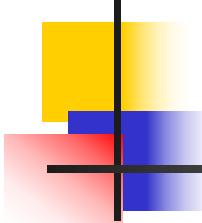
- n Distances are normally used to measure the similarity or dissimilarity between two data objects
- n Some popular ones include: *Minkowski distance*:

$$d(i, j) = \sqrt[q]{(|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{ip} - x_{jp}|^q)}$$

where $i = (x_{i1}, x_{i2}, \dots, x_{ip})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jp})$ are two p -dimensional data objects, and q is a positive integer

- n If $q = 1$, d is Manhattan distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$



Similarity and Dissimilarity Between Objects (Cont.)

n If $q = 2$, d is Euclidean distance:

$$d(i, j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

n Properties

$$\text{n } d(i, j) \geq 0$$

$$\text{n } d(i, i) = 0$$

$$\text{n } d(i, j) = d(j, i)$$

$$\text{n } d(i, j) \leq d(i, k) + d(k, j)$$

n Also one can use weighted distance, parametric Pearson product moment correlation, or other dissimilarity measures.

Binary Variables

n A contingency table for binary data

		Object <i>j</i>		<i>sum</i>
		1	0	
<i>Object i</i>	1	<i>a</i>	<i>b</i>	<i>a+b</i>
	0	<i>c</i>	<i>d</i>	<i>c+d</i>
<i>sum</i>		<i>a+c</i>	<i>b+d</i>	<i>p</i>

- n Simple matching coefficient (invariant, if the binary variable is *symmetric*): $d(i, j) = \frac{b + c}{a + b + c + d}$
- n Jaccard coefficient (noninvariant if the binary variable is *asymmetric*): $d(i, j) = \frac{b + c}{a + b + c}$

Dissimilarity between Binary Variables

n Example

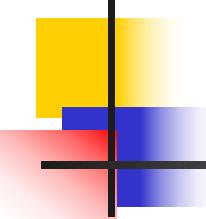
Name	Gender	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	M	Y	N	P	N	N	N
Mary	F	Y	N	P	N	P	N
Jim	M	Y	P	N	N	N	N

- n gender is a symmetric attribute
- n the remaining attributes are asymmetric binary
- n let the values Y and P be set to 1, and the value N be set to 0

$$d(\text{jack}, \text{mary}) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d(\text{jack}, \text{jim}) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d(\text{jim}, \text{mary}) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$



Nominal Variables

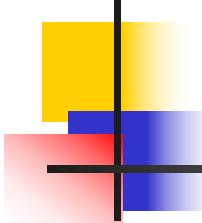
- „ A generalization of the binary variable in that it can take more than 2 states, e.g., red, yellow, blue, green
- „ Method 1: Simple matching
 - „ m : # of matches, p : total # of variables

$$d(i, j) = \frac{p - m}{p}$$

- „ Method 2: use a large number of binary variables
 - „ creating a new binary variable for each of the M nominal states

Ordinal Variables

- An ordinal variable can be discrete or continuous
 - order is important, e.g., rank
 - Can be treated like interval-scaled
 - replacing x_{if} by their rank $r_{if} \in \{1, \dots, M_f\}$
 - map the range of each variable onto $[0, 1]$ by replacing i -th object in the f -th variable by
$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$
 - compute the dissimilarity using methods for interval-scaled variables



Ratio-Scaled Variables

- n Ratio-scaled variable: a positive measurement on a nonlinear scale, approximately at exponential scale, such as Ae^{Bt} or Ae^{-Bt}
- n Methods:
 - n treat them like interval-scaled variables – *not a good choice! (why?)*
 - n apply logarithmic transformation
$$y_{if} = \log(x_{if})$$
 - n treat them as continuous ordinal data treat their rank as interval-scaled.

Variables of Mixed Types

- n A database may contain all the six types of variables
 - n symmetric binary, asymmetric binary, nominal, ordinal, interval and ratio.
- n One may use a weighted formula to combine their effects.

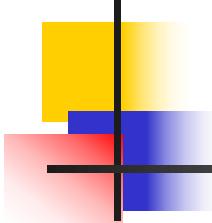
$$d(i, j) = \frac{\sum_{f=1}^p \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p \delta_{ij}^{(f)}}$$

n f is binary or nominal:

$$d_{ij}^{(f)} = 0 \text{ if } x_{if} = x_{jf}, \text{ or } d_{ij}^{(f)} = 1 \text{ o.w.}$$

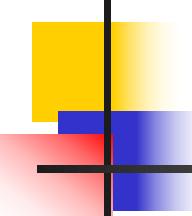
n f is interval-based: use the normalized distance
n f is ordinal or ratio-scaled

n compute ranks r_{if} and $z_{if} = \frac{r_{if} - 1}{M_f - 1}$
n and treat z_{if} as interval-scaled



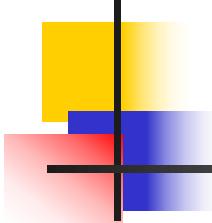
Chapter 8. Cluster Analysis

- n What is Cluster Analysis?
- n Types of Data in Cluster Analysis
- n A Categorization of Major Clustering Methods
- n Partitioning Methods
- n Hierarchical Methods
- n Density-Based Methods
- n Grid-Based Methods
- n Model-Based Clustering Methods
- n Outlier Analysis
- n Summary



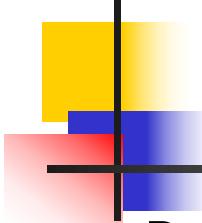
Major Clustering Approaches

- n Partitioning algorithms: Construct various partitions and then evaluate them by some criterion
- n Hierarchy algorithms: Create a hierarchical decomposition of the set of data (or objects) using some criterion
- n Density-based: based on connectivity and density functions
- n Grid-based: based on a multiple-level granularity structure
- n Model-based: A model is hypothesized for each of the clusters and the idea is to find the best fit of that model



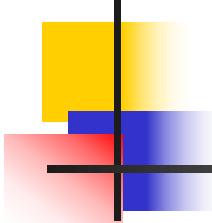
Chapter 8. Cluster Analysis

- n What is Cluster Analysis?
- n Types of Data in Cluster Analysis
- n A Categorization of Major Clustering Methods
- n **Partitioning Methods**
- n Hierarchical Methods
- n Density-Based Methods
- n Grid-Based Methods
- n Model-Based Clustering Methods
- n Outlier Analysis
- n Summary



Partitioning Algorithms: Basic Concept

- n **Partitioning method:** Construct a partition of a database D of n objects into a set of k clusters
- n Given a k , find a partition of k *clusters* that optimizes the chosen partitioning criterion
 - n Global optimal: exhaustively enumerate all partitions
 - n Heuristic methods: *k-means* and *k-medoids* algorithms
 - n *k-means* (MacQueen'67): Each cluster is represented by the center of the cluster
 - n *k-medoids* or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): Each cluster is represented by one of the objects in the cluster

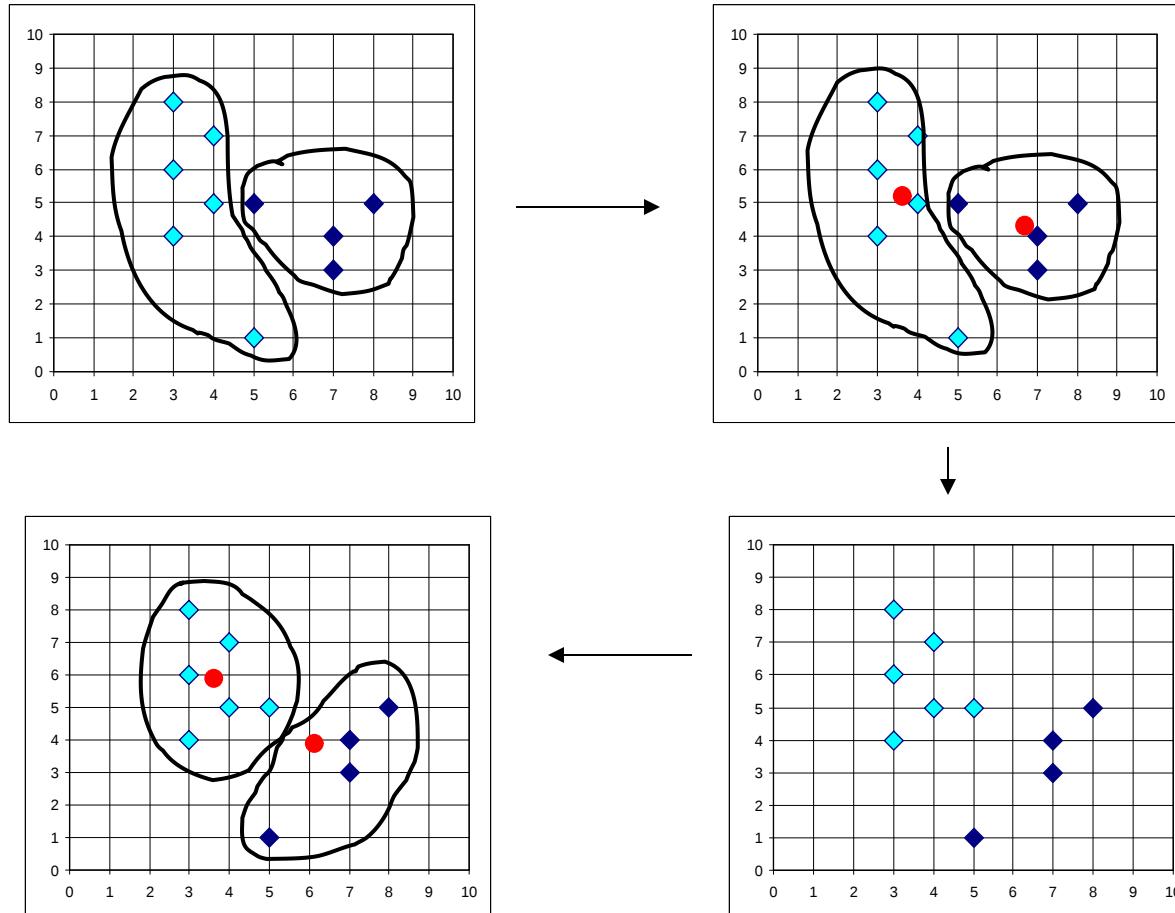


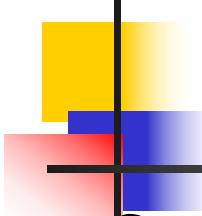
The *K*-Means Clustering Method

- Given k , the *k-means* algorithm is implemented in 4 steps:
 - Partition objects into k nonempty subsets
 - Compute seed points as the centroids of the clusters of the current partition. The centroid is the center (mean point) of the cluster.
 - Assign each object to the cluster with the nearest seed point.
 - Go back to Step 2, stop when no more new assignment.

The K-Means Clustering Method

n Example





Comments on the K-Means Method

n

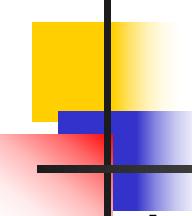
Strength

- n *Relatively efficient:* $O(tkn)$, where n is # objects, k is # clusters, and t is # iterations. Normally, $k, t \ll n$.
- n Often terminates at a *local optimum*. The *global optimum* may be found using techniques such as: *deterministic annealing* and *genetic algorithms*

n

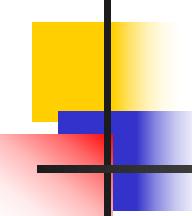
Weakness

- n Applicable only when *mean* is defined, then what about categorical data?
- n Need to specify k , the *number* of clusters, in advance
- n Unable to handle noisy data and *outliers*
- n Not suitable to discover clusters with *non-convex shapes*



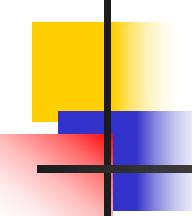
Variations of the *K-Means* Method

- n A few variants of the *k-means* which differ in
 - n Selection of the initial k means
 - n Dissimilarity calculations
 - n Strategies to calculate cluster means
- n Handling categorical data: *k-modes* (Huang'98)
 - n Replacing means of clusters with modes
 - n Using new dissimilarity measures to deal with categorical objects
 - n Using a frequency-based method to update modes of clusters
 - n A mixture of categorical and numerical data: *k-prototype* method



The *K-Medoids* Clustering Method

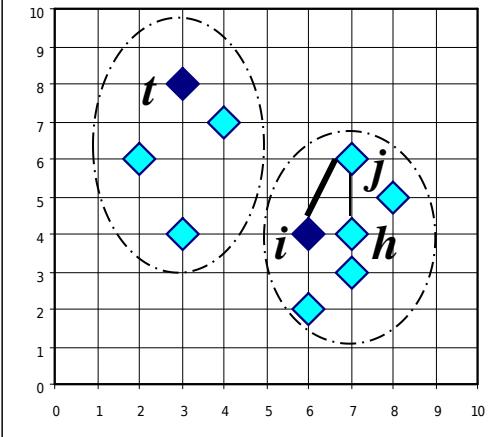
- n Find *representative* objects, called medoids, in clusters
- n *PAM* (Partitioning Around Medoids, 1987)
 - n starts from an initial set of medoids and iteratively replaces one of the medoids by one of the non-medoids if it improves the total distance of the resulting clustering
 - n *PAM* works effectively for small data sets, but does not scale well for large data sets
- n *CLARA* (Kaufmann & Rousseeuw, 1990)
- n *CLARANS* (Ng & Han, 1994): Randomized sampling
- n Focusing + spatial data structure (Ester et al., 1995)



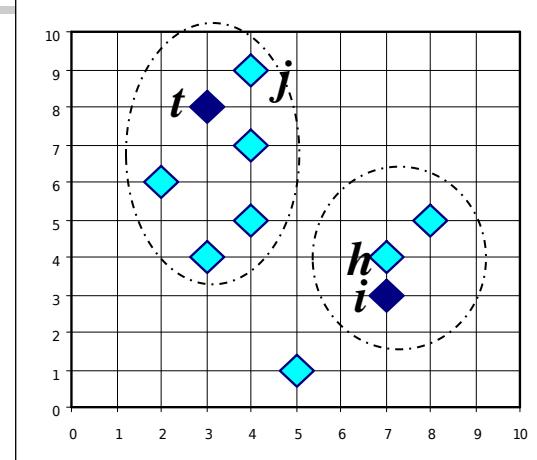
PAM (Partitioning Around Medoids) (1987)

- PAM (Kaufman and Rousseeuw, 1987), built in Splus
- Use real object to represent the cluster
 - Select k representative objects arbitrarily
 - For each pair of non-selected object h and selected object i , calculate the total swapping cost TC_{ih}
 - For each pair of i and h ,
 - If $TC_{ih} < 0$, i is replaced by h
 - Then assign each non-selected object to the most similar representative object
 - repeat steps 2-3 until there is no change

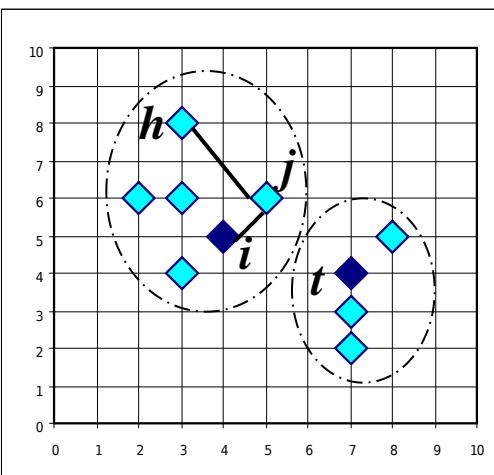
PAM Clustering: Total swapping cost $TC_{ih} = \sum_j C_{jih}$



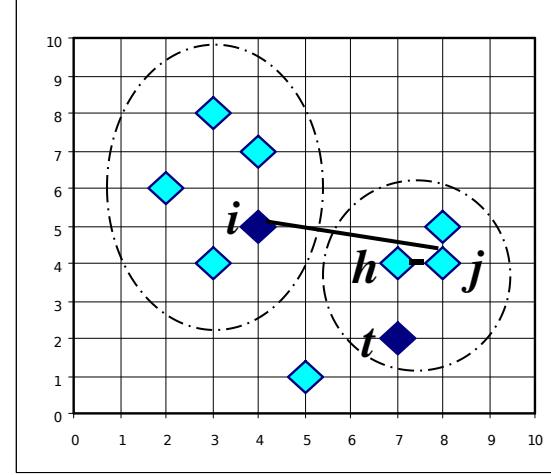
$$C_{jih} = d(j, h) - d(j, i)$$



$$C_{jih} = 0$$



$$C_{jih} = d(j, t) - d(j, i)$$



$$C_{jih} = d(j, h) - d(j, t)$$

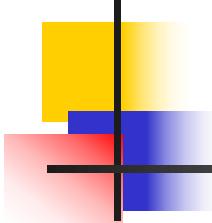
CLARA (Clustering Large Applications) (1990)

- n CLARA (Kaufmann and Rousseeuw in 1990)
 - n Built in statistical analysis packages, such as S+
 - n It draws *multiple samples* of the data set, applies *PAM* on each sample, and gives the best clustering as the output
 - n Strength: deals with larger data sets than *PAM*
 - n Weakness:
 - n Efficiency depends on the sample size
 - n A good clustering based on samples will not necessarily represent a good clustering of the whole data set if the sample is biased



CLARANS (“Randomized” CLARA) (1994)

- n CLARANS (A Clustering Algorithm based on Randomized Search) (Ng and Han'94)
- n CLARANS draws sample of neighbors dynamically
- n The clustering process can be presented as searching a graph where every node is a potential solution, that is, a set of k medoids
- n If the local optimum is found, CLARANS starts with new randomly selected node in search for a new local optimum
- n It is more efficient and scalable than both *PAM* and *CLARA*

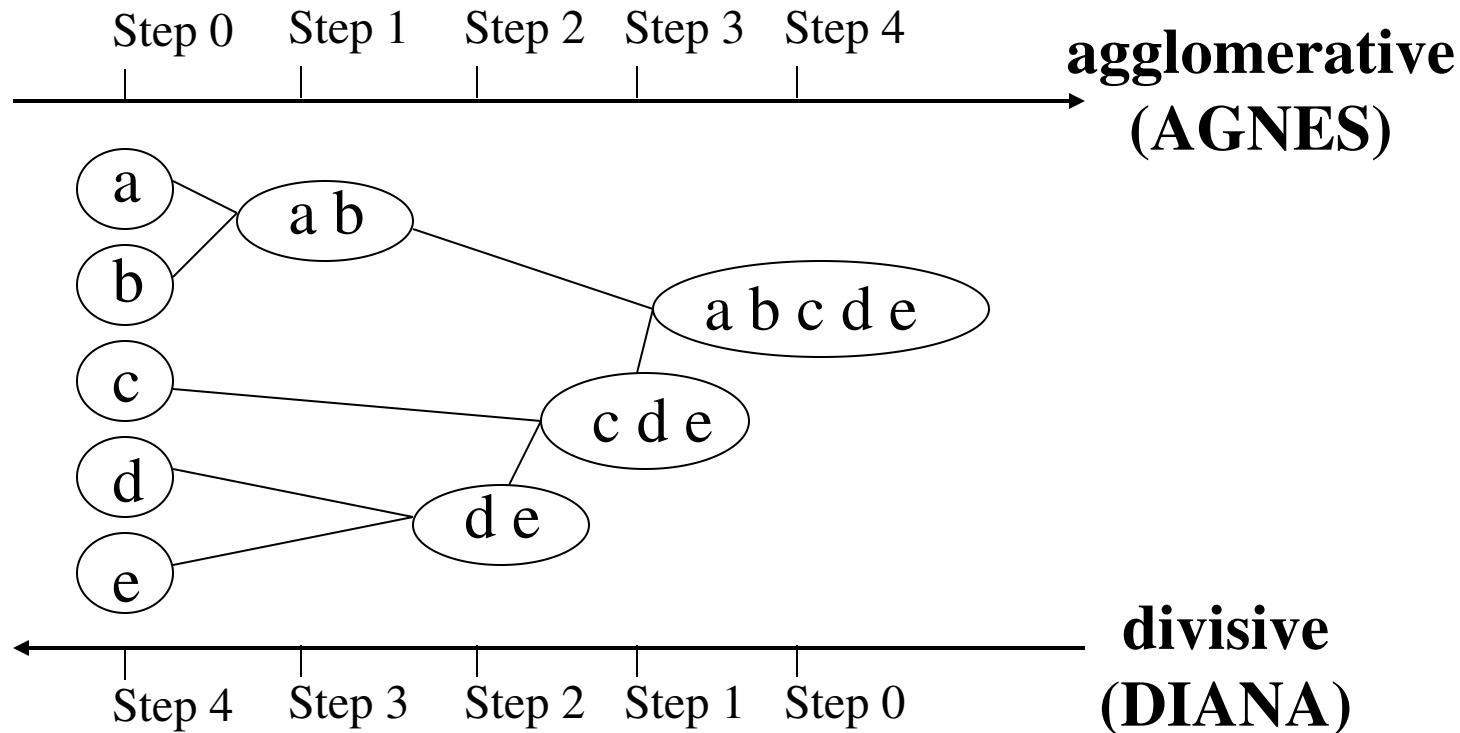


Chapter 8. Cluster Analysis

- n What is Cluster Analysis?
- n Types of Data in Cluster Analysis
- n A Categorization of Major Clustering Methods
- n Partitioning Methods
- n **Hierarchical Methods**
- n Density-Based Methods
- n Grid-Based Methods
- n Model-Based Clustering Methods
- n Outlier Analysis
- n Summary

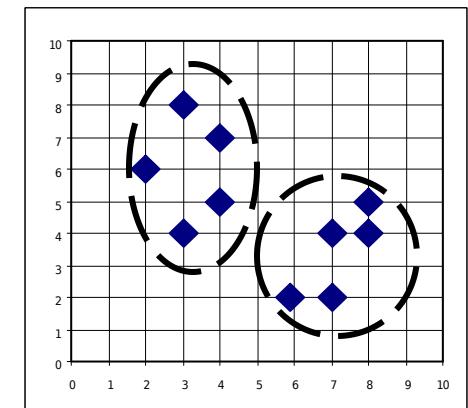
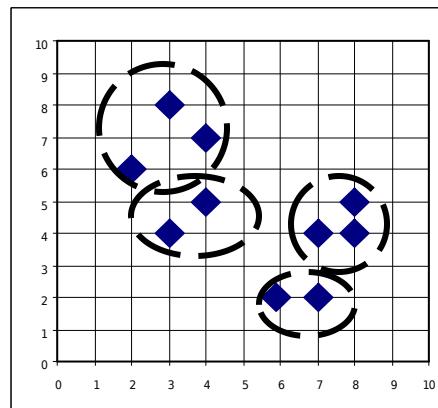
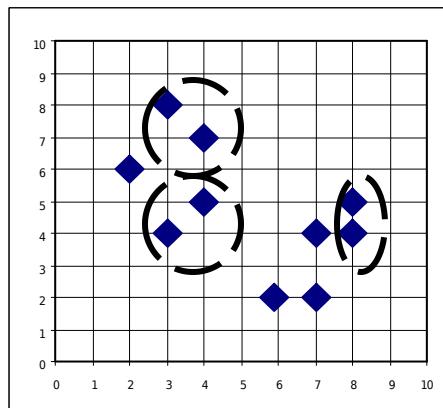
Hierarchical Clustering

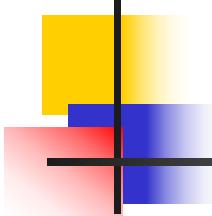
- Use distance matrix as clustering criteria. This method does not require the number of clusters k as an input, but needs a termination condition



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

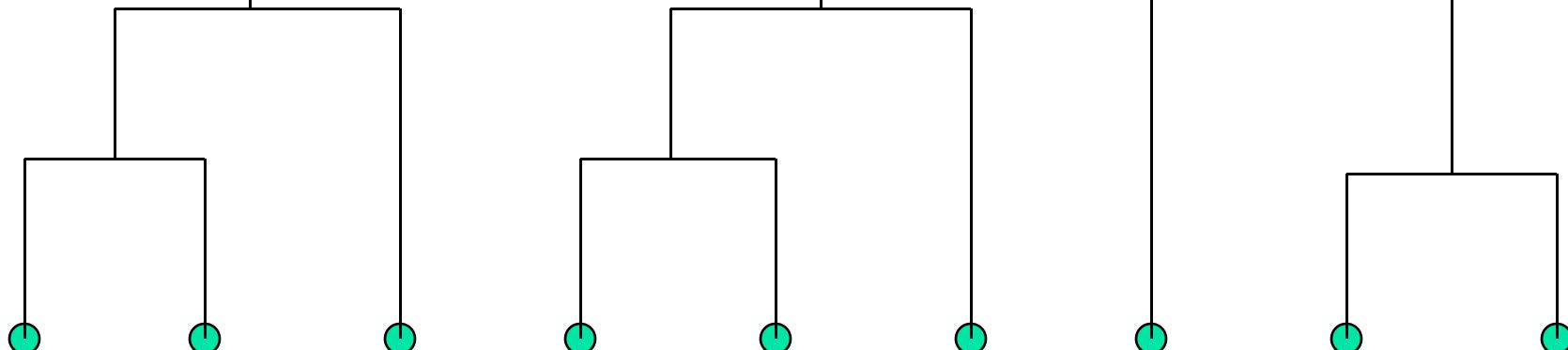




A *Dendrogram* Shows How the Clusters are Merged Hierarchically

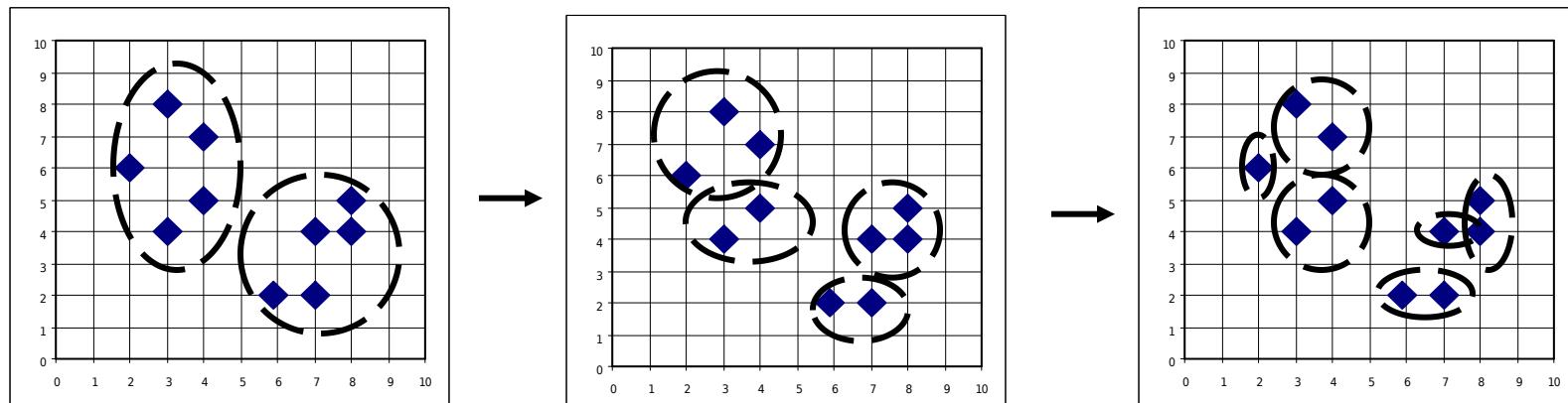
Decompose data objects into a several levels of nested partitioning (tree of clusters), called a dendrogram.

A clustering of the data objects is obtained by cutting the dendrogram at the desired level, then each connected component forms a cluster.



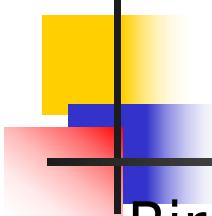
DIANA (Divisive Analysis)

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



More on Hierarchical Clustering Methods

- n Major weakness of agglomerative clustering methods
 - n do not scale well: time complexity of at least $O(n^2)$, where n is the number of total objects
 - n can never undo what was done previously
- n Integration of hierarchical with distance-based clustering
 - n BIRCH (1996): uses CF-tree and incrementally adjusts the quality of sub-clusters
 - n CURE (1998): selects well-scattered points from the cluster and then shrinks them towards the center of the cluster by a specified fraction
 - n CHAMELEON (1999): hierarchical clustering using dynamic modeling



BIRCH (1996)

- Birch: Balanced Iterative Reducing and Clustering using Hierarchies, by 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
- *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

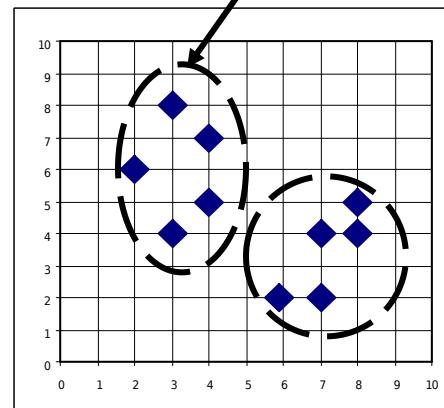
Clustering Feature: $CF = (\overrightarrow{N}, \overrightarrow{LS}, SS)$

N: Number of data points

LS: $\overrightarrow{N}_{i=1} = \overrightarrow{X}_i$

SS: $\overrightarrow{N}_{i=1} = \overrightarrow{X}_i^2$

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



(3,4)

(2,6)

(4,5)

(4,7)

(3,8)

CF Tree

Root

$B = 7$

$L = 6$

	CF_1	CF_2	CF_3	CF_6	
	$child_1$	$child_2$	$child_3$			$child_6$

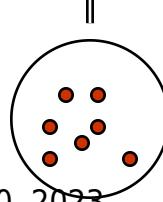
Non-leaf node

CF_1	CF_2	CF_3	CF_5	
$child_1$	$child_2$	$child_3$			$child_5$

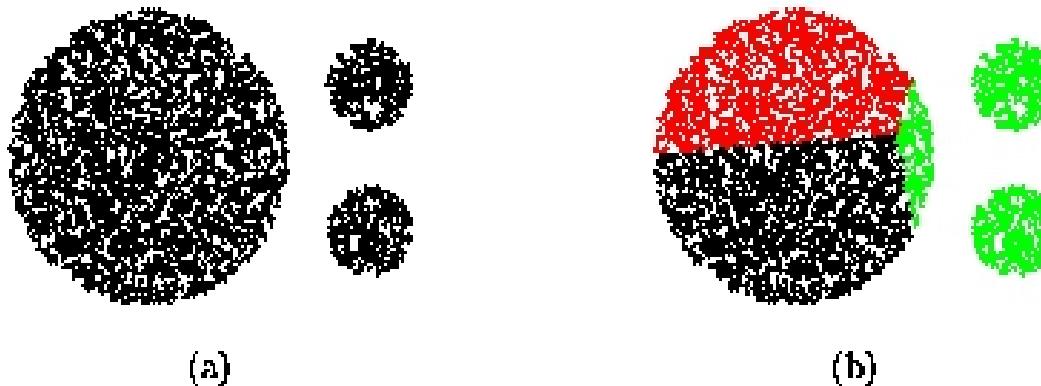
Leaf node

prev	CF_1	CF_2	CF_6	next
------	--------	--------	-------	--------	------

prev	CF_1	CF_2	CF_4	next
------	--------	--------	-------	--------	------

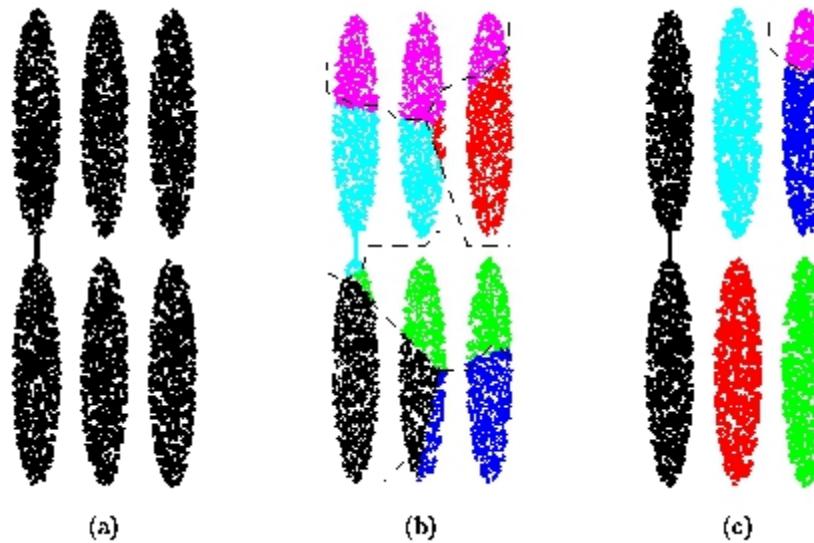


CURE (Clustering Using REpresentatives)

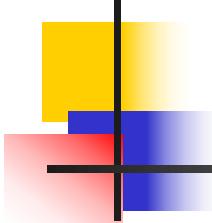


- n CURE: proposed by Guha, Rastogi & Shim, 1998
 - n Stops the creation of a cluster hierarchy if a level consists of k clusters
 - n Uses multiple representative points to evaluate the distance between clusters, adjusts well to arbitrary shaped clusters and avoids single-link effect

Drawbacks of Distance-Based Method



- n Drawbacks of square-error based clustering method
 - n Consider only one point as representative of a cluster
 - n Good only for convex shaped, similar size and density, and if k can be reasonably estimated



Cure: The Algorithm

- Draw random sample s .
- Partition sample to p partitions with size s/p
- Partially cluster partitions into s/pq clusters
- Eliminate outliers
 - By random sampling
 - If a cluster grows too slow, eliminate it.
 - Cluster partial clusters.
 - Label data in disk

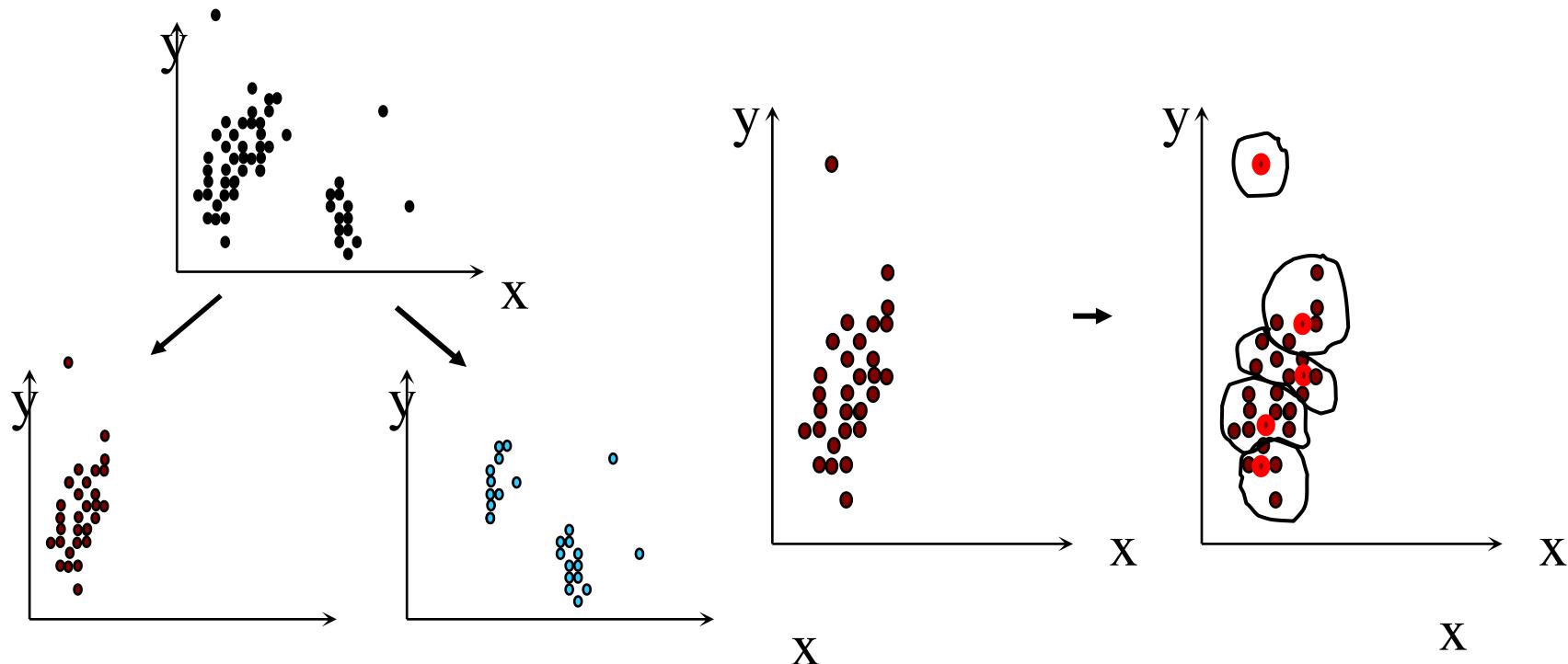
Data Partitioning and Clustering

$$\textcolor{red}{n} \quad s = 50$$

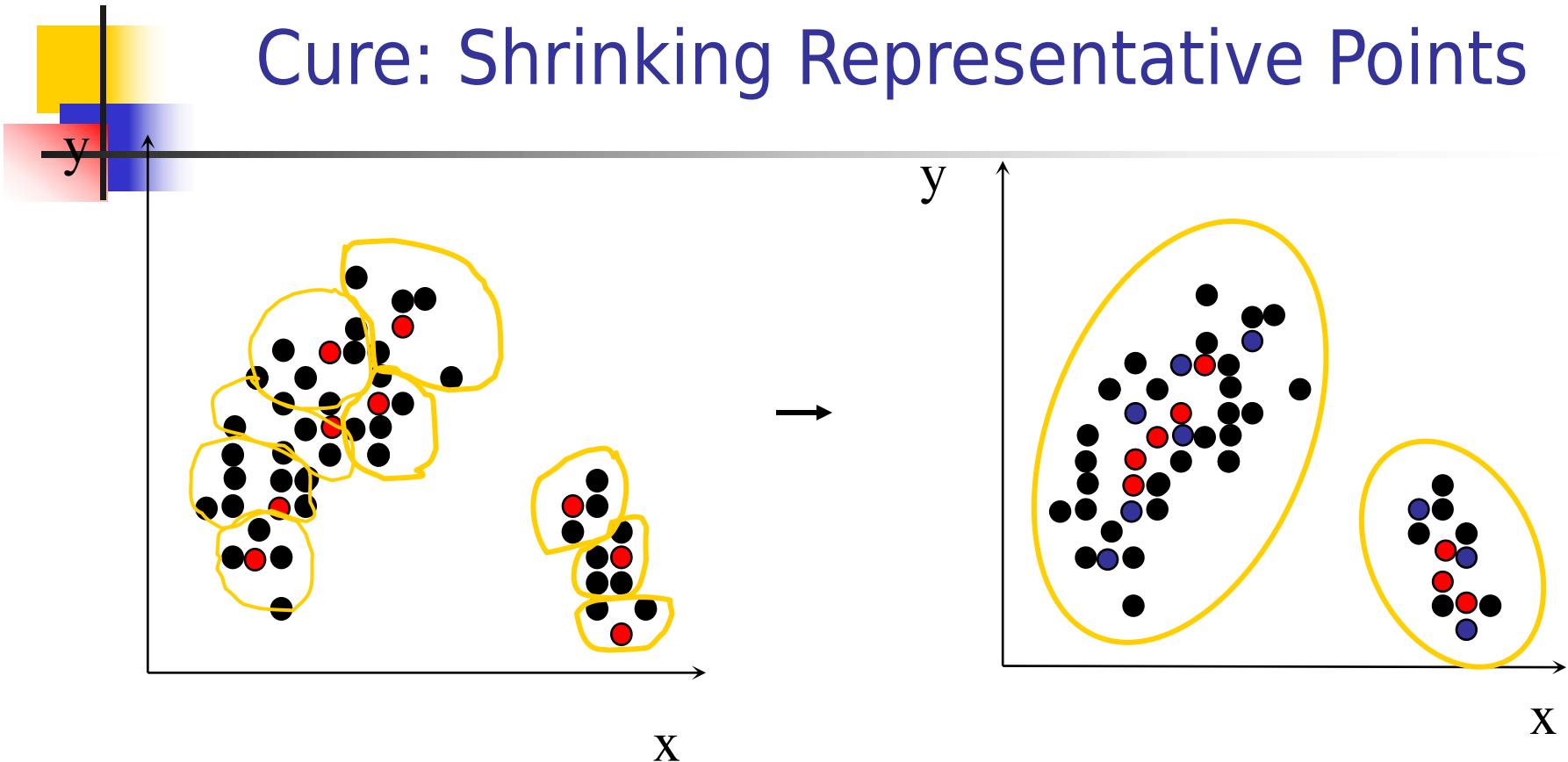
$$\textcolor{red}{n} \quad p = 2$$

$$\textcolor{red}{n} \quad s/p = 25$$

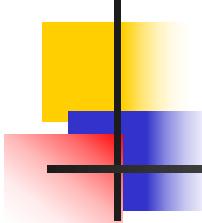
$$\textcolor{red}{n} \quad s/pq = 5$$



Cure: Shrinking Representative Points



- n Shrink the multiple representative points towards the gravity center by a fraction of .
- n Multiple representatives capture the shape of the cluster



Clustering Categorical Data: ROCK

- n ROCK: Robust Clustering using links,
by S. Guha, R. Rastogi, K. Shim (ICDE'99).
 - n Use links to measure similarity/proximity
 - n Not distance based
 - n Computational complexity: $O(n^2 + nm_m m_a + n^2 \log n)$
- n Basic ideas:
 - n Similarity function and neighbors:
$$Sim(T_1, T_2) = \frac{|T_1 \cap T_2|}{|T_1 \cup T_2|}$$

Let $T_1 = \{1,2,3\}$, $T_2 = \{3,4,5\}$

$$Sim(T_1, T_2) = \frac{|\{3\}|}{|\{1, 2, 3, 4, 5\}|} = \frac{1}{5} = 0.2$$

Rock: Algorithm

- n Links: The number of common neighbours for the two points.

$\{1,2,3\}, \{1,2,4\}, \{1,2,5\}, \{1,3,4\}, \{1,3,5\}$

$\{1,4,5\}, \{2,3,4\}, \{2,3,5\}, \{2,4,5\}, \{3,4,5\}$

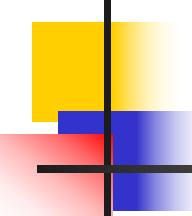
$$\{1,2,3\} \xleftarrow{3} \{1,2,4\}$$

- n Algorithm

- n Draw random sample

- n Cluster with links

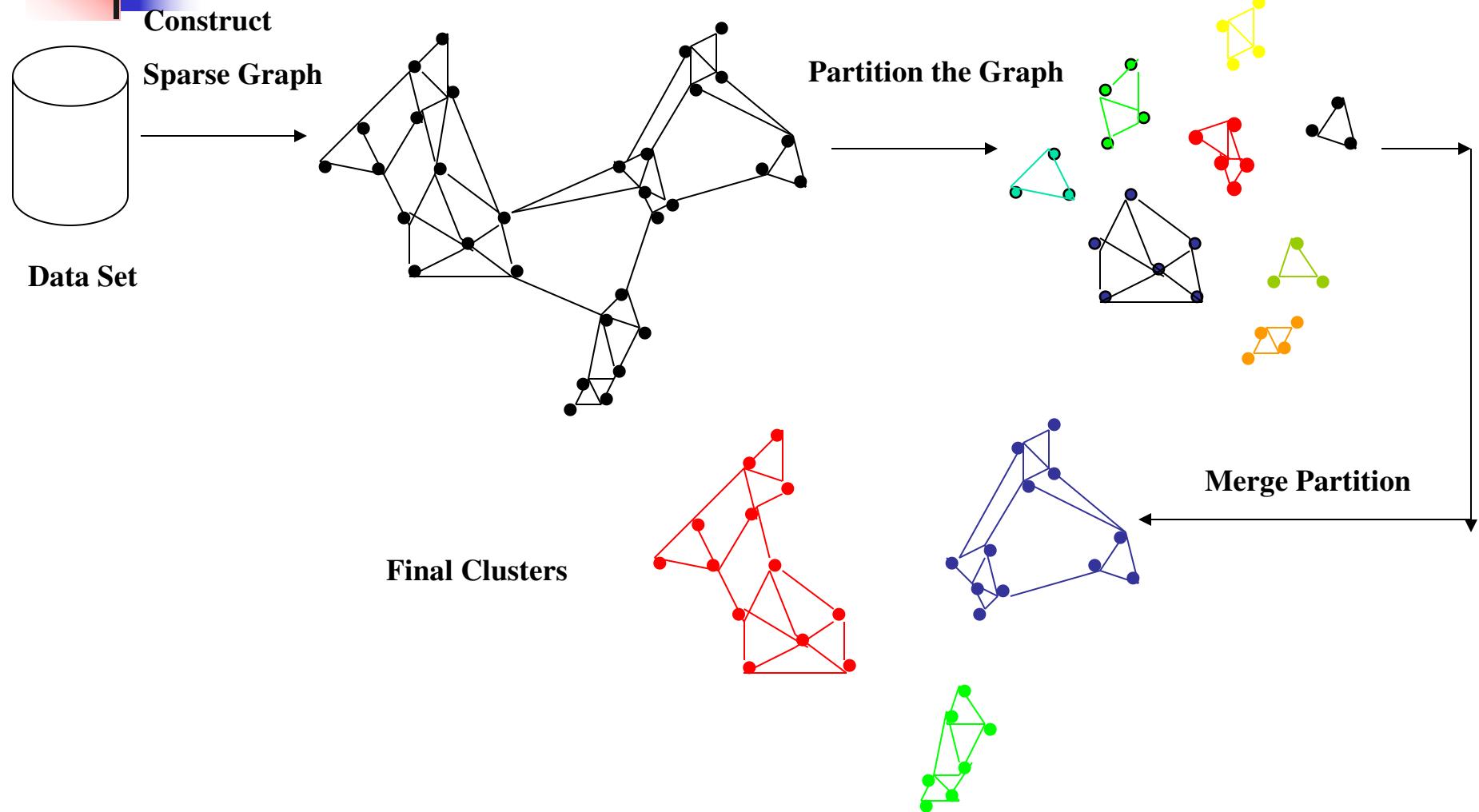
- n Label data in disk

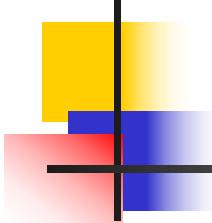


CHAMELEON

- 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
- 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

Overall Framework of CHAMELEON





Chapter 8. Cluster Analysis

- n What is Cluster Analysis?
- n Types of Data in Cluster Analysis
- n A Categorization of Major Clustering Methods
- n Partitioning Methods
- n Hierarchical Methods
- n **Density-Based Methods**
- n Grid-Based Methods
- n Model-Based Clustering Methods
- n Outlier Analysis
- n Summary

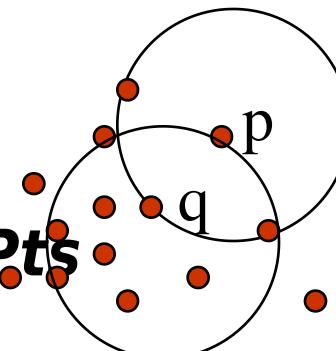
Density-Based Clustering Methods

- n Clustering based on density (local cluster criterion), such as density-connected points
- n Major features:
 - n Discover clusters of arbitrary shape
 - n Handle noise
 - n One scan
 - n Need density parameters as termination condition
- n Several interesting studies:
 - n DBSCAN: Ester, et al. (KDD'96)
 - n OPTICS: Ankerst, et al (SIGMOD'99).
 - n DENCLUE: Hinneburg & D. Keim (KDD'98)
 - n CLIQUE: Agrawal, et al. (SIGMOD'98)

Density-Based Clustering: Background

n Two parameters:

- n **Eps**: Maximum radius of the neighbourhood
- n **MinPts**: Minimum number of points in an Eps-neighbourhood of that point
- n $N_{Eps}(p)$: $\{q \text{ belongs to } D \mid dist(p,q) \leq Eps\}$
- n Directly density-reachable: A point p is directly density-reachable from a point q wrt. **Eps**, **MinPts** if
 - n 1) p belongs to $N_{Eps}(q)$
 - n 2) core point condition:
$$|N_{Eps}(q)| \geq MinPts$$



MinPts = 5

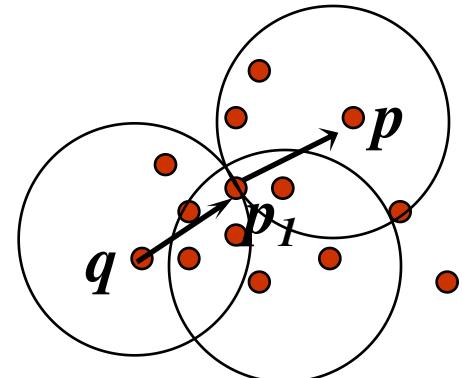
Eps = 1 cm

Density-Based Clustering: Background

(1)

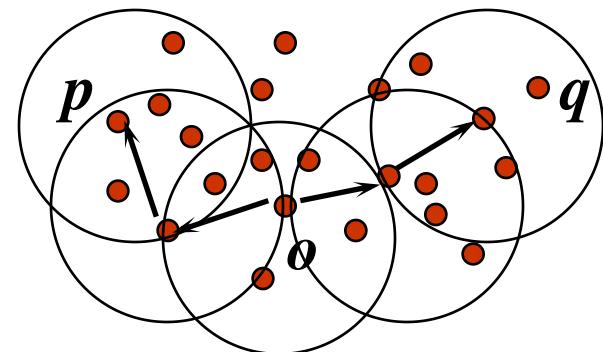
n Density-reachable:

- n A point p is density-reachable from a point q wrt. $Eps, MinPts$ 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



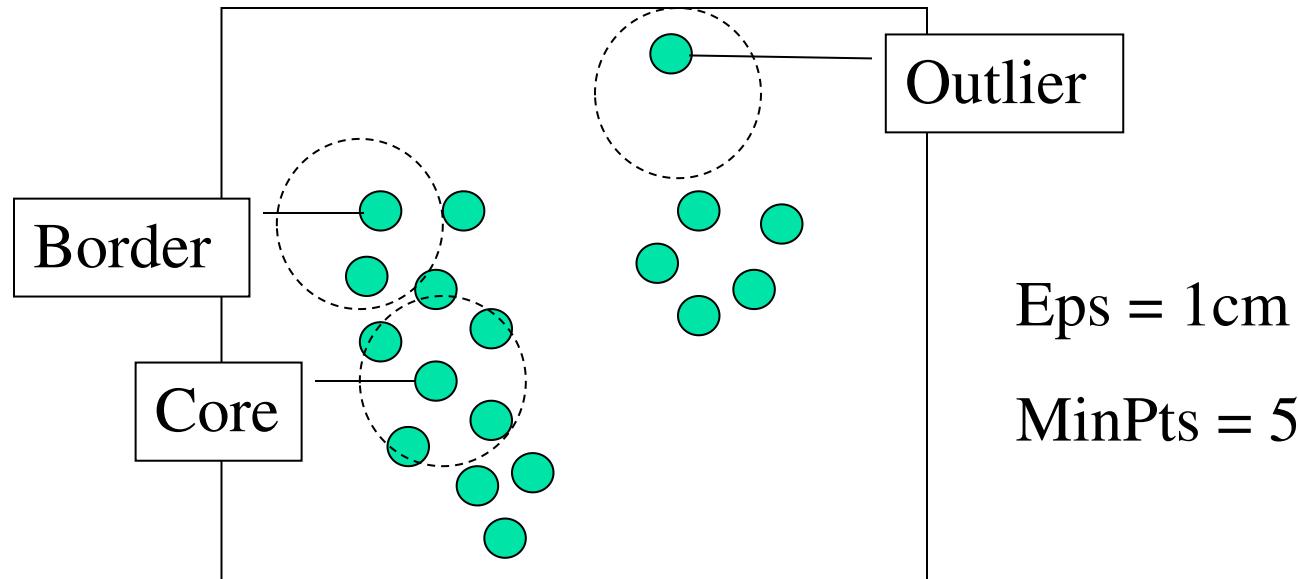
n Density-connected

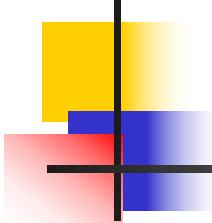
- n A point p is density-connected to a point q wrt. $Eps, MinPts$ if there is a point o such that both, p and q are density-reachable from o wrt. Eps and $MinPts$.



DBSCAN: Density Based Spatial Clustering of Applications with Noise

- n Relies on a *density-based* notion of cluster: A *cluster* is defined as a maximal set of density-connected points
- n Discovers clusters of arbitrary shape in spatial databases with noise





DBSCAN: The Algorithm

- Arbitrary select a point p
- Retrieve all points density-reachable from p wrt 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.



OPTICS: A Cluster-Ordering Method (1999)

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

OPTICS: Some Extension from DBSCAN

n Index-based:

n k = number of dimensions

n N = 20

n p = 75%

n M = N(1-p) = 5

n Complexity: $O(kN^2)$

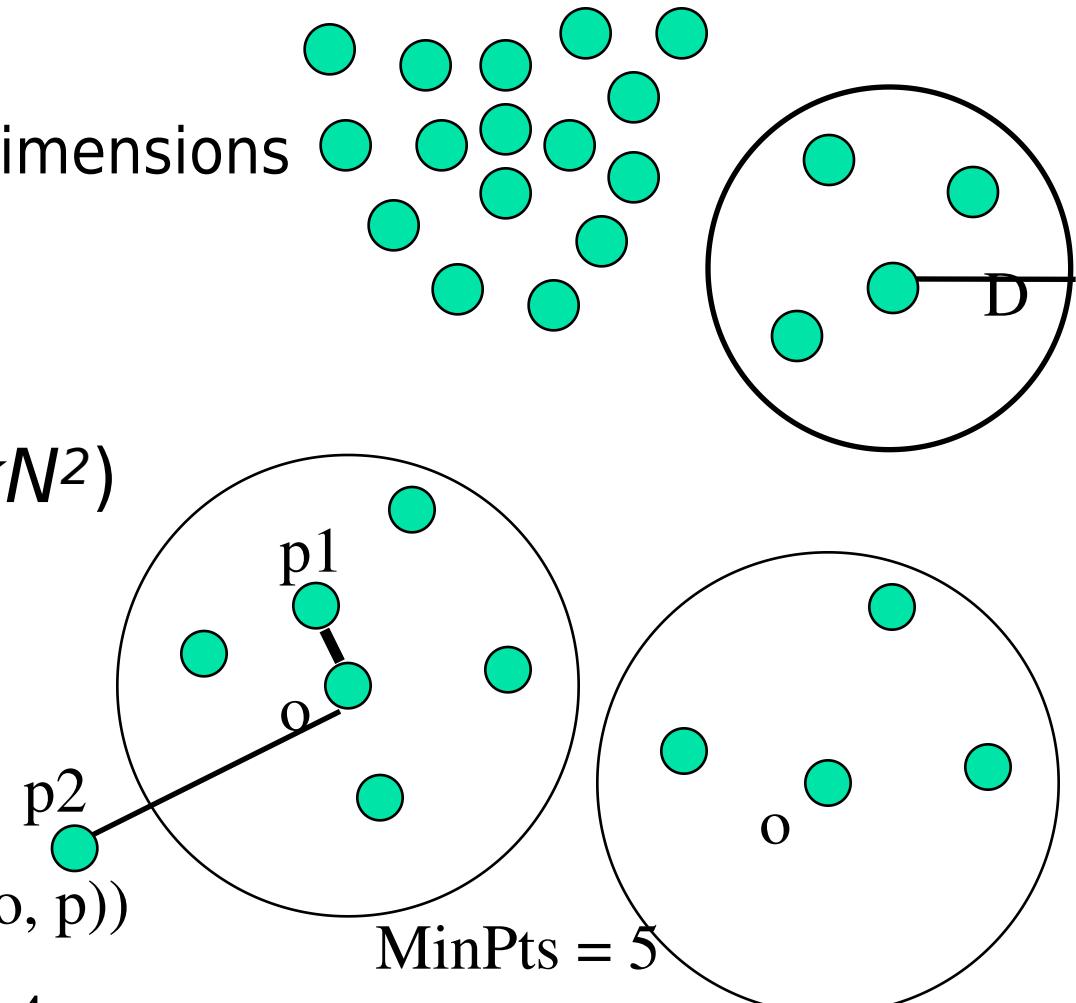
n Core Distance

n Reachability Distance

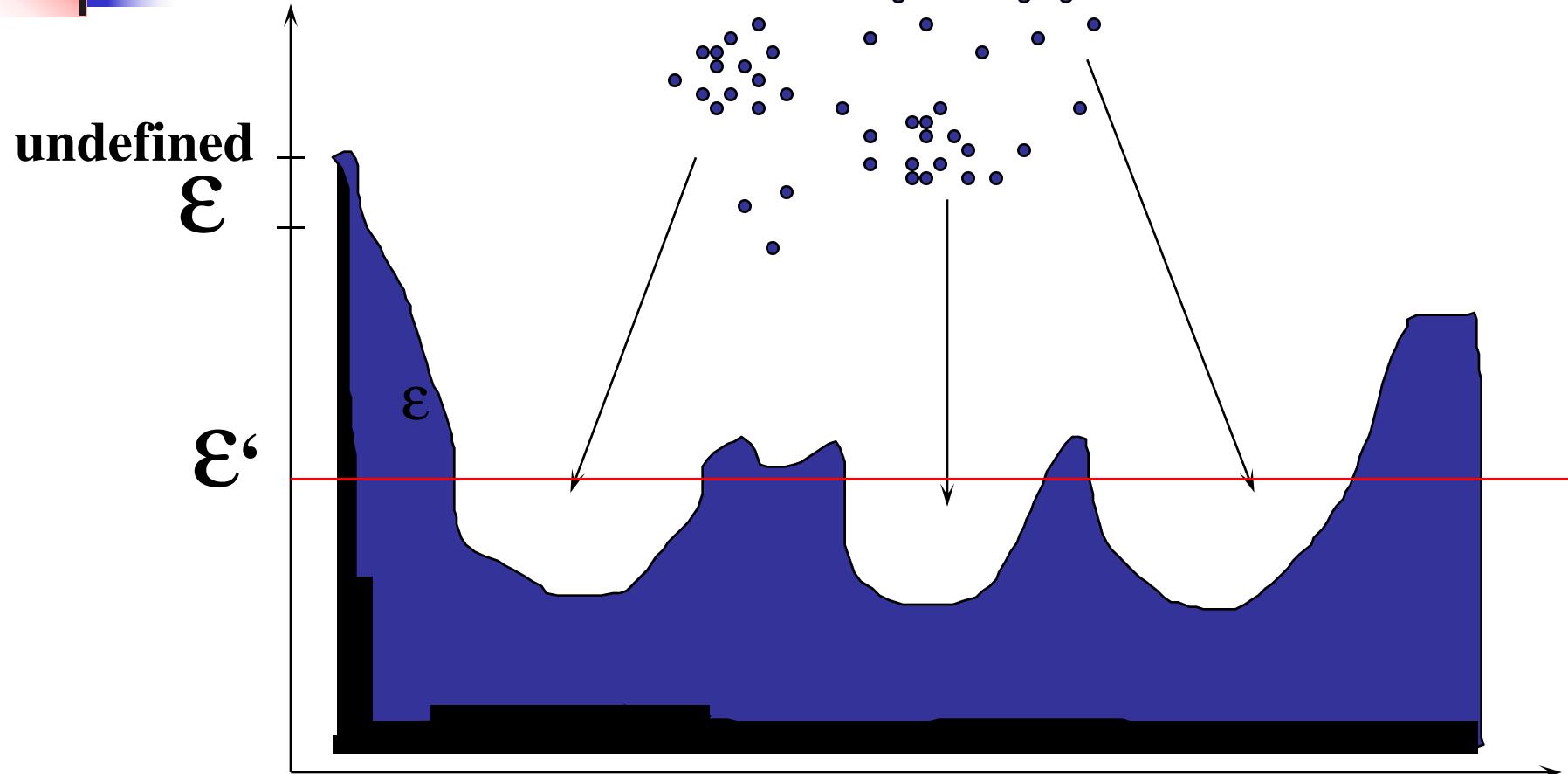
Max (core-distance (o), d (o, p))

$r(p_1, o) = 2.8\text{cm}$. $r(p_2, o) = 4\text{cm}$

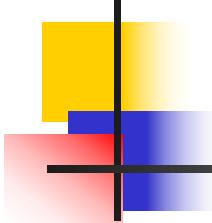
$e = 3 \text{ cm}$



Reachability-distance

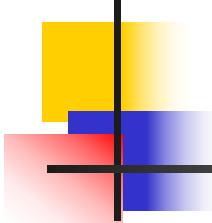


**Cluster-order
of the objects**



DENCLUE: using density functions

- n DENsity-based CLUstEring by Hinneburg & Keim (KDD'98)
- n Major features
 - n Solid mathematical foundation
 - n Good for data sets with large amounts of noise
 - n Allows a compact mathematical description of arbitrarily shaped clusters in high-dimensional data sets
 - n Significant faster than existing algorithm (faster than DBSCAN by a factor of up to 45)
 - n But needs a large number of parameters



Denclue: Technical Essence

- 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.

Gradient: The steepness of a slope

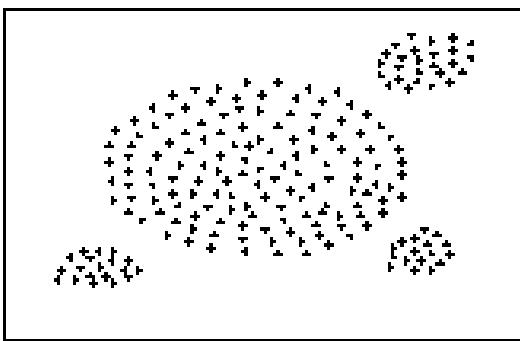
n Example

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

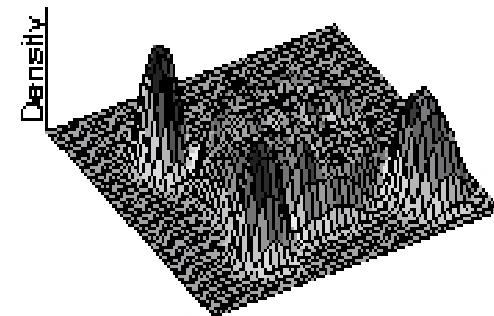
$$f_{Gaussian}^D(x) = \sum_{i=1}^N e^{-\frac{d(x, x_i)^2}{2\sigma^2}}$$

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

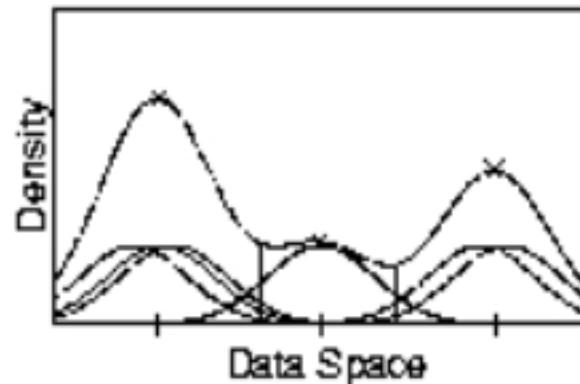
Density Attractor



(a) Data Set



(c) Gaussian



Center-Defined and Arbitrary

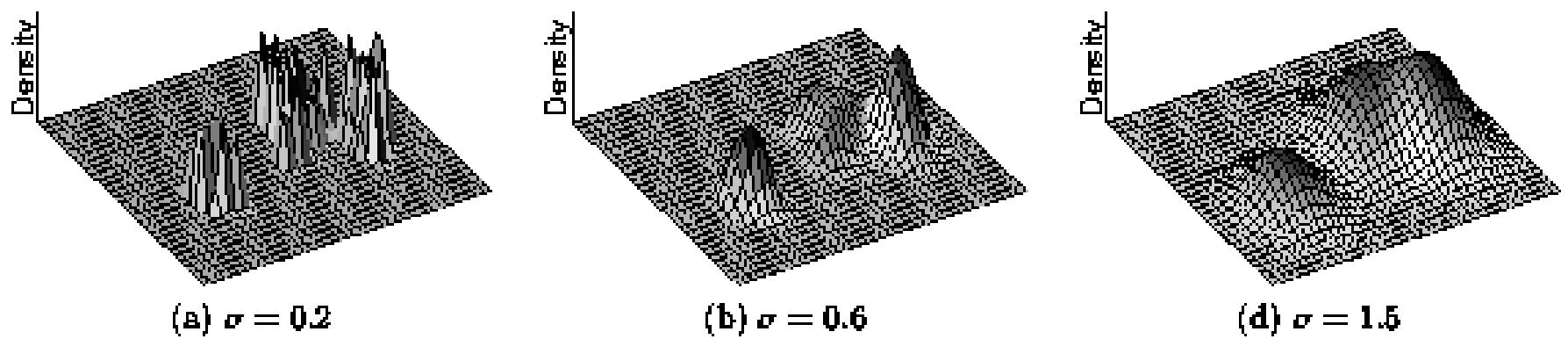


Figure 3: Example of Center-Defined Clusters for different σ

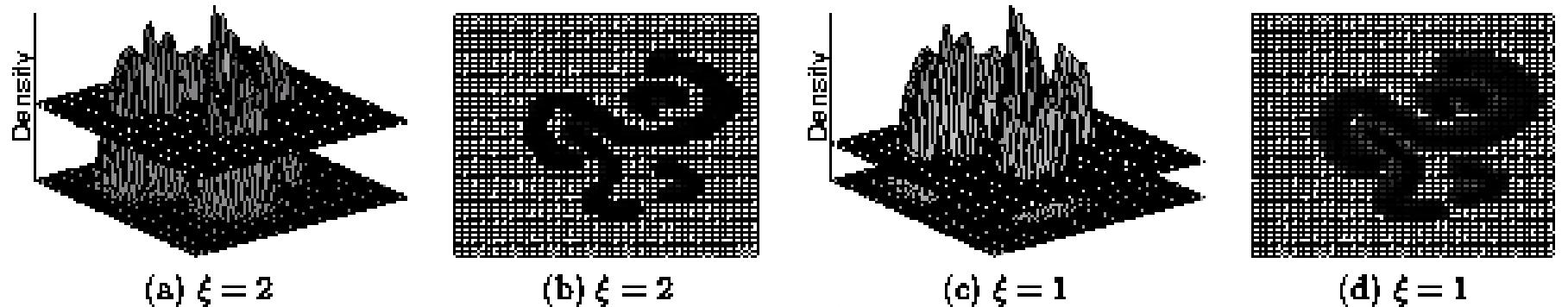
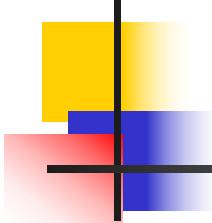
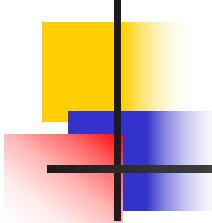


Figure 4: Example of Arbitrary-Shape Clusters for different ξ



Chapter 8. Cluster Analysis

- n What is Cluster Analysis?
- n Types of Data in Cluster Analysis
- n A Categorization of Major Clustering Methods
- n Partitioning Methods
- n Hierarchical Methods
- n Density-Based Methods
- n **Grid-Based Methods**
- n Model-Based Clustering Methods
- n Outlier Analysis
- n Summary

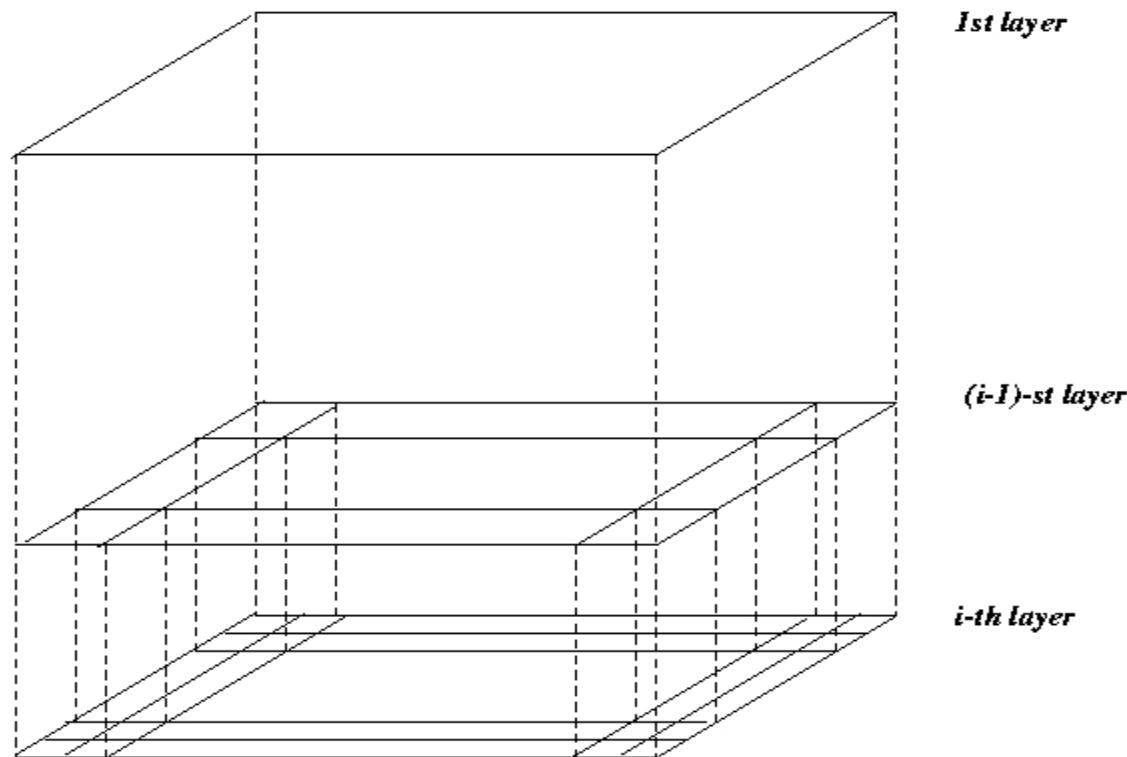


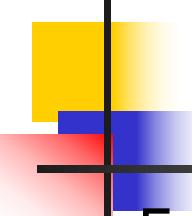
Grid-Based Clustering Method

- „ 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)

STING: A Statistical Information Grid Approach

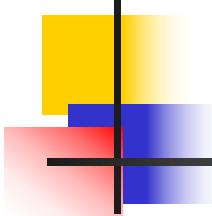
- Wang, Yang and Muntz (VLDB'97)
- The spatial area area is divided into rectangular cells
- There are several levels of cells corresponding to different levels of resolution





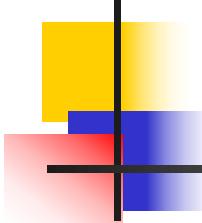
STING: A Statistical Information Grid Approach (2)

- 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: A Statistical Information Grid Approach (3)

- 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

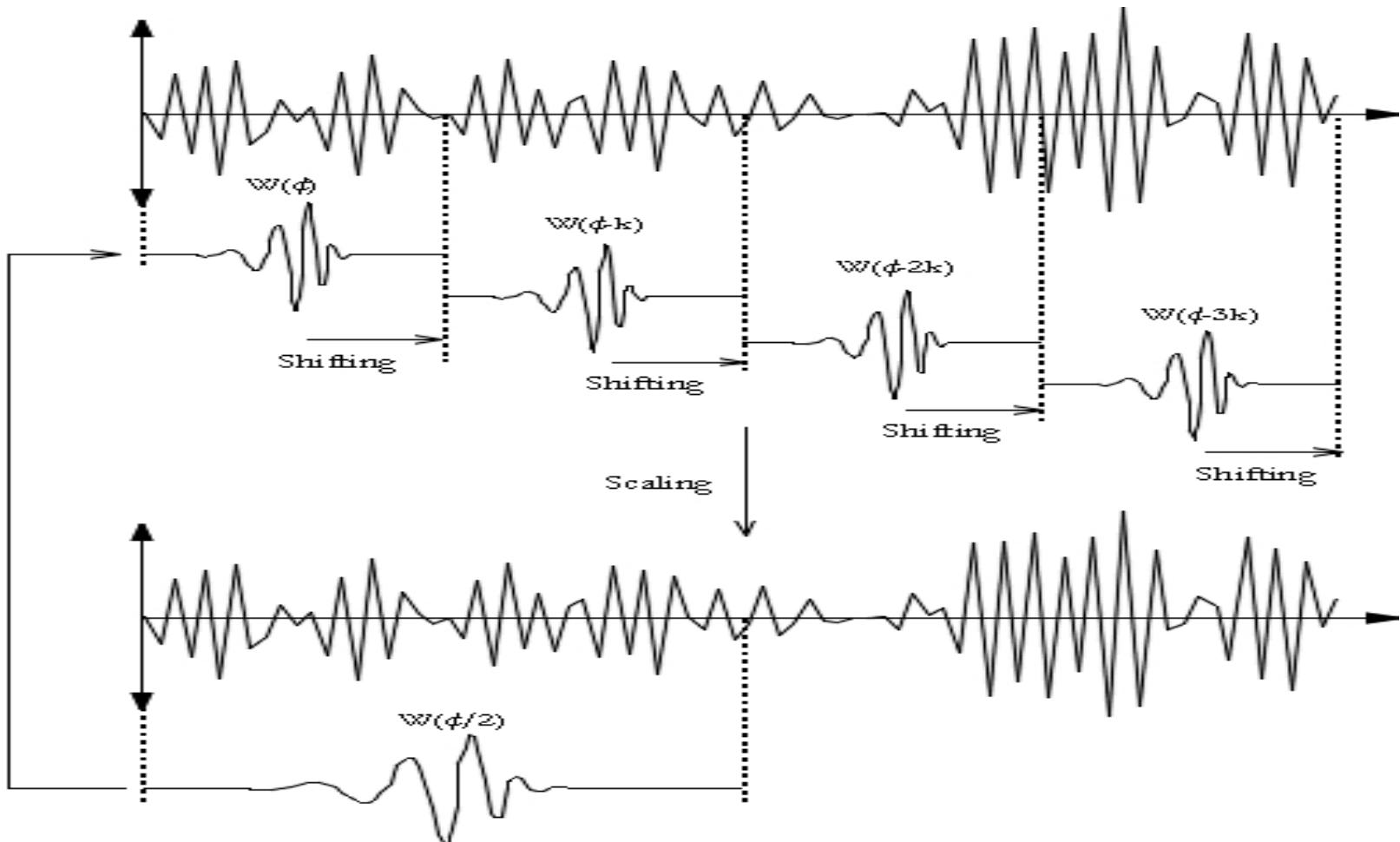


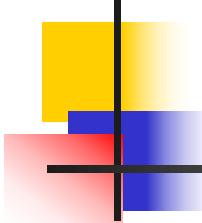
WaveCluster (1998)

- „ Sheikholeslami, Chatterjee, and Zhang (VLDB'98)
- „ A multi-resolution clustering approach which applies wavelet transform to the feature space
 - „ A wavelet transform is a signal processing technique that decomposes a signal into different frequency sub-band.
- „ Both grid-based and density-based
- „ Input parameters:
 - „ # of grid cells for each dimension
 - „ the wavelet, and the # of applications of wavelet transform.

What is Wavelet (1)?

Repeat Shifting Operation

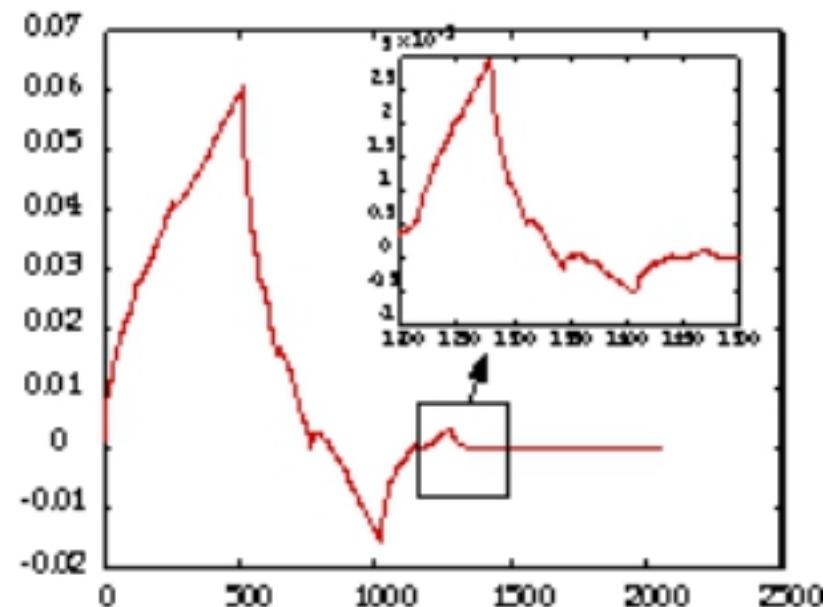
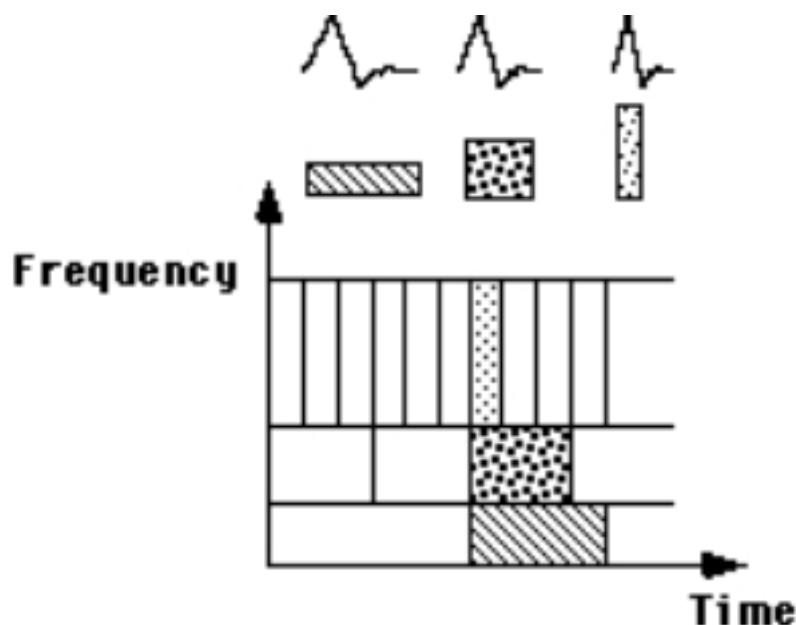




WaveCluster (1998)

- n How to apply wavelet transform to find clusters
 - n Summaries the data by imposing a multidimensional grid structure onto data space
 - n These multidimensional spatial data objects are represented in a n-dimensional feature space
 - n Apply wavelet transform on feature space to find the dense regions in the feature space
 - n Apply wavelet transform multiple times which result in clusters at different scales from fine to coarse

What Is Wavelet (2)?



Quantization

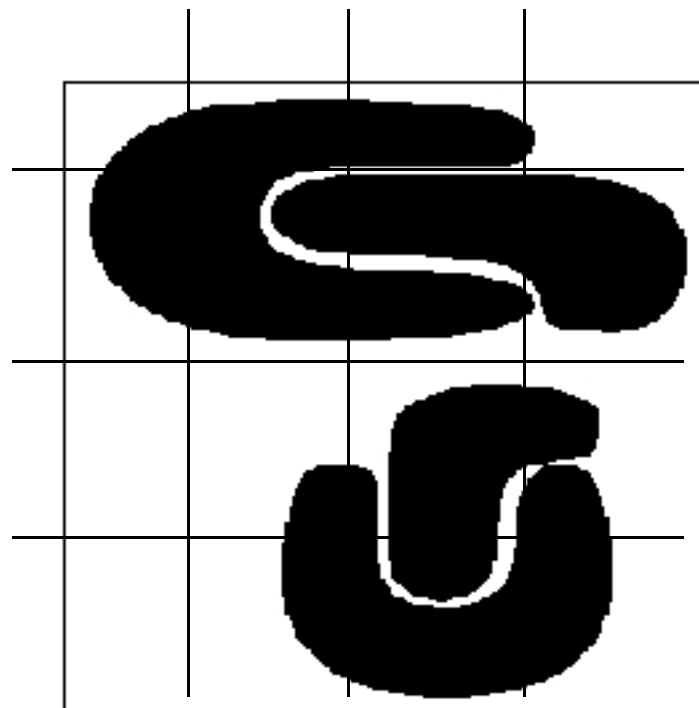


Figure 1: A sample 2-dimensional feature space.

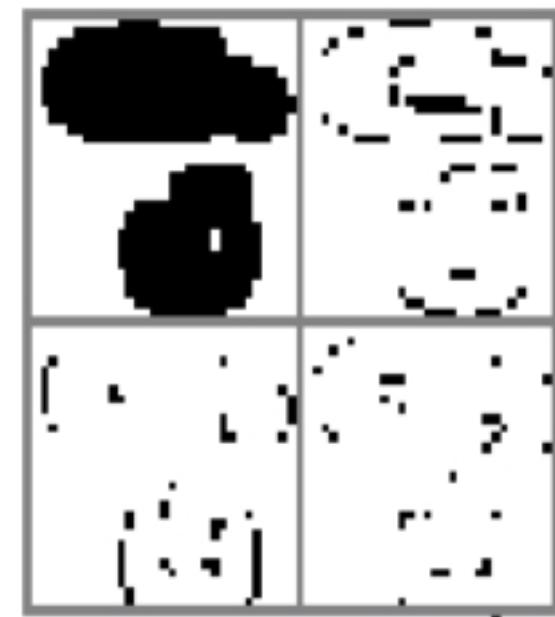
Transformation



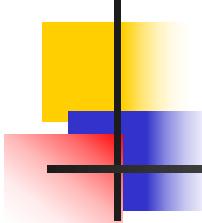
a)



b)

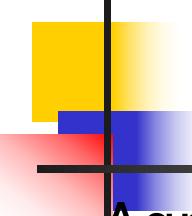


c)



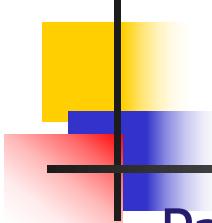
WaveCluster (1998)

- n Why is wavelet transformation useful for clustering
 - n Unsupervised clustering
 - It uses hat-shape filters to emphasize region where points cluster, but simultaneously to suppress weaker information in their boundary
 - n Effective removal of outliers
 - n Multi-resolution
 - n Cost efficiency
- n Major features:
 - n Complexity $O(N)$
 - n Detect arbitrary shaped clusters at different scales
 - n Not sensitive to noise, not sensitive to input order
 - n Only applicable to low dimensional data



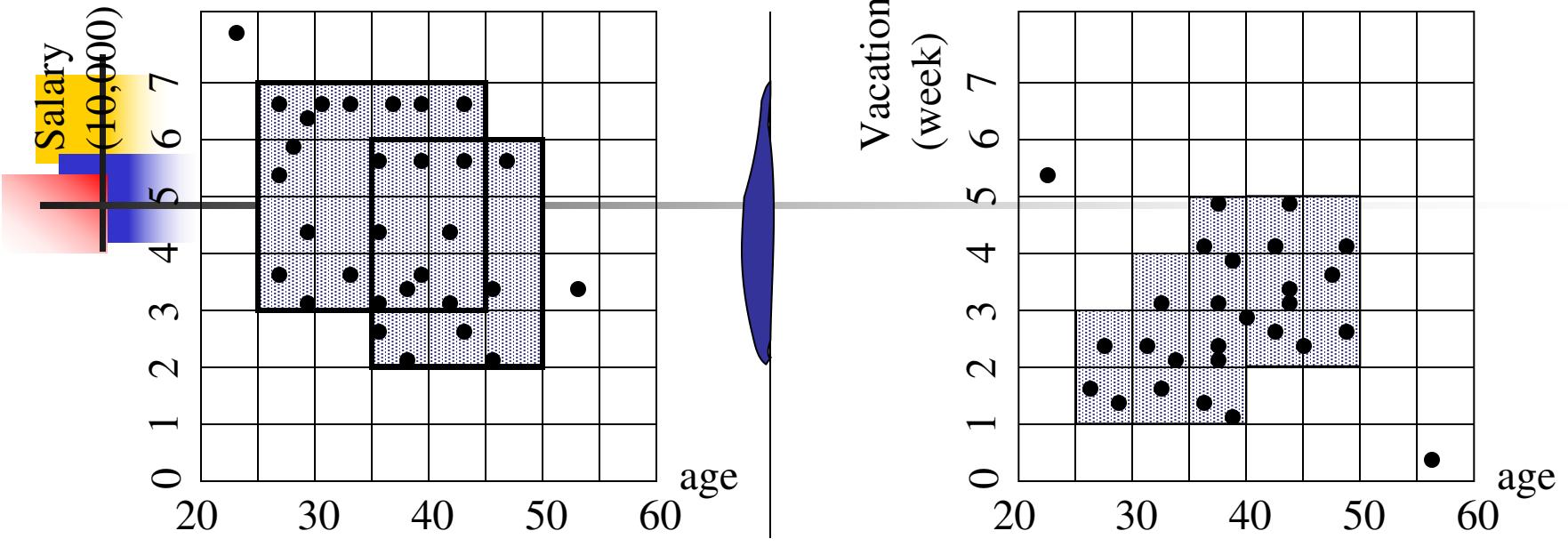
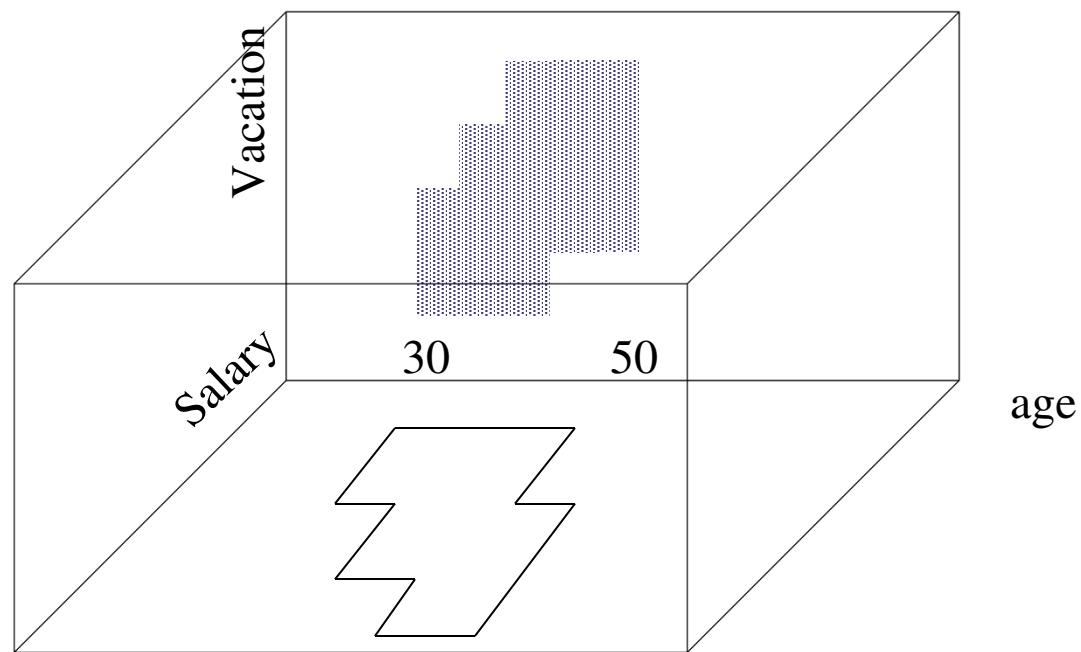
CLIQUE (Clustering In QUEst)

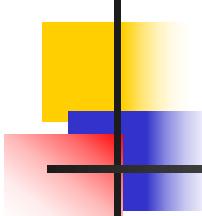
- n Agrawal, Gehrke, Gunopulos, Raghavan (SIGMOD'98).
- n Automatically identifying subspaces of a high dimensional data space that allow better clustering than original space
- n CLIQUE can be considered as both density-based and grid-based
 - n It partitions each dimension into the same number of equal length interval
 - n It partitions an m-dimensional data space into non-overlapping rectangular units
 - n A unit is dense if the fraction of total data points contained in the unit exceeds the input model parameter
 - n A cluster is a maximal set of connected dense units within a subspace



CLIQUE: The Major Steps

- n Partition the data space and find the number of points that lie inside each cell of the partition.
- n Identify the subspaces that contain clusters using the Apriori principle
- n Identify clusters:
 - n Determine dense units in all subspaces of interests
 - n Determine connected dense units in all subspaces of interests.
- n Generate minimal description for the clusters
 - n Determine maximal regions that cover a cluster of connected dense units for each cluster
 - n Determination of minimal cover for each cluster


 $= 3$




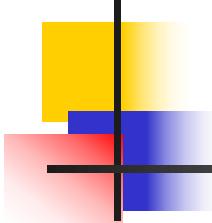
Strength and Weakness of CLIQUE

n Strength

- n It *automatically* finds subspaces of the highest dimensionality such that high density clusters exist in those subspaces
- n It is *insensitive* to the order of records in input and does not presume some canonical data distribution
- n It scales *linearly* with the size of input and has good scalability as the number of dimensions in the data increases

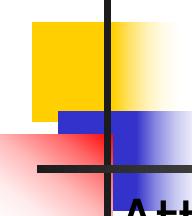
n Weakness

- n The accuracy of the clustering result may be degraded at the expense of simplicity of the method



Chapter 8. Cluster Analysis

- n What is Cluster Analysis?
- n Types of Data in Cluster Analysis
- n A Categorization of Major Clustering Methods
- n Partitioning Methods
- n Hierarchical Methods
- n Density-Based Methods
- n Grid-Based Methods
- n **Model-Based Clustering Methods**
- n Outlier Analysis
- n Summary

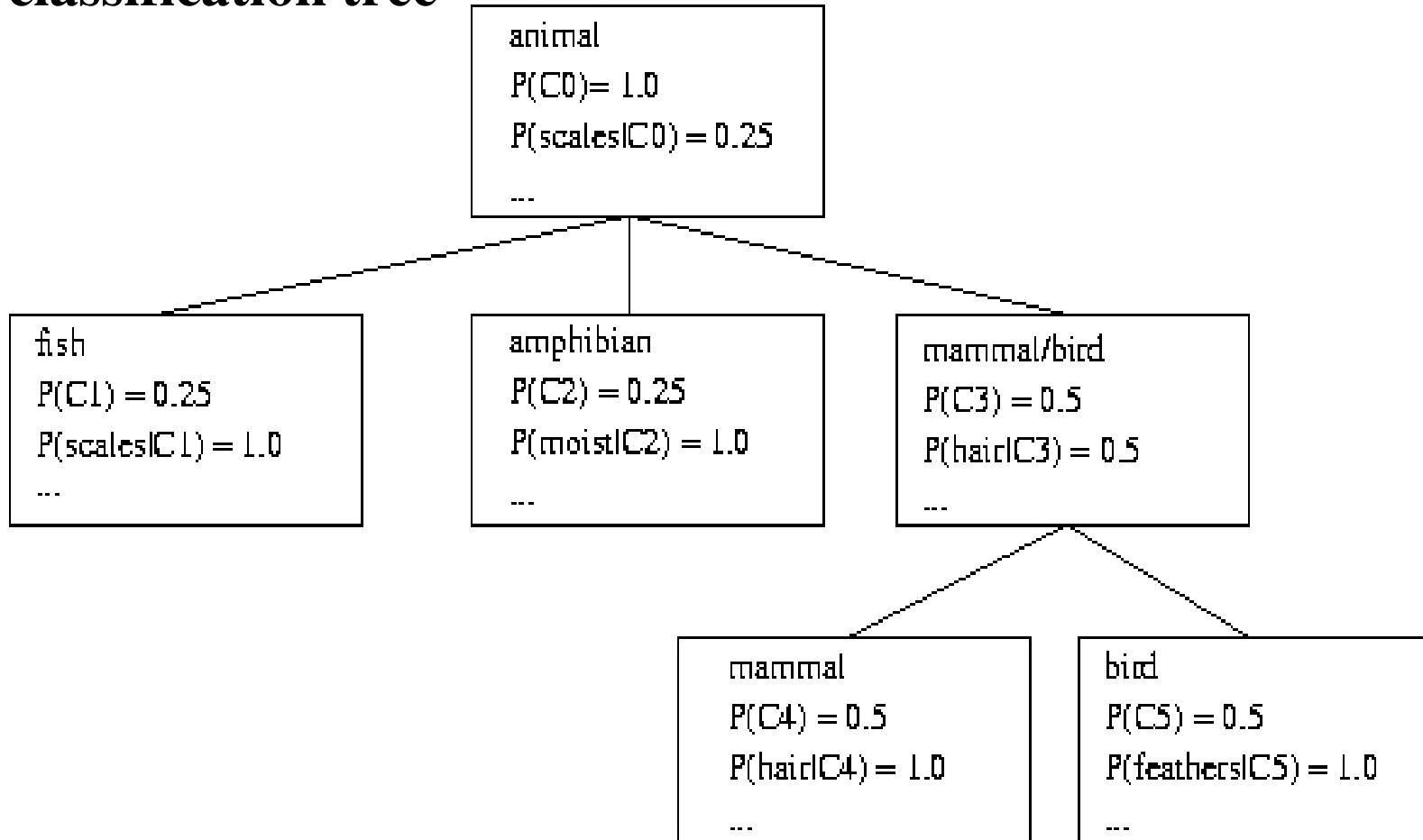


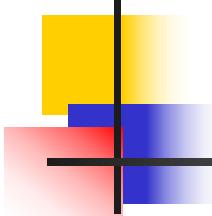
Model-Based Clustering Methods

- n Attempt to optimize the fit between the data and some mathematical model
- n Statistical and AI approach
 - n Conceptual clustering
 - n A form of clustering in machine learning
 - n Produces a classification scheme for a set of unlabeled objects
 - n Finds characteristic description for each concept (class)
 - n COBWEB (Fisher'87)
 - n A popular a simple method of incremental conceptual learning
 - n Creates a hierarchical clustering in the form of a **classification tree**
 - n Each node refers to a concept and contains a probabilistic description of that concept

COBWEB Clustering Method

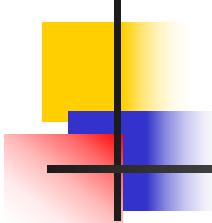
A classification tree





More on Statistical-Based Clustering

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

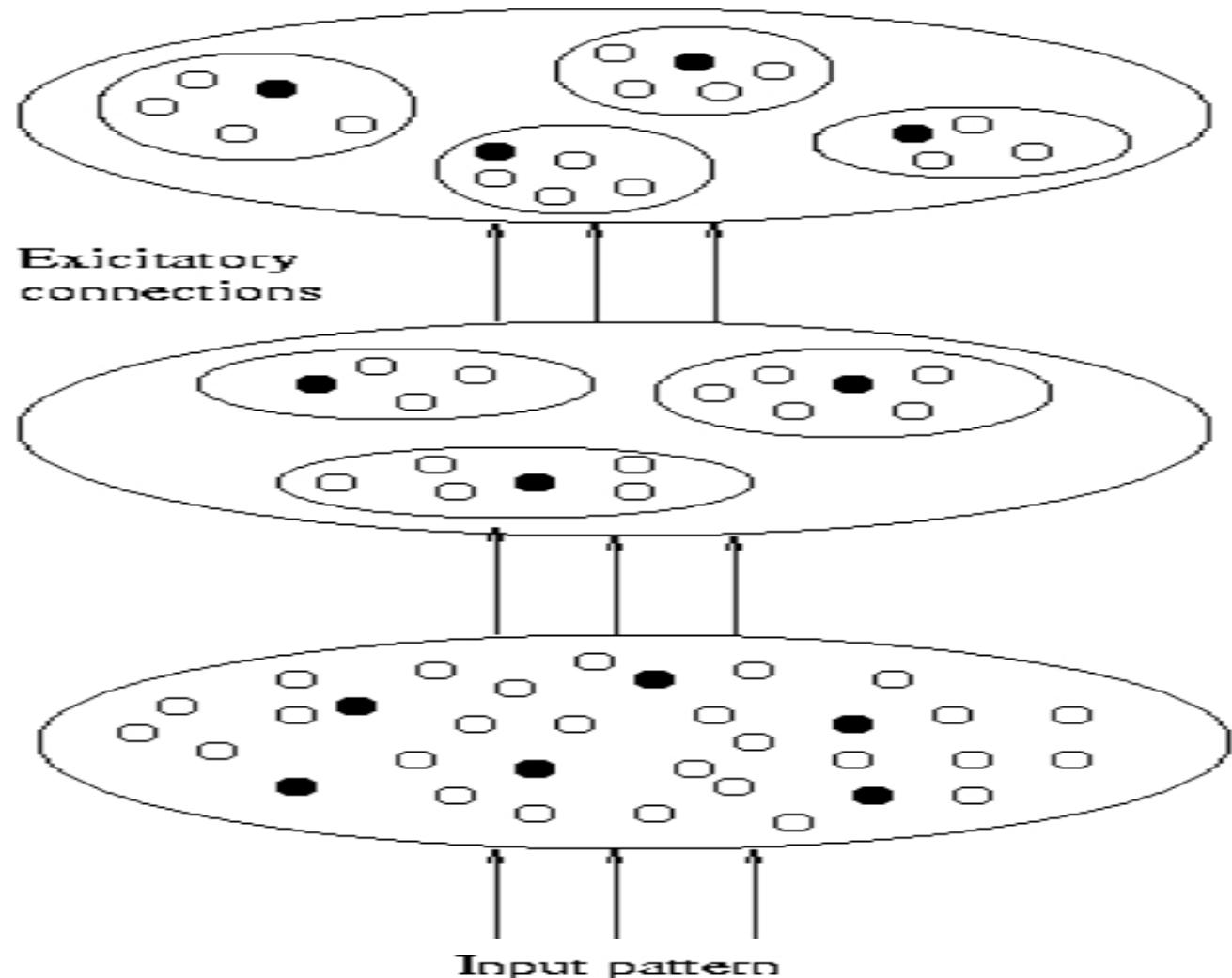


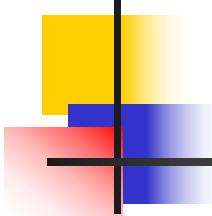
Other Model-Based Clustering Methods

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

Model-Based Clustering Methods

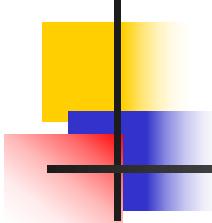
Layer 3
Inhibitory
clusters





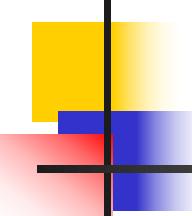
Self-organizing feature maps (SOMs)

- „ Clustering is also 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



Chapter 8. Cluster Analysis

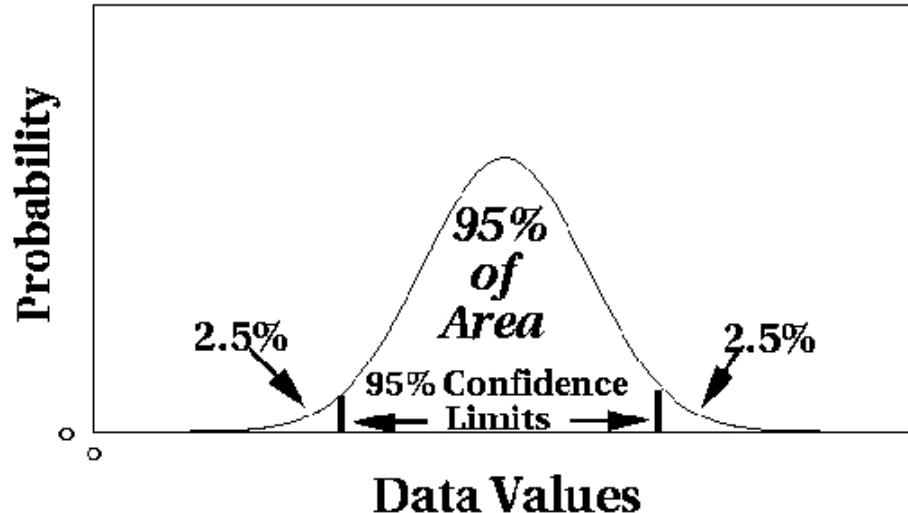
- n What is Cluster Analysis?
- n Types of Data in Cluster Analysis
- n A Categorization of Major Clustering Methods
 - n Partitioning Methods
 - n Hierarchical Methods
 - n Density-Based Methods
 - n Grid-Based Methods
 - n Model-Based Clustering Methods
- n **Outlier Analysis**
- n Summary



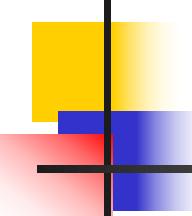
What Is Outlier Discovery?

- n What are outliers?
 - n The set of objects are considerably dissimilar from the remainder of the data
 - n Example: Sports: Michael Jordon, Wayne Gretzky, ...
- n Problem
 - n Find top n outlier points
- n Applications:
 - n Credit card fraud detection
 - n Telecom fraud detection
 - n Customer segmentation
 - n Medical analysis

Statistical Approaches

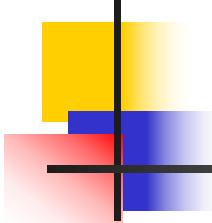


- f Assume a model underlying distribution that generates data set (e.g. normal distribution)
- n Use discordancy tests depending on
 - n data distribution
 - n distribution parameter (e.g., mean, variance)
 - n number of expected outliers
- n Drawbacks
 - n most tests are for single attribute
 - n In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



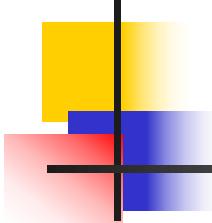
Outlier Discovery: Distance-Based Approach

- 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



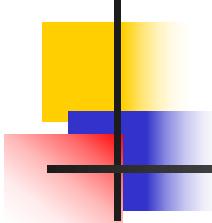
Outlier Discovery: Deviation-Based Approach

- 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 in large multidimensional data



Chapter 8. Cluster Analysis

- n What is Cluster Analysis?
- n Types of Data in Cluster Analysis
- n A Categorization of Major Clustering Methods
 - n Partitioning Methods
 - n Hierarchical Methods
 - n Density-Based Methods
 - n Grid-Based Methods
 - n Model-Based Clustering Methods
- n Outlier Analysis
- n **Summary**

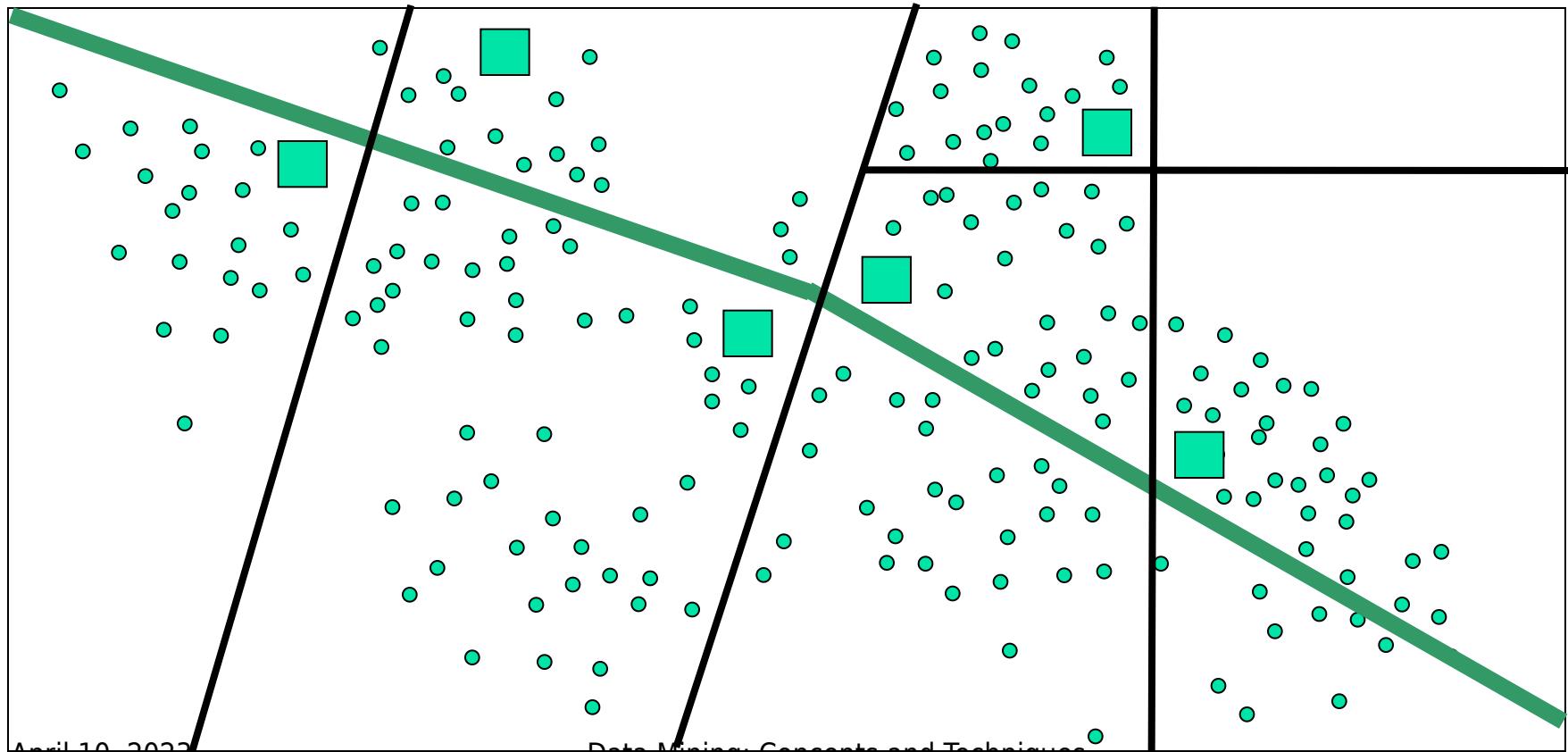


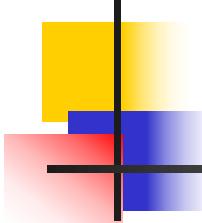
Problems and Challenges

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

Constraint-Based Clustering Analysis

- n Clustering analysis: less parameters but more user-desired constraints, e.g., an ATM allocation problem





Summary

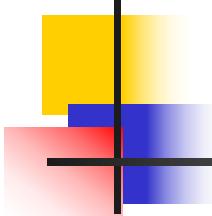
- n Cluster analysis groups objects based on their similarity and has wide applications
- n Measure of similarity can be computed for various types of data
- n Clustering algorithms can be categorized into partitioning methods, hierarchical methods, density-based methods, grid-based methods, and model-based methods
- n Outlier detection and analysis are very useful for fraud detection, etc. and can be performed by statistical, distance-based or deviation-based approaches
- n There are still lots of research issues on cluster analysis, such as constraint-based clustering

References (1)

- R. Agrawal, J. Gehrke, D. Gunopulos, and P. Raghavan. Automatic subspace clustering of high dimensional data for data mining applications. SIGMOD'98
- M. R. Anderberg. Cluster Analysis for Applications. Academic Press, 1973.
- M. Ankerst, M. Breunig, H.-P. Kriegel, and J. Sander. Optics: Ordering points to identify the clustering structure, SIGMOD'99.
- P. Arabie, L. J. Hubert, and G. De Soete. Clustering and Classification. World Scientific, 1996
- M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases. KDD'96.
- M. Ester, H.-P. Kriegel, and X. Xu. Knowledge discovery in large spatial databases: Focusing techniques for efficient class identification. SSD'95.
- D. Fisher. Knowledge acquisition via incremental conceptual clustering. Machine Learning, 2:139-172, 1987.
- D. Gibson, J. Kleinberg, and P. Raghavan. Clustering categorical data: An approach based on dynamic systems. In Proc. VLDB'98.
- S. Guha, R. Rastogi, and K. Shim. Cure: An efficient clustering algorithm for large databases. SIGMOD'98.
- A. K. Jain and R. C. Dubes. Algorithms for Clustering Data. Prentice Hall, 1988.

References (2)

- n L. Kaufman and P. J. Rousseeuw. Finding Groups in Data: an Introduction to Cluster Analysis. John Wiley & Sons, 1990.
- n E. Knorr and R. Ng. Algorithms for mining distance-based outliers in large datasets. VLDB'98.
- n G. J. McLachlan and K.E. Bkasford. Mixture Models: Inference and Applications to Clustering. John Wiley and Sons, 1988.
- n P. Michaud. Clustering techniques. Future Generation Computer systems, 13, 1997.
- n R. Ng and J. Han. Efficient and effective clustering method for spatial data mining. VLDB'94.
- n E. Schikuta. Grid clustering: An efficient hierarchical clustering method for very large data sets. Proc. 1996 Int. Conf. on Pattern Recognition, 101-105.
- n G. Sheikholeslami, S. Chatterjee, and A. Zhang. WaveCluster: A multi-resolution clustering approach for very large spatial databases. VLDB'98.
- n W. Wang, Yang, R. Muntz, STING: A Statistical Information grid Approach to Spatial Data Mining, VLDB'97.
- n T. Zhang, R. Ramakrishnan, and M. Livny. BIRCH : an efficient data clustering method for very large databases. SIGMOD'96.



<http://www.cs.sfu.ca/~han>



Thank you !!!