

Chapter 7. Cluster Analysis

1. What is Cluster Analysis?
2. Types of Data in Cluster Analysis
3. A Categorization of Major Clustering Methods
4. Partitioning Methods
5. Hierarchical Methods
6. Density-Based Methods
7. Clustering High-Dimensional Data
8. Constraint-Based Clustering
9. Outlier Analysis
10. Summary



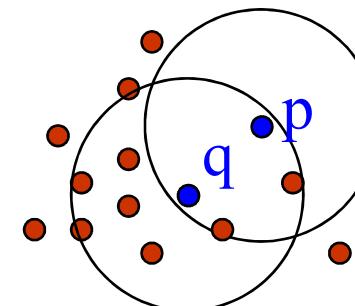
Density-Based Clustering Methods

- Clustering based on **density** (local cluster criterion), such as density-connected points, rather than just a distance
- Major features:
 - Discover clusters of arbitrary shape
 - Handle noise
 - One scan, thus being efficient
 - Need density parameters as termination condition
- Several interesting studies:
 - DBSCAN: Ester, et al. (KDD'96)
 - OPTICS: Ankerst, et al (SIGMOD'99).
 - CLIQUE: Agrawal, et al. (SIGMOD'98) (more grid-based)



Density-Based Clustering: Basic Concepts

- Two parameters:
 - *Eps*: Maximum radius of the neighborhood
 - *MinPts*: Minimum number of points in an *Eps*-neighborhood of a given 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* and *MinPts* if
 - p belongs to $N_{Eps}(q)$
 - **core point condition**:
 $|N_{Eps}(q)| \geq MinPts$
 - Note: *Not symmetric*



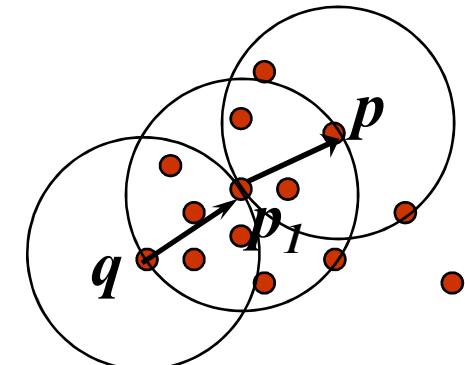
MinPts = 5

Eps = 1 cm

Density-Reachable and Density-Connected

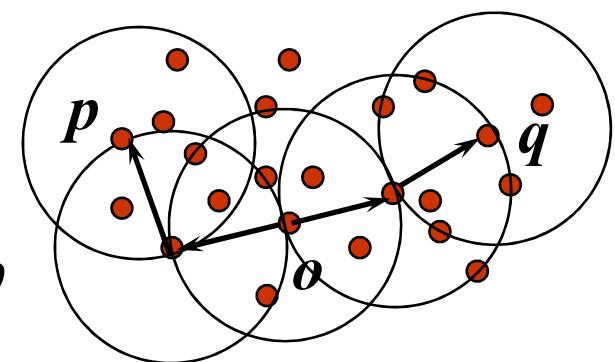
- Density-reachable:

- A point p is **density-reachable** from a point q w.r.t. Eps and $MinPts$ 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



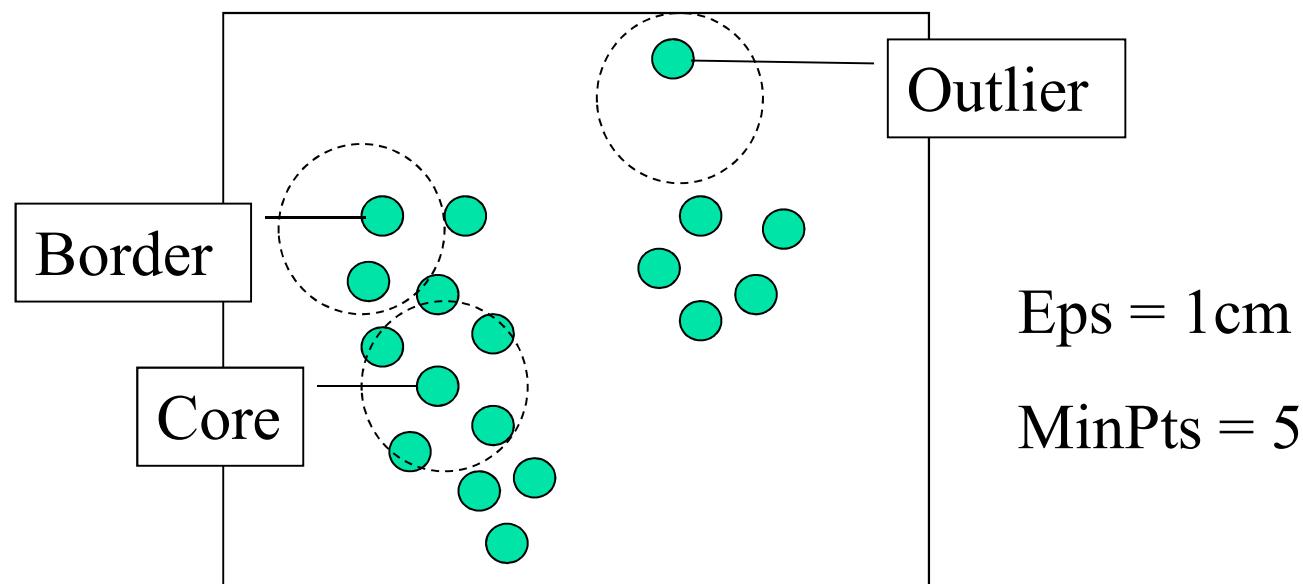
- Density-connected

- A point p is **density-connected** to a point q w.r.t. Eps and $MinPts$ 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 an *arbitrary shape* in spatial databases with noise



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



DBSCAN: Sensitive to Parameters

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

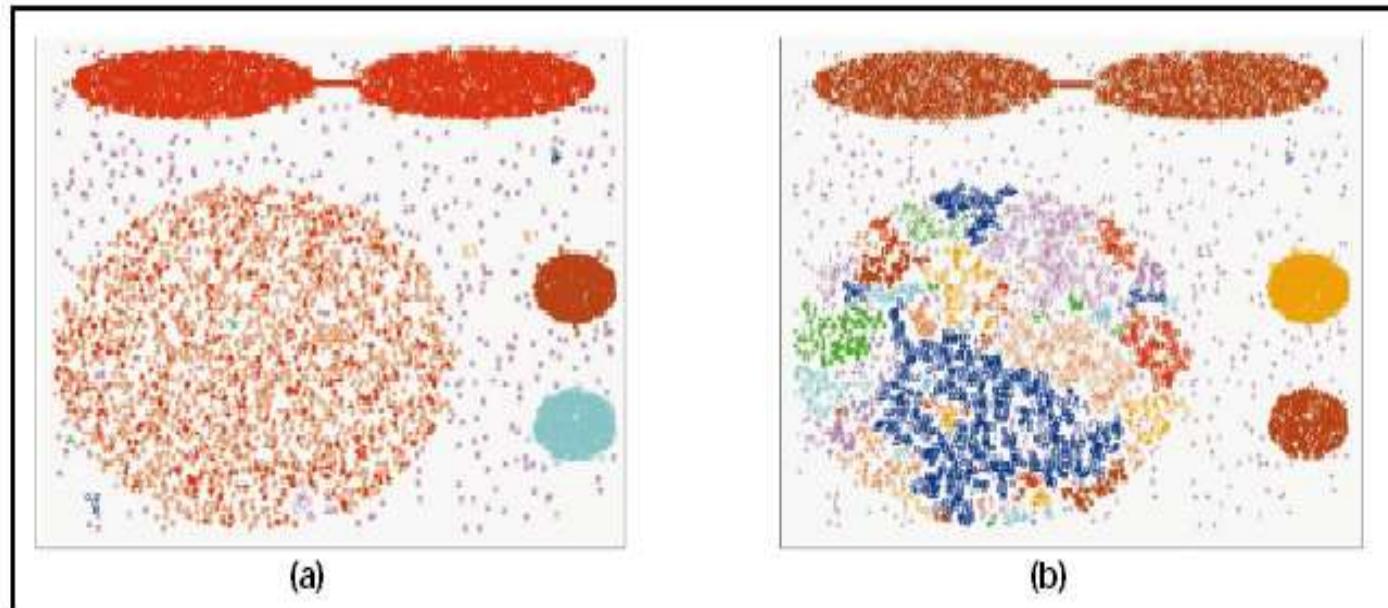
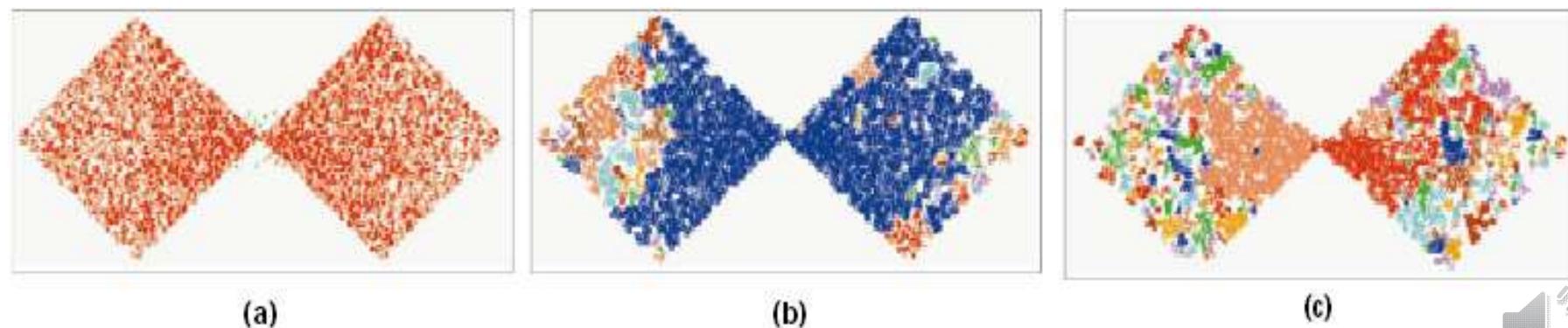
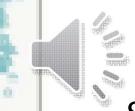
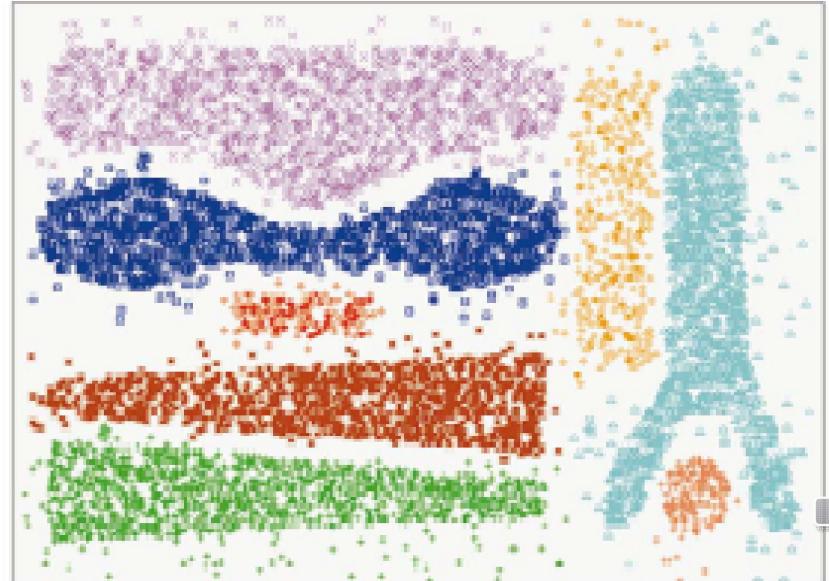
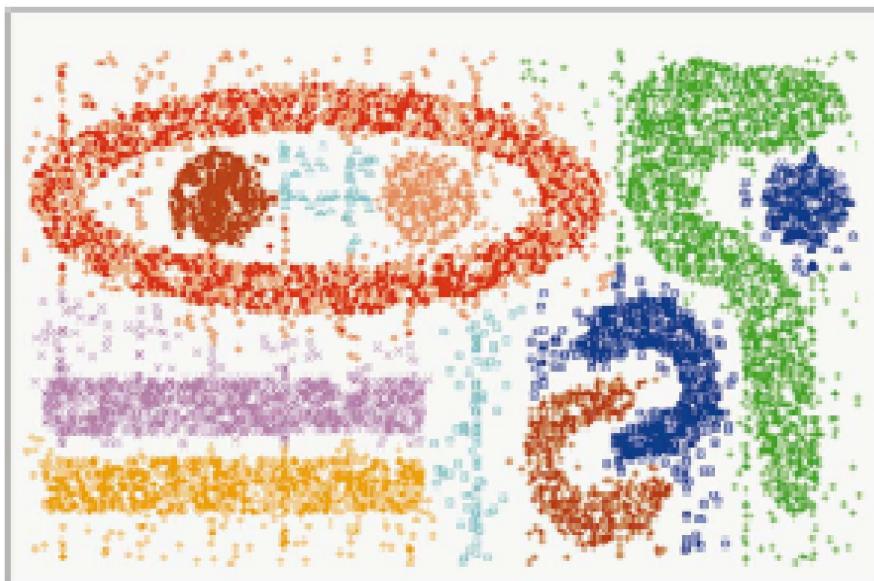
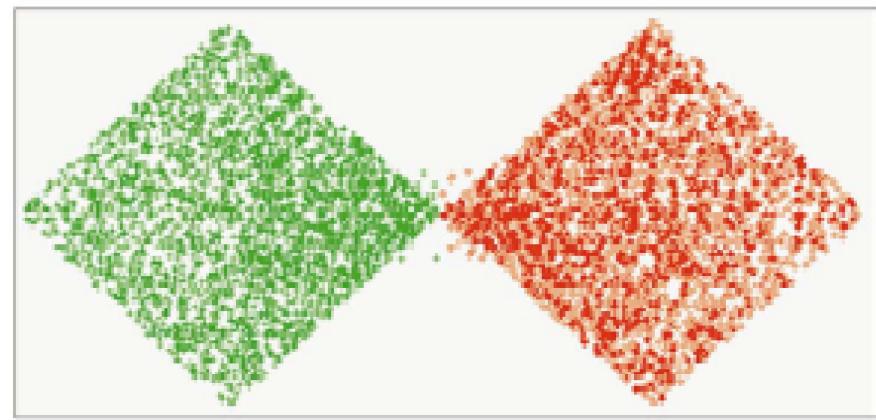
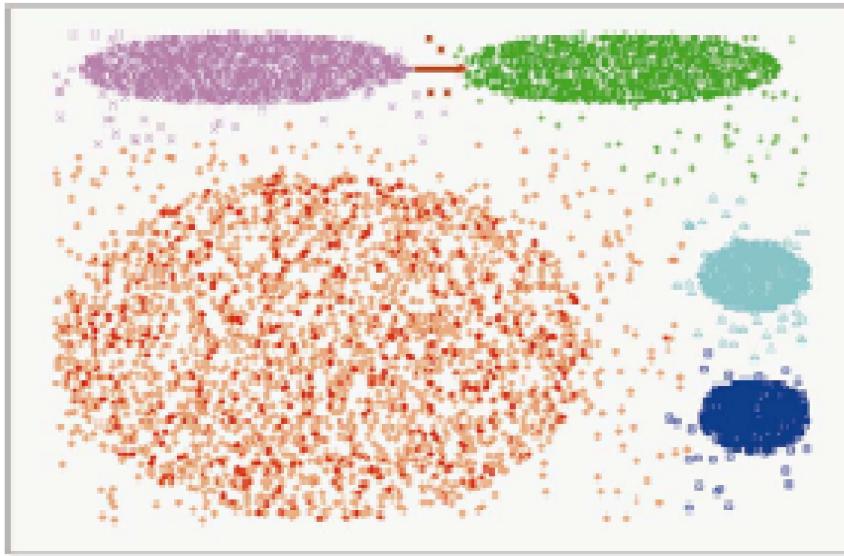


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



CHAMELEON (Clustering Complex Objects)

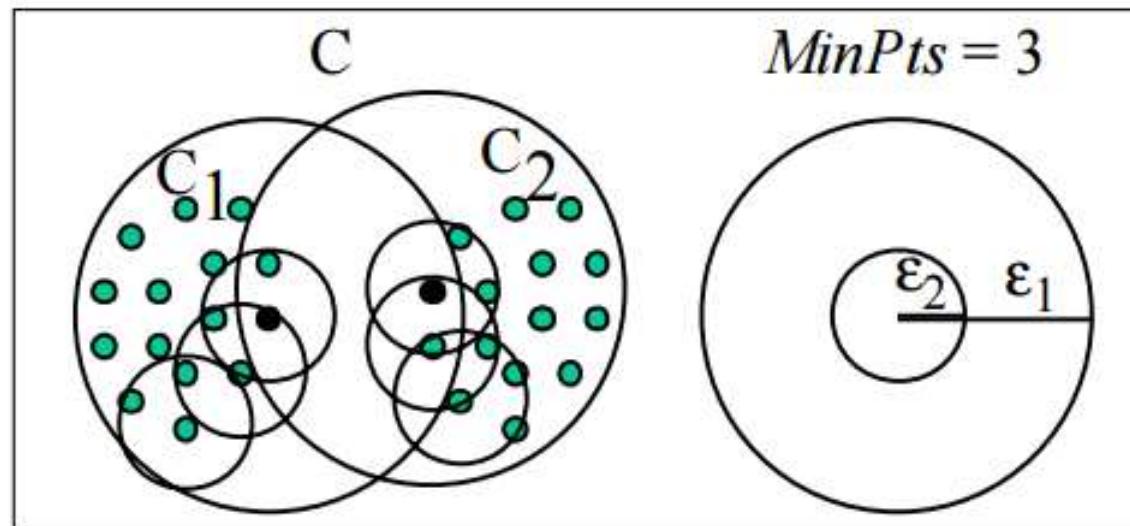


OPTICS: A Cluster-Ordering Method (1999)

- OPTICS: Ordering Points To Identify the Clustering Structure
 - Ankerst, Breunig, Kriegel, and Sander (SIGMOD'99)
 - Produces a special order of objects in the database w.r.t. its density-based clustering structure
 - This cluster-ordering contains the information equivalent to different density-based clustering structure corresponding to a broad range of parameter settings (*Eps*)
 - Good for both automatic and interactive cluster analysis, including finding intrinsic clustering structure
 - Can be represented graphically or using visualization techniques



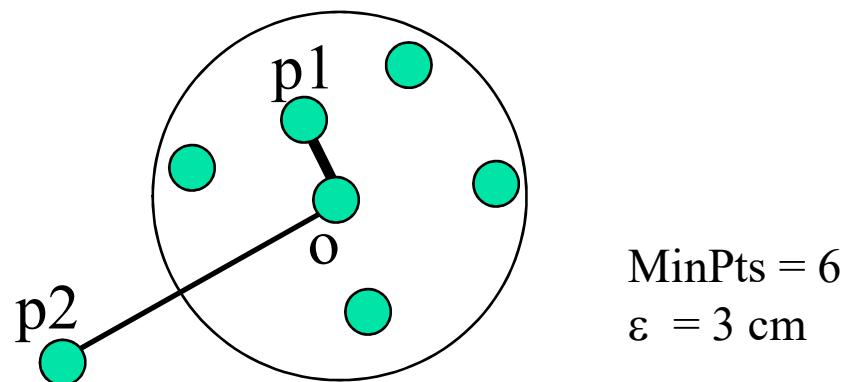
OPTICS: A Cluster-Ordering Method (1999)



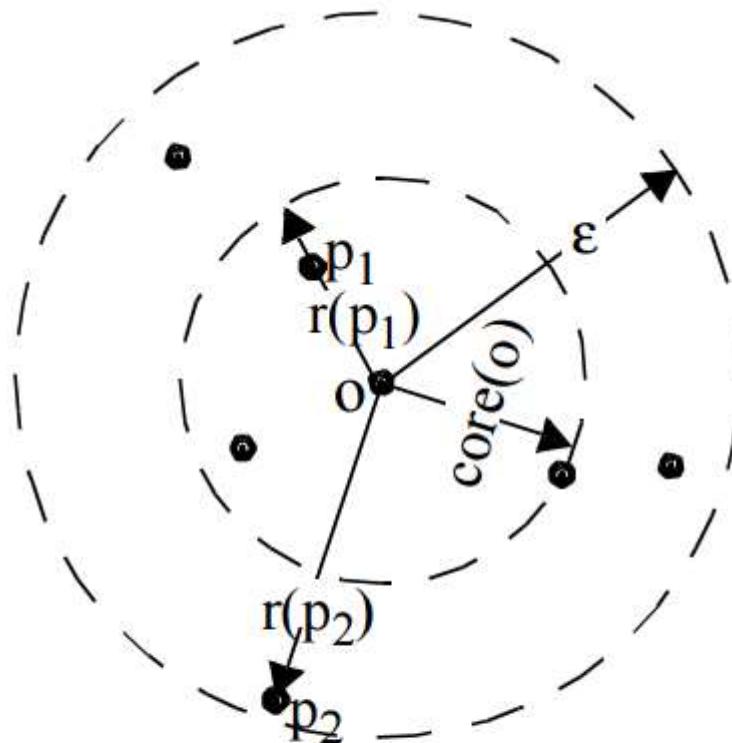
**Figure 3. Illustration of “nested”
density-based clusters**

OPTICS: Some Extension from DBSCAN

- Core Distance (of o)
 - Distance to make the object a core
- Reachability Distance (of p from o)
 - $r(p, o) = \max(\text{core-distance}(o), d(o, p))$
 - Ex: $r(p_1, o) = 2.8\text{cm}$, $r(p_2, o) = 4\text{cm}$

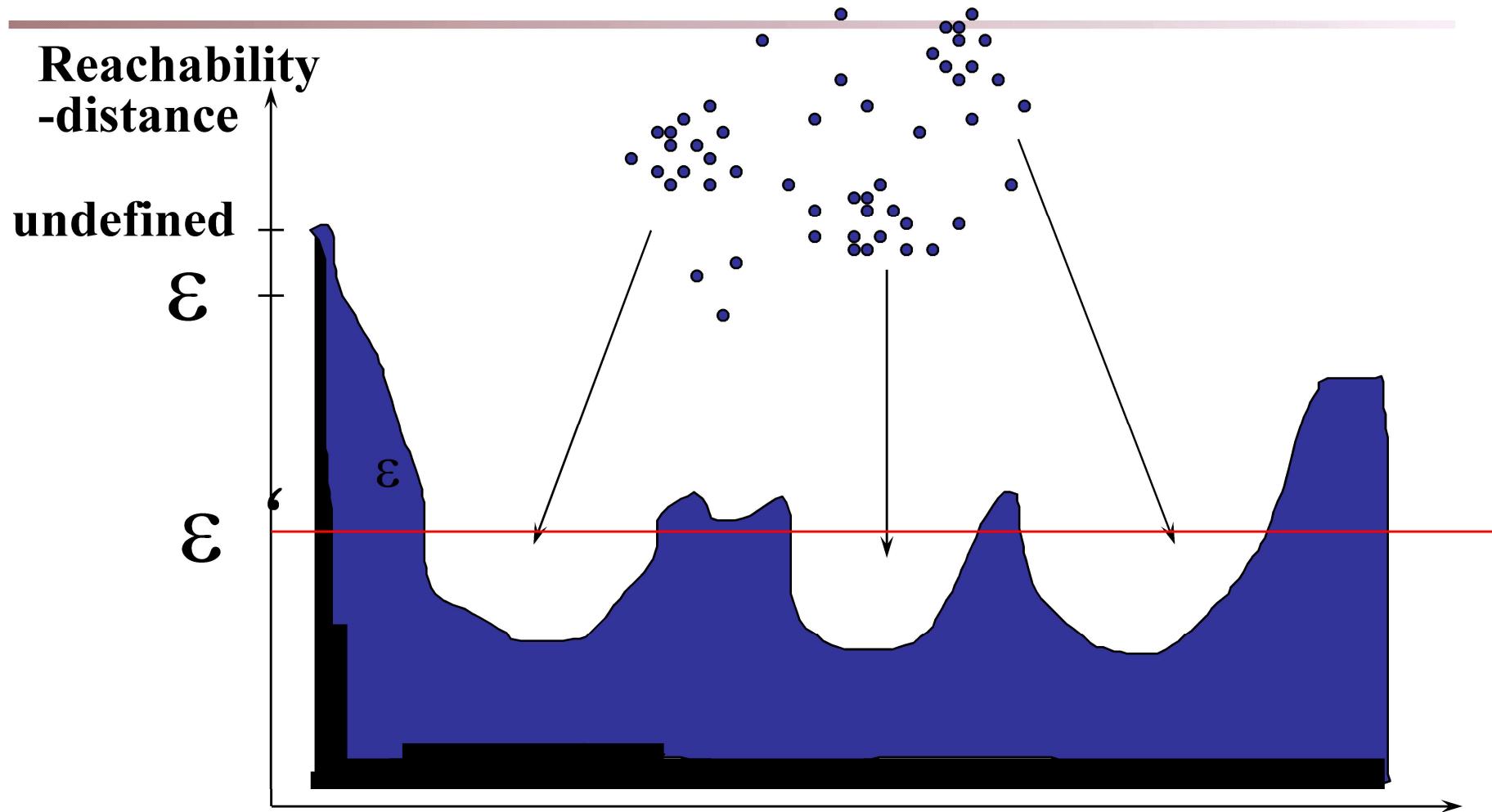


OPTICS: A Cluster-Ordering Method (1999)



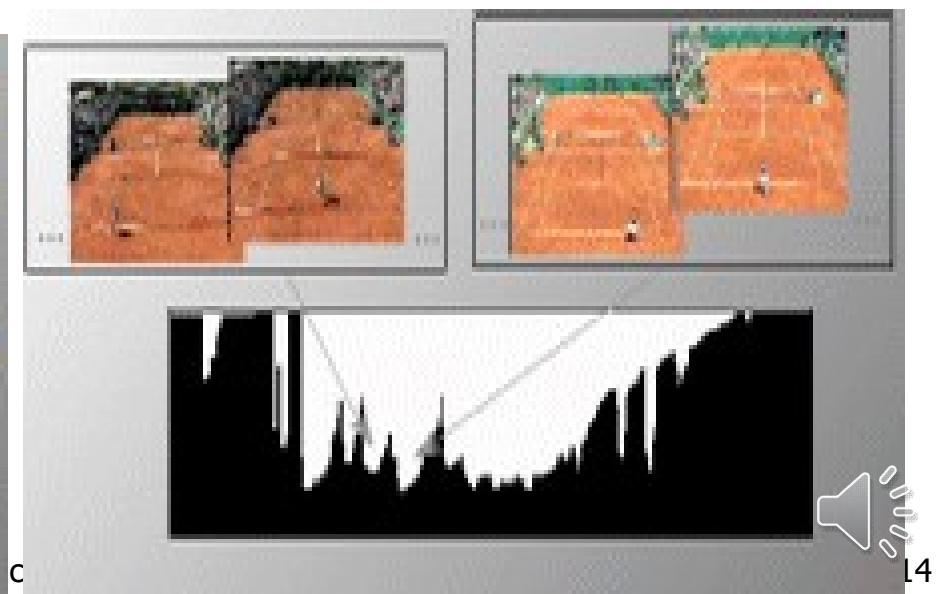
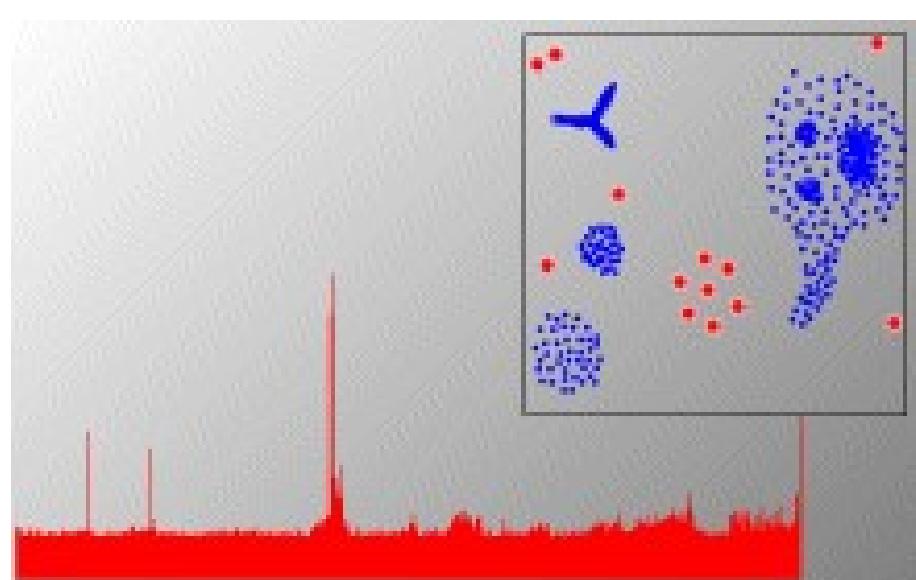
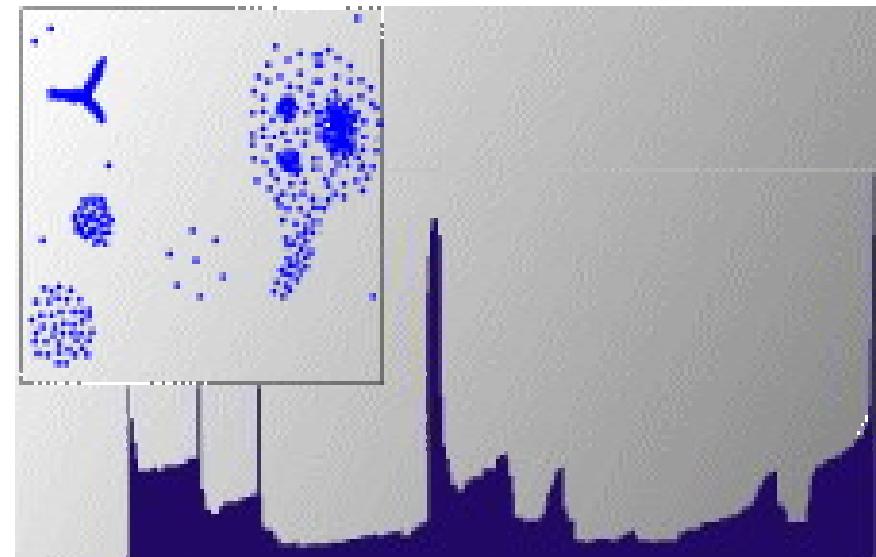
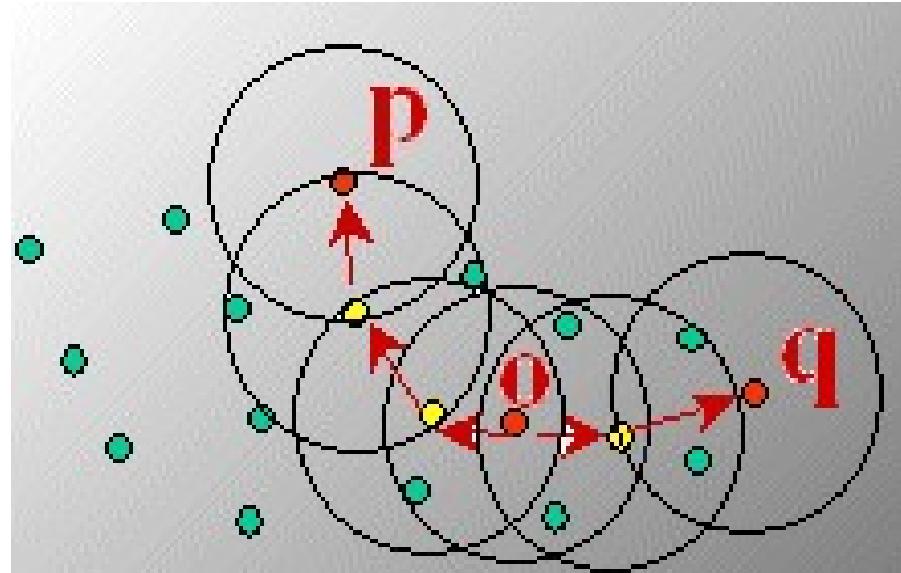
**Figure 4. Core-distance(o),
reachability-distances
 $r(p_1, o)$, $r(p_2, o)$ for $MinPts=4$**





**Cluster-order
of the objects**  13

Density-Based Clustering: OPTICS & Its Applications



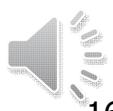
OPTICS: A Cluster-Ordering Method (1999)

- OPTICS: Ordering Points To Identify the Clustering Structure
 - Ankerst, Breunig, Kriegel, and Sander (SIGMOD'99)
 - Produces a special order of objects in the database w.r.t. its density-based clustering structure
 - This cluster-ordering contains the information equivalent to different density-based clustering structure corresponding to a broad range of parameter settings (*Eps*)
 - Good for both automatic and interactive cluster analysis, including finding intrinsic clustering structure
 - Can be represented graphically or using visualization techniques



Chapter 7. Cluster Analysis

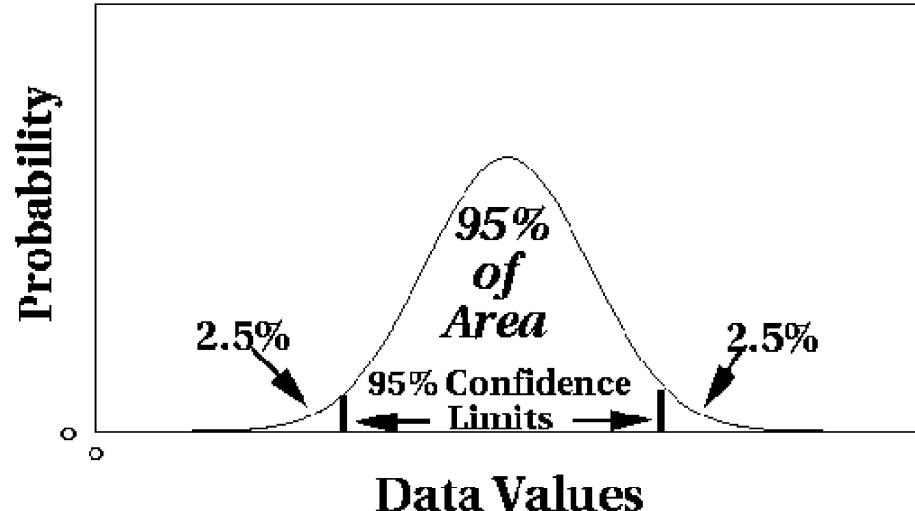
1. What is Cluster Analysis?
2. Types of Data in Cluster Analysis
3. A Categorization of Major Clustering Methods
4. Partitioning Methods
5. Hierarchical Methods
6. Density-Based Methods
7. Clustering High-Dimensional Data
8. Constraint-Based Clustering
9. Outlier Analysis
10. Summary



What Is Outlier Discovery?

- What are outliers?
 - The set of objects are considerably dissimilar from the remainder of the data
 - Example: Sports: Michael Jordon, Wayne Gretzky, ...
- Problem: Define and find outliers in large data sets
- Applications:
 - Credit card fraud detection
 - Telecom fraud detection
 - Customer segmentation
 - Medical analysis

Outlier Discovery: Statistical Approaches



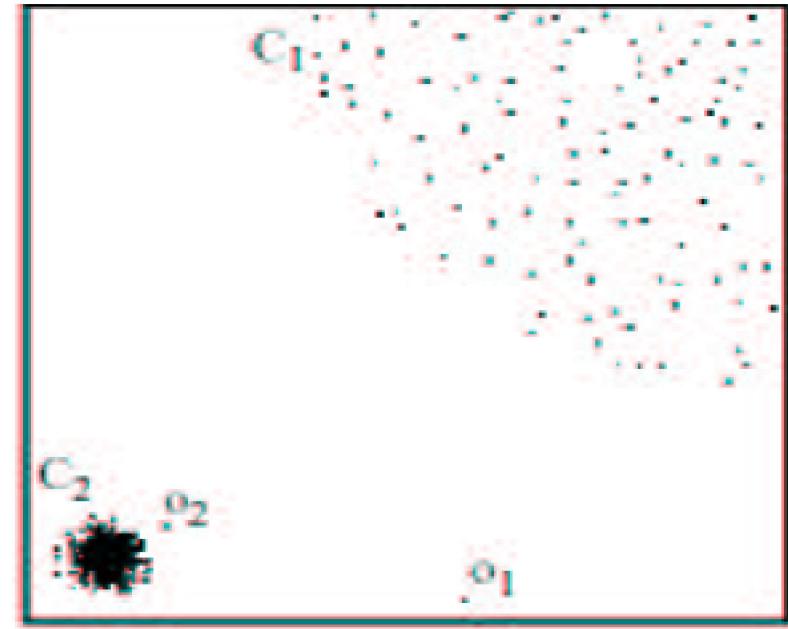
- Assume a model underlying distribution that generates data set (e.g. normal distribution)
- Use discordancy tests depending on
 - data distribution
 - distribution parameter (e.g., mean, variance)
 - number of expected outliers
- Drawbacks
 - most tests are for a *single attribute*
 - In many cases, 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

Density-Based Local Outlier Detection

- Distance-based outlier detection is based on global distance distribution
- It encounters difficulties to identify outliers *if data is not uniformly distributed*
- Ex. C_1 contains 400 loosely distributed points, C_2 has 100 tightly condensed points, 2 outlier points o_1, o_2
- Distance-based method cannot identify o_2 as an outlier
- Need the concept of *a local outlier*



- Local outlier factor (LOF)
 - Assume outlier is not crisp
 - Each point has a LOF

Chapter 7. Cluster Analysis

1. What is Cluster Analysis?
2. Types of Data in Cluster Analysis
3. A Categorization of Major Clustering Methods
4. Partitioning Methods
5. Hierarchical Methods
6. Density-Based Methods
7. Clustering High-Dimensional Data
8. Constraint-Based Clustering
9. Outlier Analysis
10. Summary



Summary

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

Problems and Challenges

- Considerable progress has been made in scalable clustering methods
 - Partitioning: k-means, k-medoids, CLARANS
 - Hierarchical: BIRCH, ROCK, CHAMELEON
 - Density-based: DBSCAN, OPTICS, DenClue
 - Constraint-based: COD, constrained-clustering
- Current clustering techniques do not address all the requirements adequately, still an active area of research

OPTICS: Some Extension from DBSCAN

- Index-based:
 - $k = \text{number of dimensions}$
 - $N = 20$
 - $p = 75\%$
 - $M = N(1-p) = 5$
- Complexity: $O(kN^2)$

■ Core Distance

Distance to make the object a core

■ Reachability Distance

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

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

