

# **Data Mining Fundamentals**

## **Chapter 10. Cluster Analysis: Basic Concepts and Methods**

# Density-Based and Grid-Based Clustering Methods

---

## □ Density-Based Clustering

- Basic Concepts

- DBSCAN: A Density-Based Clustering Algorithm

- OPTICS: Ordering Points To Identify Clustering Structure

## □ Grid-Based Clustering Methods

- Basic Concepts

- STING: A Statistical Information Grid Approach

- CLIQUE: Grid-Based Subspace Clustering

# Density-Based Clustering Methods

---

- ❑ Clustering based on density (a local cluster criterion), such as density-connected points
- ❑ Major features:
  - ❑ Discover clusters of arbitrary shape
  - ❑ Handle noise
  - ❑ One scan (only examine the local region to justify density)
  - ❑ 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) (also, grid-based)

# DBSCAN: A Density-Based Spatial Clustering Algorithm

DBSCAN (M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, KDD'96)

Discovers clusters of arbitrary shape: Density-Based Spatial Clustering of Applications with Noise

A *density-based* notion of cluster

A *cluster* is defined as a maximal set of density-connected points

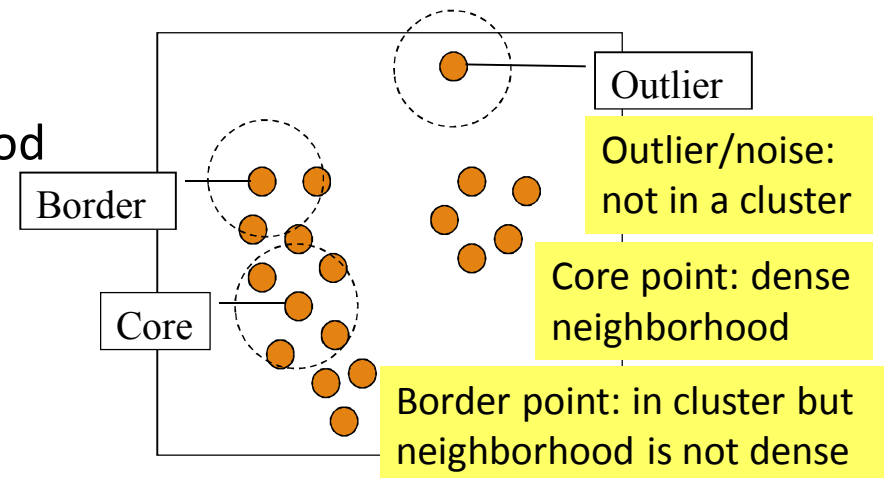
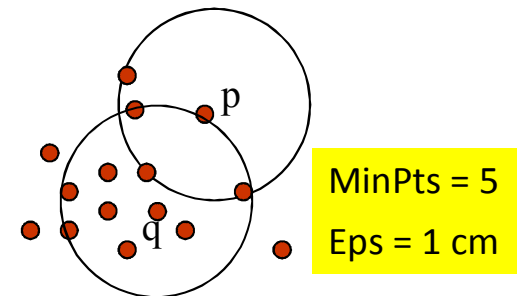
Two parameters:

**Eps** ( $\epsilon$ ): Maximum radius of the neighborhood

**MinPts**: Minimum number of points in the Eps-neighborhood of a point

The Eps( $\epsilon$ )-neighborhood of a point  $q$ :

$N_{Eps}(q) = \{p \text{ belongs to } D \mid \text{dist}(p, q) \leq Eps\}$



# DBSCAN: Density-Reachable and Density-Connected

## □ Directly density-reachable:

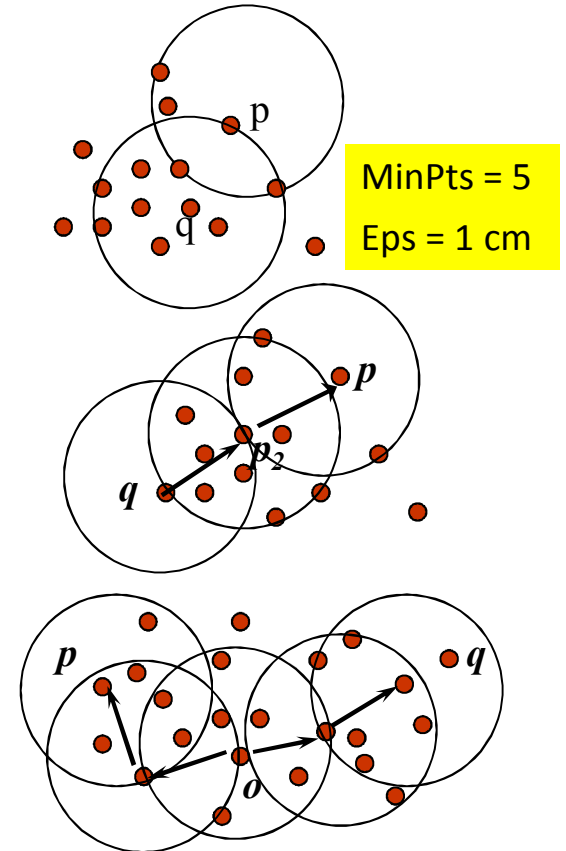
- A point  $p$  is **directly density-reachable** from a point  $q$  w.r.t.  $Eps$  ( $\epsilon$ ),  $MinPts$  if
  - $p$  belongs to  $N_{Eps}(q)$
  - **core point** condition:  $|N_{Eps}(q)| \geq MinPts$

## □ Density-reachable:

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

## □ Density-connected:

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



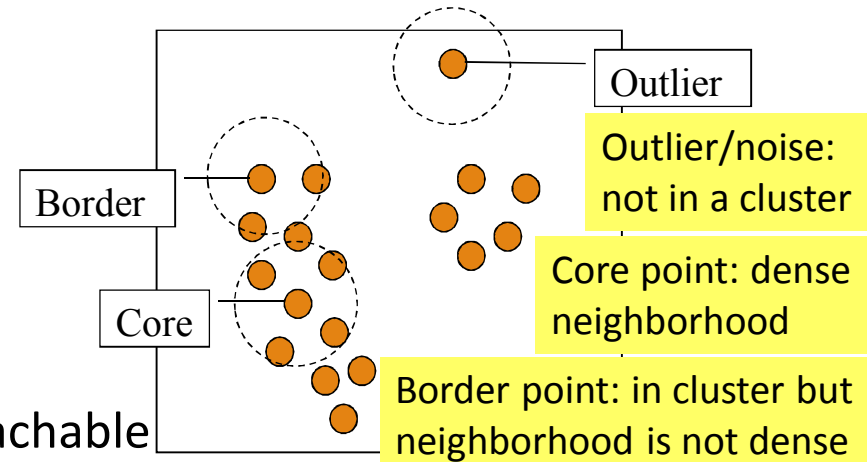
# DBSCAN: The Algorithm

## Algorithm

- Arbitrarily 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

## Computational complexity

- If a spatial index is used, the computational complexity of DBSCAN is  $O(n \log n)$ , where  $n$  is the number of database objects
- Otherwise, the complexity is  $O(n^2)$



# DBSCAN Is Sensitive to the Setting of Parameters

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

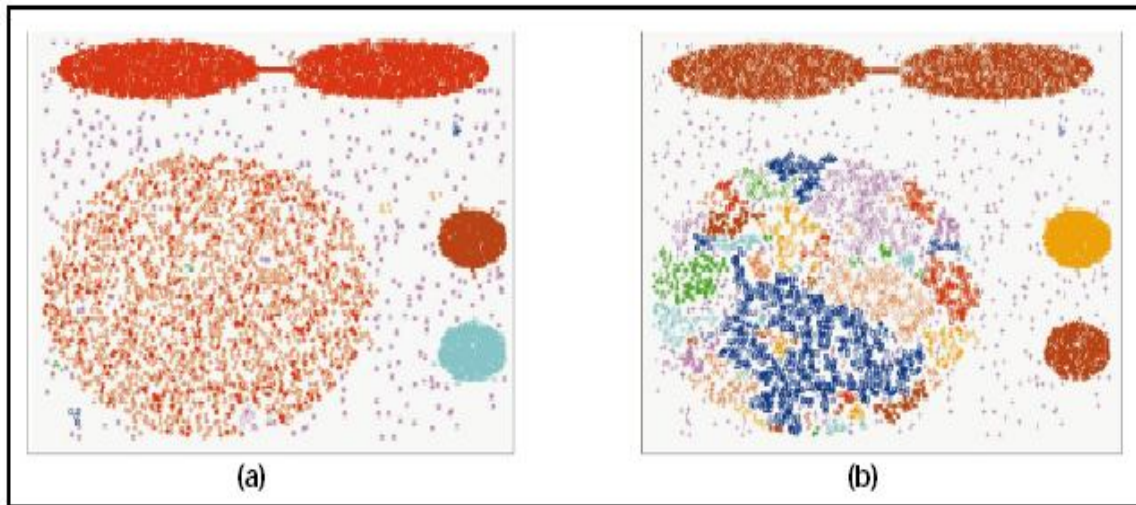
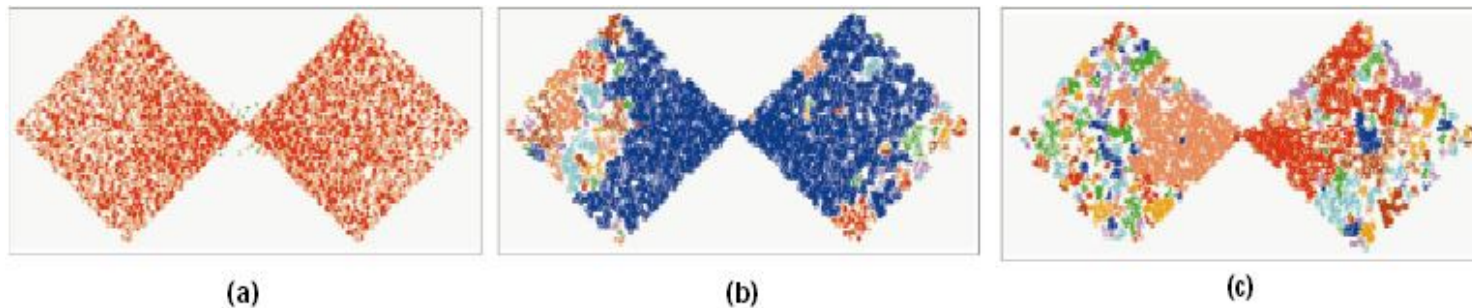


Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



Ack. Figures from G. Karypis, E.-H. Han, and V. Kumar, *COMPUTER*, 32(8), 1999

# Grid-Based Clustering Methods

---

- ❑ Grid-Based Clustering: Explore multi-resolution grid data structure in clustering
  - ❑ Partition the data space into a finite number of cells to form a grid structure
  - ❑ Find clusters (dense regions) from the cells in the grid structure
- ❑ Features and challenges of a typical grid-based algorithm
  - ❑ Efficiency and scalability: # of cells  $\ll$  # of data points
  - ❑ Uniformity: Uniform, hard to handle highly irregular data distributions
  - ❑ Locality: Limited by predefined cell sizes, borders, and the density threshold
  - ❑ Curse of dimensionality: Hard to cluster high-dimensional data
- ❑ Methods to be introduced
  - ❑ **STING** (a Statistical Information Grid approach) (Wang, Yang and Muntz, VLDB'97)
  - ❑ **CLIQUE** (Agrawal, Gehrke, Gunopulos, and Raghavan, SIGMOD'98)
    - ❑ Both grid-based and subspace clustering



# STING: A Statistical Information Grid Approach

- ❑ STING (Statistical Information Grid) (Wang, Yang and Muntz, VLDB'97)
- ❑ The spatial area is divided into rectangular cells at different levels of resolution, and these cells form a tree structure
- ❑ A cell at a high level contains a number of smaller cells of the next lower level
- ❑ Statistical information of each cell is calculated and stored beforehand and is used to answer queries
- ❑ Parameters of higher level cells can be easily calculated from that of lower level cell, including
  - ❑ *count, mean, s*(standard deviation), *min, max*
  - ❑ type of distribution—*normal, uniform, etc.*

