



# Cluster Analysis

—Density-Based Methods—

徐华

清华大学 计算机系 智能技术与系统国家重点实验室

xuhua@tsinghua.edu.cn

1

## Cluster Analysis



- ◉ What is Cluster Analysis?
- ◉ Types of Data in Cluster Analysis
- ◉ A Categorization of Major Clustering Methods
- ◉ Partitioning Methods
- ◉ Hierarchical Methods
- ◉ **Density-Based Methods**
- ◉ Grid-Based Methods
- ◉ Model-Based Clustering Methods
- ◉ Outlier Analysis
- 2 ◉ **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)

3



## Density Concepts



- ◉ Core object (中心对象CO)–object with at least 'M' objects within a radius 'E-neighborhood'
- ◉ Directly density reachable (直接密度可达DDR)–x is CO, y is in x' s 'E-neighborhood'
- ◉ Density reachable (密度可达) – there exists a chain of DDR objects from x to y
- ◉ Density connected objects(密度相连) - there exists a O, p and q are density reachable to O respectively.
- ◉ Density based cluster–density connected objects maximum w.r.t. reachability

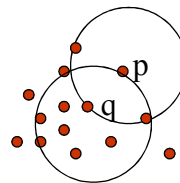
4



## Density-Based Clustering: Background



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



MinPts = 5

Eps = 1 cm

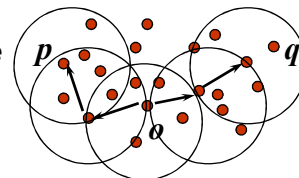
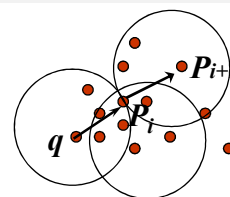
5



## Density-Based Clustering: Background (II)



- ◉ Density-reachable:
  - ◆ 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  $p_i$
- ◉ Density-connected
  - ◆ 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*.



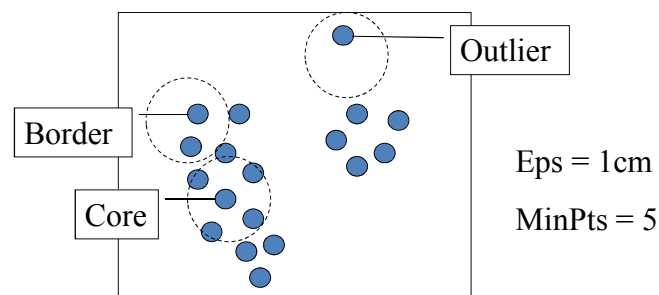
6



## DBSCAN: Density Based Spatial Clustering of Applications with Noise



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



7



## 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.

8





**Thanks!**

