

# Clustering

M.M. Pedram  
[pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)  
Tarbiat Moallem University of Tehran  
(Fall 2011)

1

## Density-Based Clustering Methods

Lecture slides taken/modified from:  
– Jiawei Han ([http://www.sai.cs.uiuc.edu/~hanj/DM\\_Book.html](http://www.sai.cs.uiuc.edu/~hanj/DM_Book.html))  
– Vipin Kumar (<http://www-users.cs.umn.edu/~kumar/csci5980/index.html>)

2

## 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)

▶ 3

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

3

## DBSCAN (M. Ester, et al. 1996)

- ❖ **DBSCAN**: **D**ensity **B**ased **S**patial **C**lustering of **A**pplications with **N**oise
- ❖ Outliers will not effect creation of cluster.
- ❖ Have, in principle, the ability to discover clusters of *arbitrarily shaped clusters*.
- ❖ One scan
- ❖ Need density parameters as termination condition
- ❖ Time complexity less than  $O(N^2)$ .

▶ 4

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

4

## DBSCAN

### Nomenclature:

- ❖  $ball_\epsilon(q)$  : the hyper-sphere (n-dimensional ball) of radius  $\epsilon$  centered at  $q$ .
- ❖  $N_\epsilon(q)$  : the set of points lying in  $ball_\epsilon(q)$ .  

$$N_\epsilon(q) = \{x \mid dist(x,q) \leq \epsilon\} \quad \text{or} \quad N_\epsilon(q) = \{x \mid x \in ball_\epsilon(q)\}$$
- ❖  $|N_\epsilon(q)|$  : the number of points lying in  $ball_\epsilon(q)$ .

### 2 input Parameters:

- ❖  $\epsilon$  : Maximum radius of the neighborhood.
- ❖ **MinPts** : minimum number of points in cluster ( $\epsilon$ -neighborhood of that point).

▶ 5

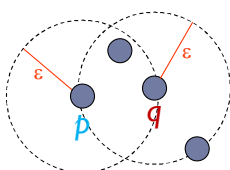
TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

5

## DBSCAN Concepts

- ❖  **$\epsilon$ -neighborhood**: Points within  $\epsilon$  distance (radius) of a point.
- ❖ **Density**: number of points within a specified radius ( $\epsilon$ ).
- ❖ **Core point**:  $\epsilon$ -neighborhood dense enough (*MinPts*),
  - or, If the  $\epsilon$ -neighborhood contains at least a minimum number of points *Minpts*, then the object is a core object.
  - or,  $q$  is core point if  $|N_\epsilon(q)| \geq MinPts$ .
  - These are points that are at the interior of a cluster.

 $\epsilon$ -neighborhood of  $q$  $\epsilon$ -neighborhood of  $p$  $q$  is a core point ( $MinPts = 4$ ) $p$  is not a core point.

▶ 6

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

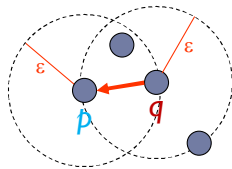
6

## DBSCAN Concepts

❖ **Directly density-reachable**: point  $p$  is directly density-reachable from a point  $q$  w.r.t.  $\epsilon$ ,  $MinPts$  if

1.  $p$  belongs to  $\epsilon$ -neighborhood of  $q$ ,
  - or, the distance is at most  $\epsilon$ ,
  - or,  $p \in ball_\epsilon(q) \equiv p \in N_\epsilon(q)$ .
2.  $q$  is a core point,

DDR is an asymmetric relation!



$MinPts = 4$

$p$  is DDR from  $q$ .

$q$  is not DDR from  $p$ !

► 7

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

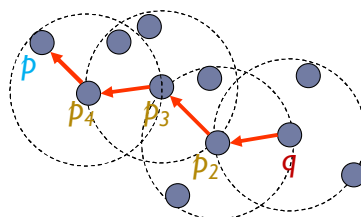
7

## DBSCAN Concepts

❖ **Density-reachable**: A point  $p$  is density-reachable from a point  $q$  w.r.t.  $\epsilon$ ,  $MinPts$  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  $p_i$ .

Or, point  $p$  is density-reachable from  $q$ , if there is a path (chain of points) from  $p$  to  $q$  consisting of only core points.

DR is an asymmetric relation!



$MinPts = 4$

$p$  is DR from  $q$ .

$q$  is not DR from  $p$ !

$p$  is not core.

► 8

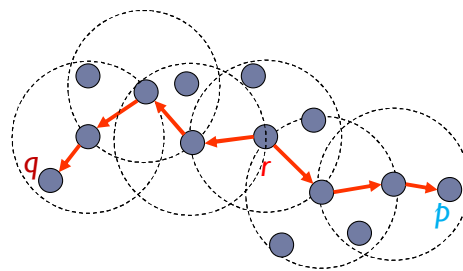
TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

8

## DBSCAN Concepts

- ❖ **Density-connectivity**: point  $p$  is density-connected to point  $q$  w.r.t.  $\epsilon$ ,  $MinPts$  if there is a point  $r$  such that both,  $p$  and  $q$  are **density-reachable** from  $r$  w.r.t.  $\epsilon$  and  $MinPts$ .



DC is an symmetric relation!

$MinPts = 4$

$p$  and  $q$  are density-connected.

► 9

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

9

## DBSCAN Concepts

- ❖ DBSCAN relies on a *density-based* notion of cluster.
- ❖ **Cluster**: a cluster  $C$  is a non-empty set of **density-connected** points that is **maximal** w.r.t. density-reachability.
  - **Maximality**: For all  $p, q$ ; if  $p \in C$  and if  $q$  is density-reachable from  $p$  w.r.t.  $\epsilon$  and  $MinPts$ , then also  $q \in C$ .
- ❖ **Border point**: density-reachable from a core.
  - or, has fewer than  $MinPts$  within  $\epsilon$ -neighborhood, but is in the neighborhood of a core point.

► 10

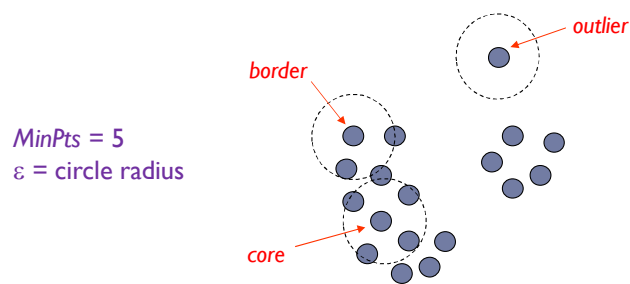
TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

10

## DBSCAN Density Concepts

- ❖ **Noise (outlier) point** : is any point that is not a core point nor a border point.
  - ▶ or, not directly density-reachable from at least one core,
  - ▶ or, Let  $C_1, \dots, C_m$  be the clusters, every point not contained in any cluster is considered to be *noise*.
- ❖ cluster contains *core points* as well as *border points*.



▶ 11

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

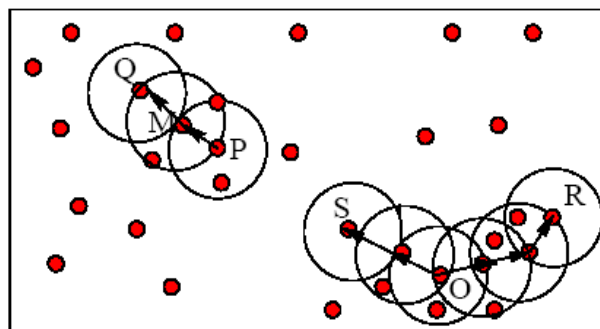
11

## DBSCAN Density Concepts

### Example:

- ❖ Core objects : M, P, O, R  
(each one is in an  $\epsilon$ -neighborhood containing at least 3 points.)

MinPts = 3  
 $\epsilon$  = circle radius



▶ 12

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

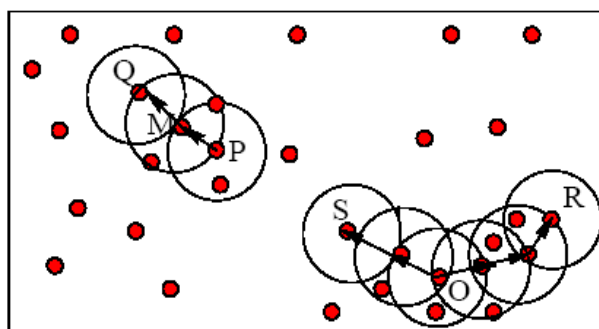
12

## DBSCAN Density Concepts

### Example:

- ❖ Q is directly density reachable from M.
- ❖ M is directly density reachable from P and vice versa.

MinPts = 3  
 $\epsilon$  = circle radius



► 13

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

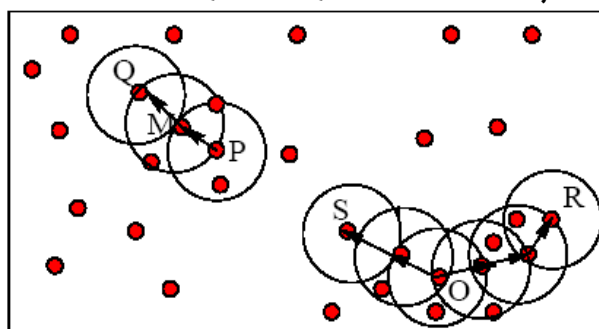
13

## DBSCAN Density Concepts

### Example:

- ❖ Q is (indirectly) density reachable from P since Q is directly density reachable from M and M is directly density reachable from P.
- ❖ P is not density reachable from Q since Q is not a core object.

MinPts = 3  
 $\epsilon$  = circle radius



► 14

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

14

## DBSCAN – thoughts behind the algorithm

### Note 1:

- ❖ If  $q$  is a core point and  $D_q$  is the set of points that are density reachable from  $q$ , then  $D_q$  is a cluster.

### Note 2:

- ❖ If  $C$  is a cluster and  $q$  is a core point in  $C$ , then  $C$  equals to the set of the points  $p$  that are *density reachable* from  $q$ .

⇒ A cluster is uniquely defined by any of its core points.

▶ 15

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

15

## DBSCAN algorithm

- ❖ DBSCAN searches for clusters by checking the  $\varepsilon$ -neighborhood of each point in the database. If the  $\varepsilon$ -neighborhood of a point  $p$  contains more than  $MinPts$ , a new cluster with  $p$  as a **core** point created.
- ❖ DBSCAN then iteratively collects **directly density-reachable** objects from these core objects, which may involve the merge of a few density-reachable clusters.
- ❖ The process terminates when no new point can be added to any cluster.

▶ 16

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

16



## DBSCAN Algorithm

- ❖ Arbitrary select a point  $p$
- ❖ Retrieve all points density-reachable from  $p$  w.r.t.  $\varepsilon$  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.

▶ 17

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

17

## DBSCAN Algorithm

```

 $D_{unprocessed} \leftarrow D$  //  $D$ : data-set
 $no\_of\_clusters \leftarrow 0$ 
while  $D_{unprocessed} \neq \emptyset$  do
    Arbitrarily select a  $p \in D_{unprocessed}$ 
    if  $p$  is a non-core point, then
        ▪ Mark  $p$  as noise point
        ▪  $D_{unprocessed} \leftarrow D_{unprocessed} - \{p\}$ 
    else //  $p$  is a core point
        ▪  $no\_of\_clusters++$ 
        ▪  $D_{DR}(p) \leftarrow$  Determine all Density-Reachable points in  $D$  from  $p$ 
          //note: The border points that may have been marked as noise, now belong to  $D_{DR}(p)$ .
        ▪  $Cluster_{no\_of\_clusters} \leftarrow \{p\} + D_{DR}(p)$ 
        ▪  $D_{unprocessed} \leftarrow D_{unprocessed} - Cluster_{no\_of\_clusters}$ 
    end-if
end-while

```

▶ 18

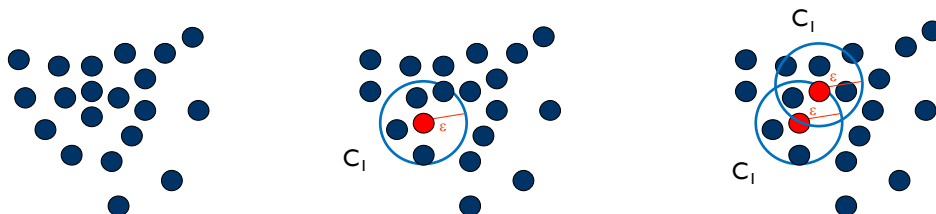
TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

18

## DBSCAN

MinPts = 4



► 19

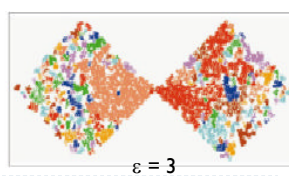
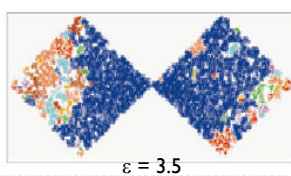
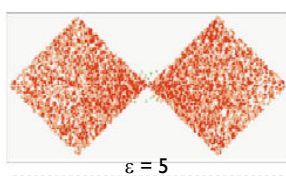
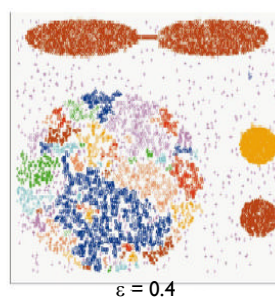
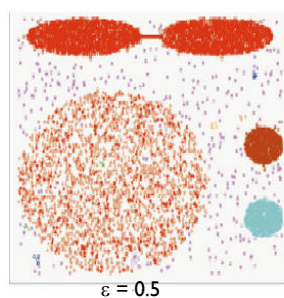
TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

19

## DBSCAN

**DBSCAN is Sensitive to Parameters.** MinPts = 4



► 20

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

20

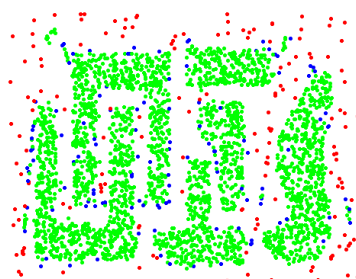
## DBSCAN

❖ Core, Border and Noise Points:

$$MinPts = 4, \quad \varepsilon = 10$$



Original Points



Point types: core, border  
and noise

► 21

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

21

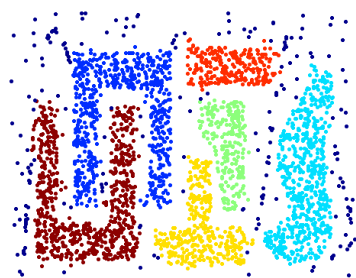
## DBSCAN

**When DBSCAN works well:**

- ❖ Resistant to Noise
- ❖ Can handle clusters of different shapes and sizes



Original Points



Clusters

► 22

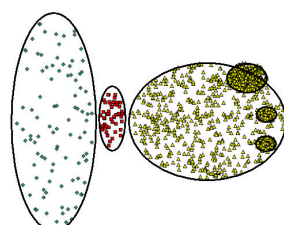
TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

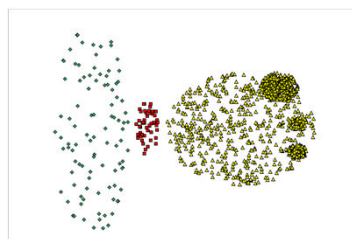
22

## DBSCAN

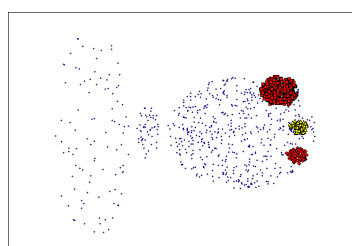
When DBSCAN does not work well:



Original Points



$MinPts = 4,$   
 $\epsilon = 9.75$



$MinPts = 4,$   
 $\epsilon = ?$

- Varying densities
- High-dimensional data

► 23

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

23

## DBSCAN: Determining $\epsilon$ and $MinPts$

- ❖ **Idea** : for points in a cluster, their  $k$ -th nearest neighbors are at roughly the same distance.
- ❖ Noise points have the  $k$ -th nearest neighbor at farther distance.
- ❖ So, plot sorted distance of every point to its  $k$ -th nearest neighbor.
- ❖  $k$ -dist : Distance from a point to its  $k$ -th nearest neighbor.
  - For points that belong to some clusters, the value of  $k$ -dist will be small if  $k$  is not larger than cluster size.
  - For points that are not in a cluster such as noise points, the  $k$ -dist will be relatively large.

► 24

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

24

## DBSCAN: Determining $\epsilon$ and $MinPts$

### Method:

- ❖ Compute  $k$ -dist for all points for some  $k$ .
- ❖ Sort them in increasing order and plot sorted values.
- ❖ A sharp change at the value of  $k$ -dist that corresponds to suitable value of  $\epsilon$  and the value of  $k$  as  $MinPts$ .
  - ▶ Points for which  $k$ -dist is less than  $\epsilon$  will be labeled as **core** points while other points will be labeled as **noise** or **border** points.

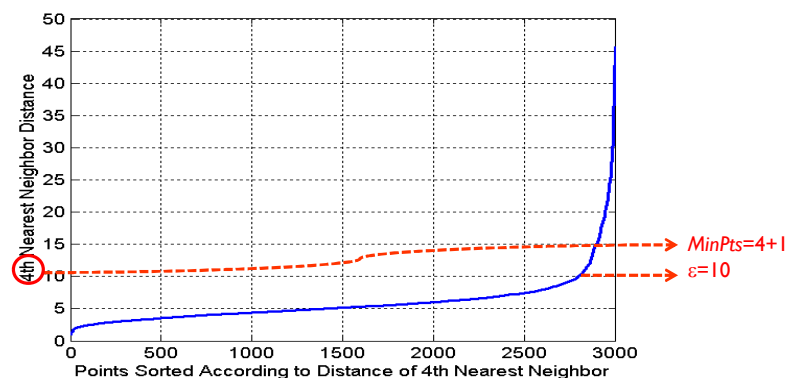
▶ 25

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

25

## DBSCAN: Determining $\epsilon$ and $MinPts$



▶ 26

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

26

## DBSCAN: Determining $\varepsilon$ and *MinPts*

- ❖ If  $k$  is too small  $\Rightarrow$  Even a small number of closely spaced points that are noise or outliers will be incorrectly labeled as clusters
- ❖ If  $k$  is too large  $\Rightarrow$  small clusters (of size less than  $k$ ) are likely to be labeled as noise
  
- ❖ If a spatial index (ex,  $kd$ -tree,  $R^*$ -tree) is used, the computational complexity of DBSCAN is  $O(n \log n)$ , where  $n$  is the number of database objects. Otherwise, it is  $O(n^2)$ .

► 27

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

27

## Conclusion

- ❖ DBSCAN is sensitive to the choice of input parameters.
- ❖ DBSCAN is a density-based clustering algorithm which is able to discover clusters of arbitrary shapes. Distance is not the metric unlike the case of hierarchical methods.
- ❖ Parameter setting is done empirically:
  - ❖ They should be selected such that it would be able to detect the least “dense” cluster, thus experimentation with several values for parameters should be done.
- ❖ High dimensional data – more noticeable (pronounced).
- ❖ High dimensional data clustering structures are not generally characterized by global density parameters like  $\varepsilon$  & *MinPts*,
- ❖ OPTICS as a solution!

► 28

TMU, M.M. Pedram, [pedram@tmu.ac.ir](mailto:pedram@tmu.ac.ir)

Data Mining, Spring 2011

28