# Cluster Analysis – Density based Algorithms
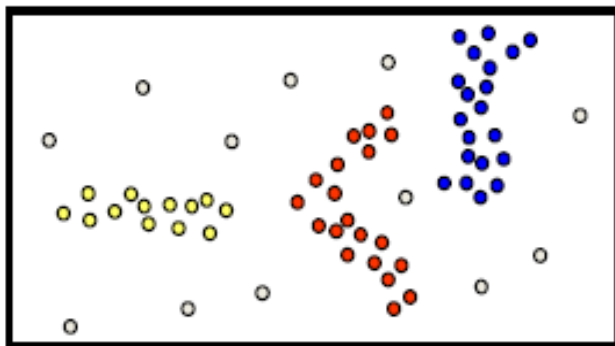
# Density-based Clustering

- **Basic idea**

  – Clusters are dense regions in the data space, separated by regions of lower object density

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

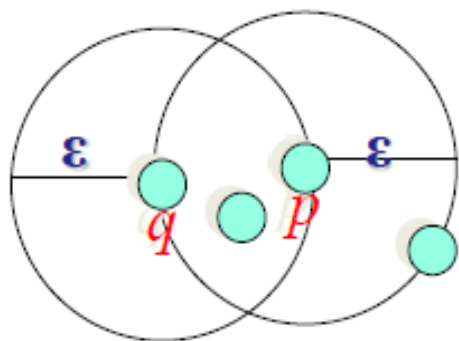  – Discovers clusters of arbitrary shape

- **Method**

  – DBSCAN

# Density Definition

- ε-Neighborhood – Objects within a radius of ε from an object.

$$N_\varepsilon(p) : \{q \mid d(p,q) \le \varepsilon\}$$

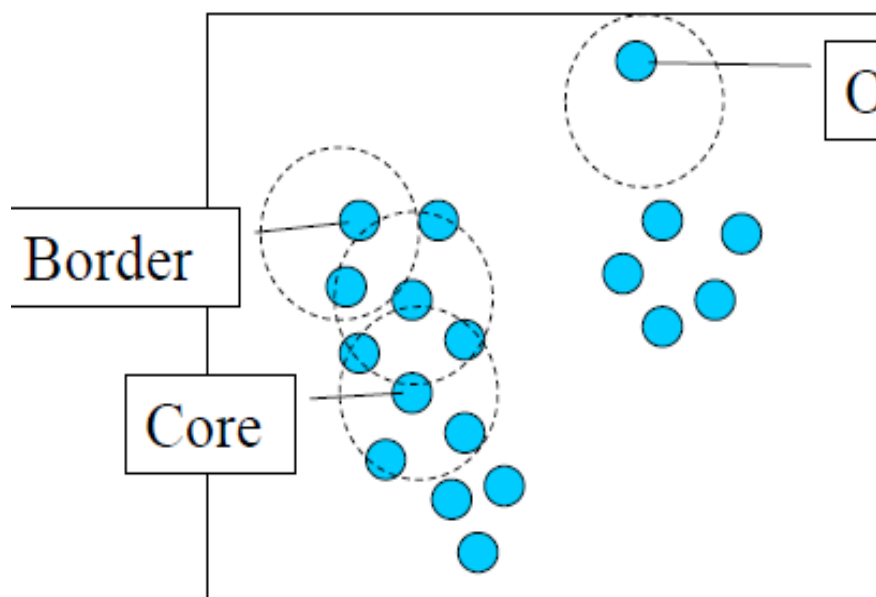- "High density" - ε-Neighborhood of an object contains at least *MinPts* of objects.



ε-Neighborhood of $p$

ε-Neighborhood of $q$

*Density of $p$ is "high" (MinPts = 4)*

*Density of $q$ is "low" (MinPts = 3 )*

# Core, Border & Outlier



Outlier

Border

Core

$\varepsilon = 1\text{unit}, \text{MinPts} = 5$

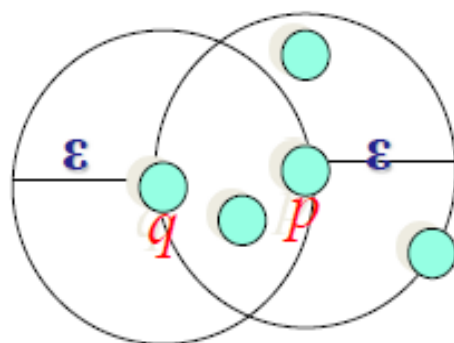Given $\varepsilon$ and *MinPts*, categorize the objects into three exclusive groups.

A point is a core point if it has more than a specified number of points (MinPts) within Eps—These are points that are at the interior of a cluster.

A border point has fewer than MinPts within Eps, but is in the neighborhood of a core point.

A noise point is any point that is not a core point nor a border point.
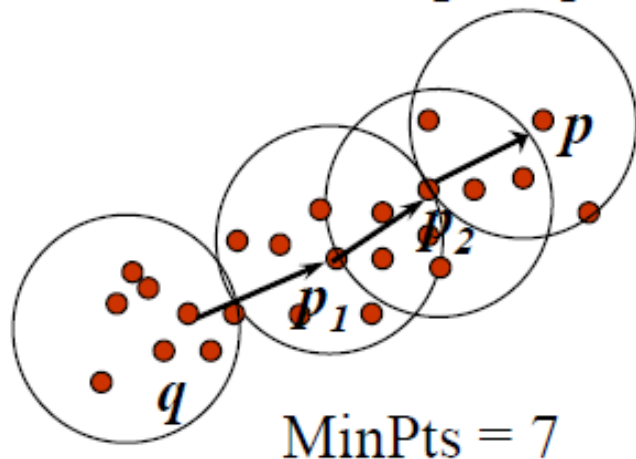
# Density-reachability

- Directly density-reachable

  - An object $q$ is directly density-reachable from object $p$ if $p$ is a core object and $q$ is in p's ε-neighborhood.



MinPts = 4

- $q$ is directly density-reachable from $p$
- $p$ is not directly density-reachable from $q$
- Density-reachability is asymmetric

# Density-reachability

- Density-Reachable (directly and indirectly):

  - A point $p$ is directly density-reachable from $p_2$

  - $p_2$ is directly density-reachable from $p_1$

  - $p_1$ is directly density-reachable from $q$

  - $p \leftarrow p_2 \leftarrow p_1 \leftarrow q$ form a chain



MinPts = 7

- $p$ is (indirectly) density-reachable from $q$

- $q$ is not density-reachable from $p$
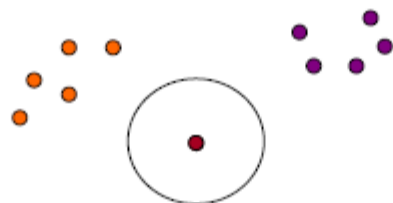
# DBSCAN Algorithm: Example

- **Parameter**

  - $\varepsilon$ = 2 cm

  - *MinPts* = 3



```
for each o ∈ D do
    if o is not yet classified then
        if o is a core-object then
            collect all objects density-reachable from o
            and assign them to a new cluster.
        else
            assign o to NOISE
```
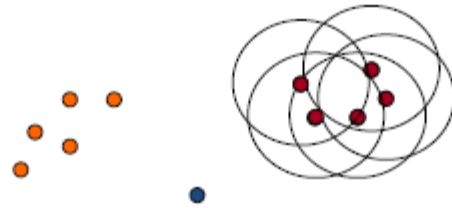
# DBSCAN Algorithm: Example

- **Parameter**
  - $\varepsilon$ = 2 cm
  - *MinPts* = 3

```
for each o ∈ D do
    if o is not yet classified then
        if o is a core-object then
            collect all objects density-reachable from o
            and assign them to a new cluster.
        else
            assign o to NOISE
```

# DBSCAN Algorithm: Example

- **Parameter**

  - $\varepsilon$ = 2 cm

  - *MinPts* = 3



```
for each o ∈ D do
    if o is not yet classified then
        if o is a core-object then
            collect all objects density-reachable from o
            and assign them to a new cluster.
        else
            assign o to NOISE
```

# DBSCAN: Sensitive to Parameters

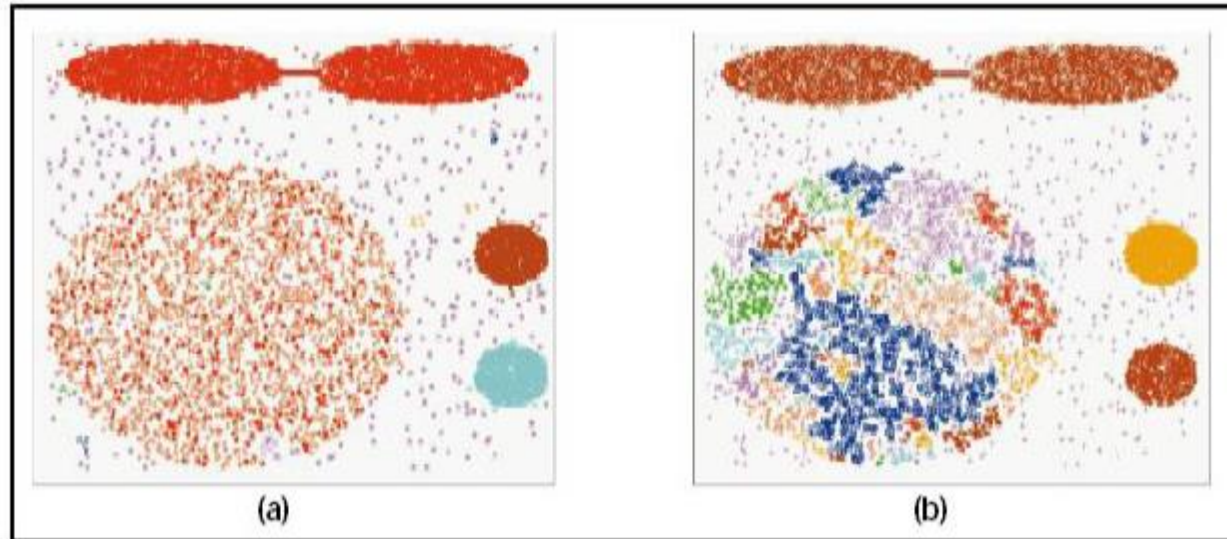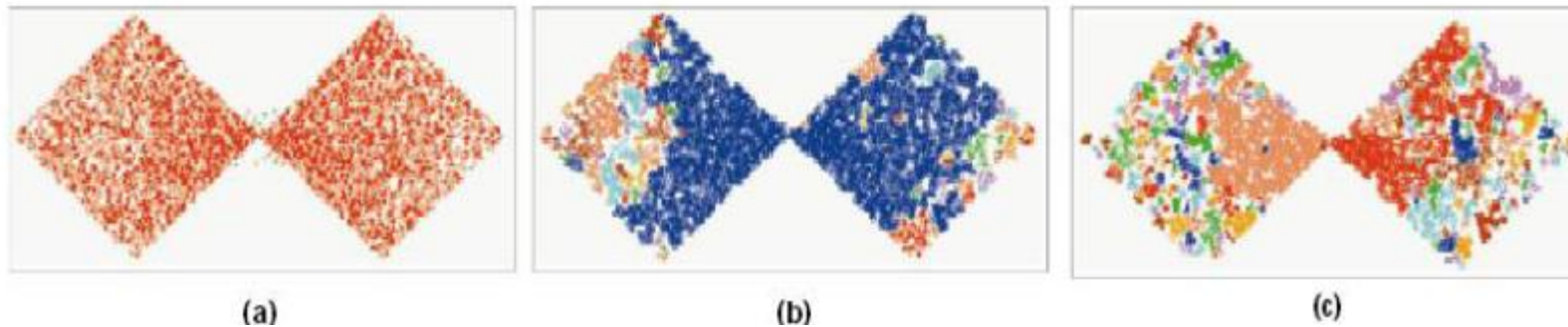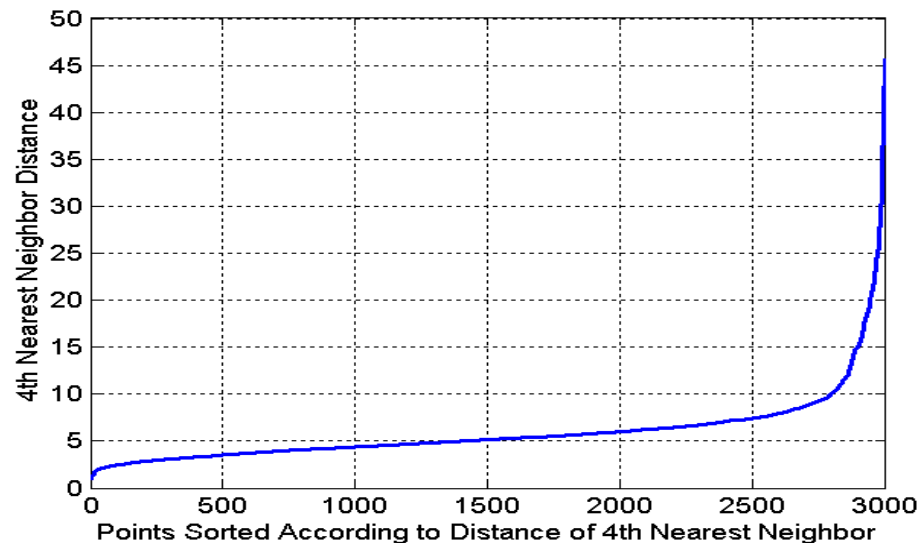Figure 8. DBScan results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.

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



(a)

(b)

(a)

(b)
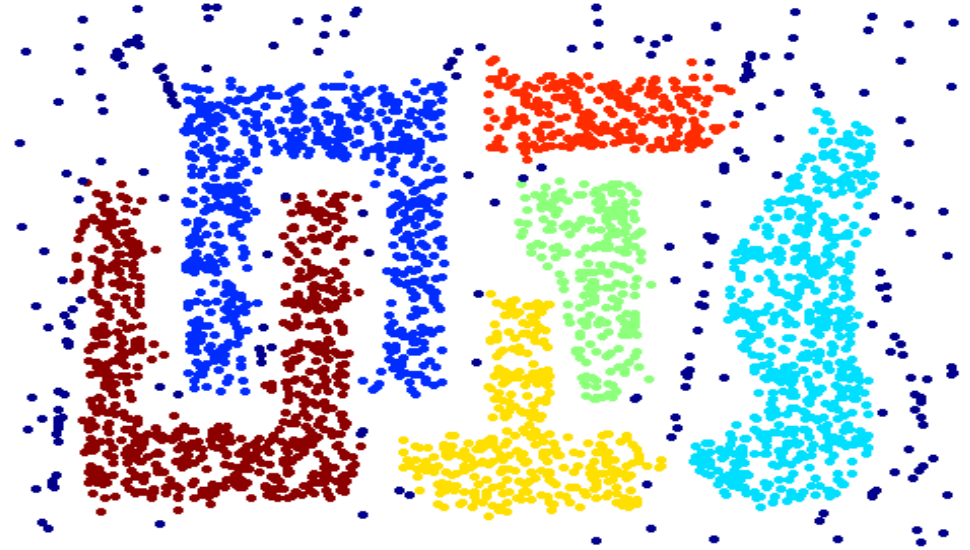
(c)

# DBSCAN: determining EPS and MinPts

- Idea is that 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
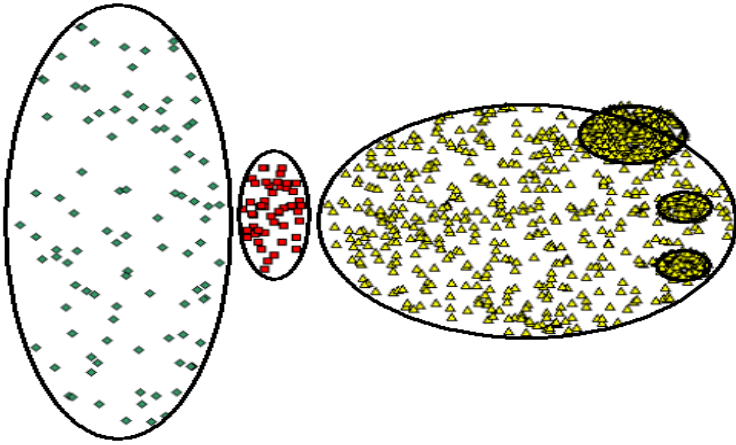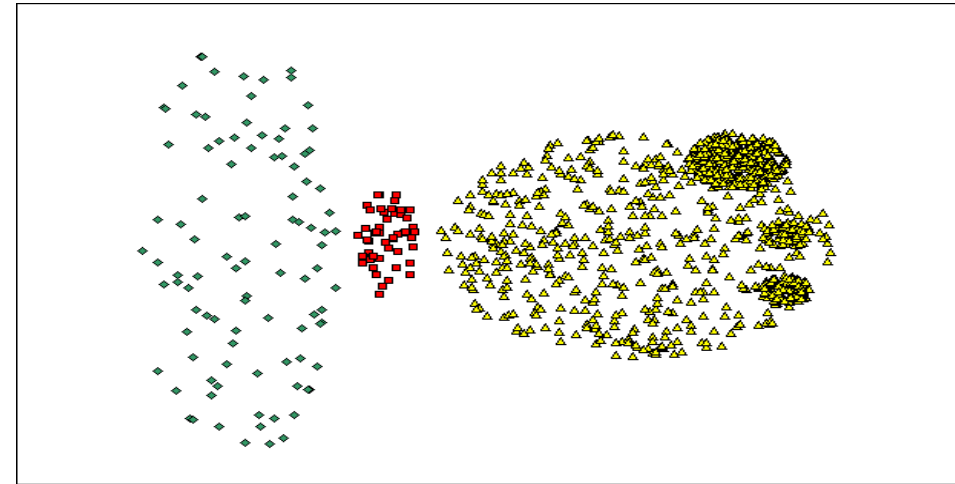
# When DBSCAN works well



**original points**

**clusters**

- resistant to noise
- can handle clusters of different shapes and sizes
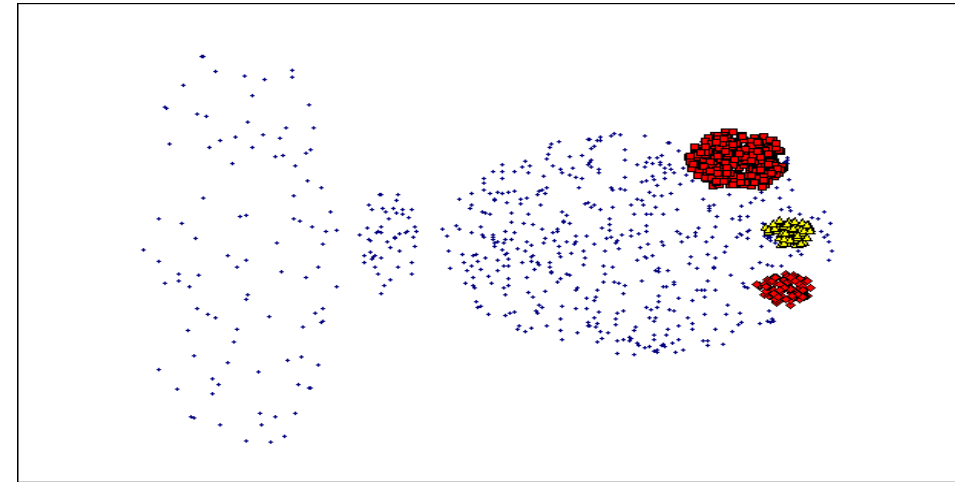
# When DBSCAN does NOT work well



(MinPts=4, Eps=9.75).



(MinPts=4, Eps=9.92)

**original points**

- varying densities
- high-dimensional data