# Clustering II: Spatial Clustering

# Roadmap

o Density Based Clustering

o DBSCAN

   o Concepts

   o Algorithm

   o Comments

o Take-home messages

# Density-based Approaches*

- Why Density-Based Clustering methods?
  - Discover clusters of arbitrary shape
  - Clusters – Dense regions of objects separated by regions of low density

- DBSCAN – the first density based clustering

- Other methods:
  - OPTICS – density based cluster-ordering
  - DENCLUE – a general density-based description of cluster and clustering

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

- Proposed by Ester, Kriegel, Sander, and Xu (KDD96)

- Relies on a density-based notion of cluster: A cluster is defined as a maximal set of density-connected points.

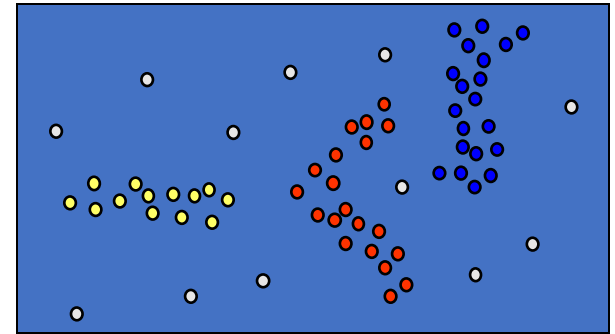- Discovers clusters of arbitrary shape in spatial databases with noise

Visualization tool:
https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/
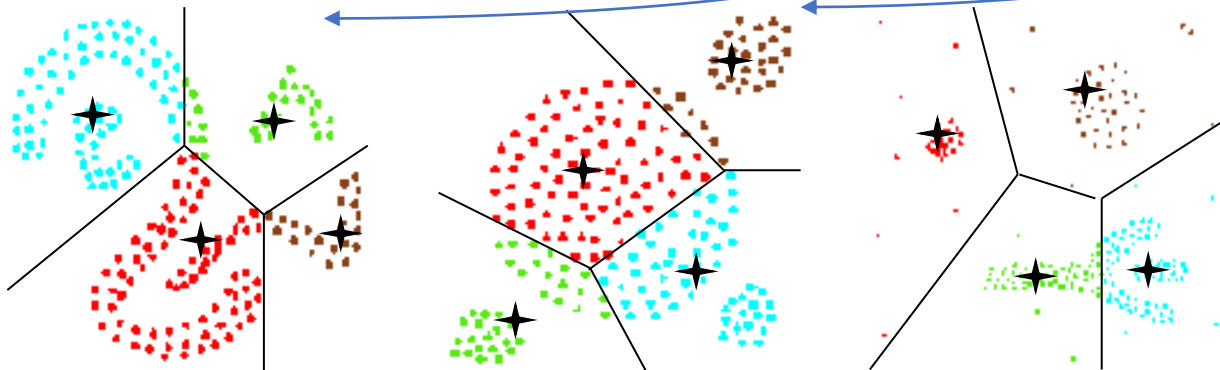
# Density-Based Clustering

✳ *Basic Idea*:

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



• Why Density-Based Clustering?

Results of a *k*-medoid algorithm for *k*=4
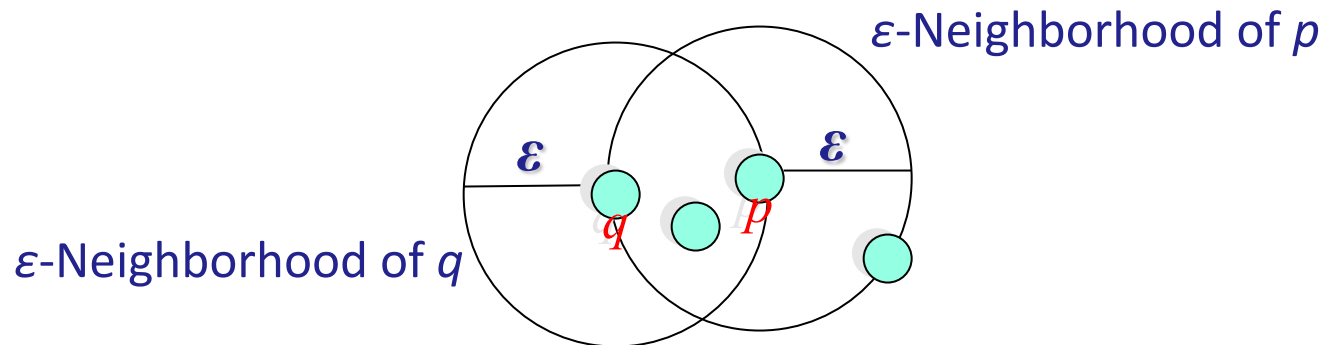
*Are these reasonable?*



Different density-based approaches exist (see Textbook & Papers)
Here, we discuss the ideas underlying the DBSCAN algorithm

# Density Based Clustering: Basic Concept

- Intuition for the formalization of the basic idea
  - For any point in a cluster, the local point density around that point has to exceed some threshold
  - The set of points from one cluster is spatially connected

- Local point density at a point *p* defined by two parameters
  - $\varepsilon$ – radius for the neighborhood of point p:
    $N_{\varepsilon}(p) := \{q$ in data set $D \mid dist(p, q) \leq \varepsilon\}$
  - *MinPts* – minimum number of points in the given neighbourhood $N(p)$
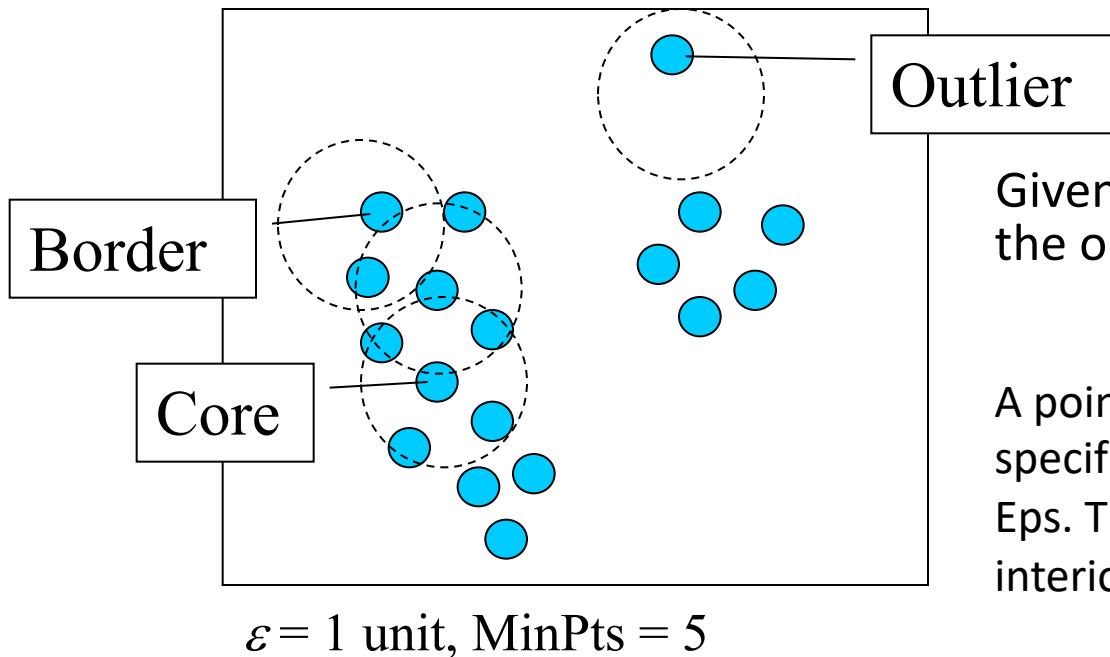
# $\varepsilon$-Neighborhood

- $\varepsilon$-Neighborhood – Objects within a radius of $\varepsilon$ from an object.

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

- "High density" – $\varepsilon$-Neighborhood of an object contains at least *MinPts* of objects.



$\varepsilon$-Neighborhood of p

$\varepsilon$

$\varepsilon$

$q$

$p$

$\varepsilon$-Neighborhood of q

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

*Density of q* is "low" (MinPts = 4)

# Core, Border & Outlier



Outlier

Border

Core

$\varepsilon = 1$ unit, MinPts = 5

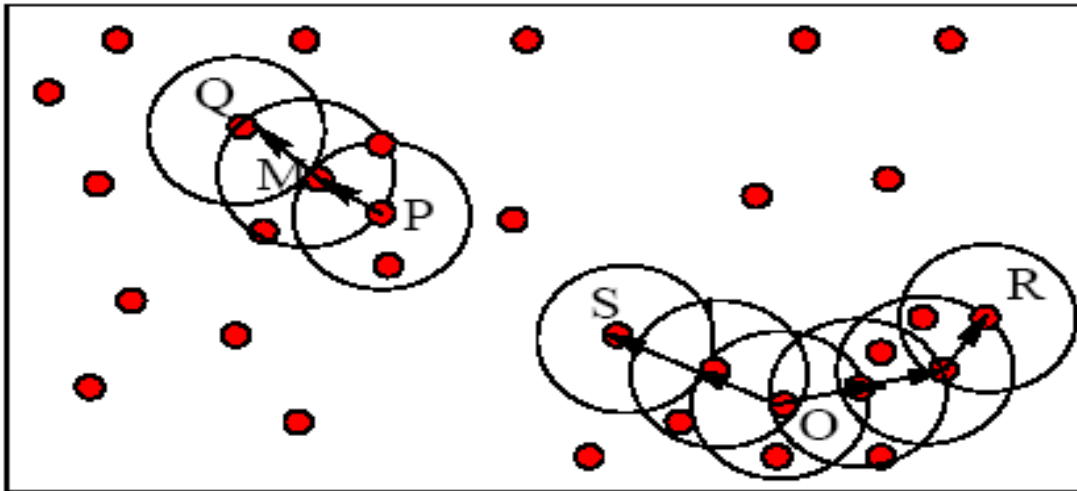Given $\varepsilon$ and *MinPts*, we can categorize the objects into three exclusive groups.

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/outlier is any point that is not a core point nor a border point.

# Example

- M, P, O, and R are core objects (out of Q, M, P, S, O, R) since each of them is in an $\varepsilon$-neighborhood containing at least 3 points
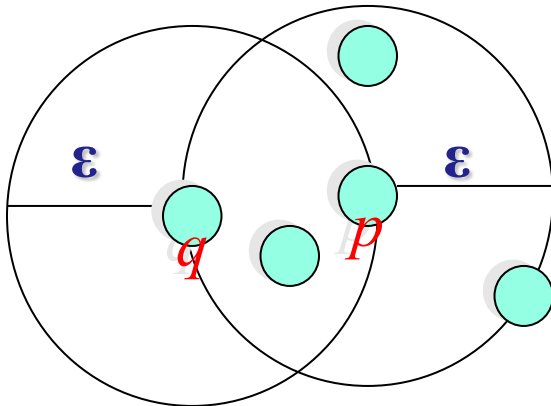


MinPts = 3

$\varepsilon$ = radius of the circles

# 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 $\varepsilon$-neighborhood.
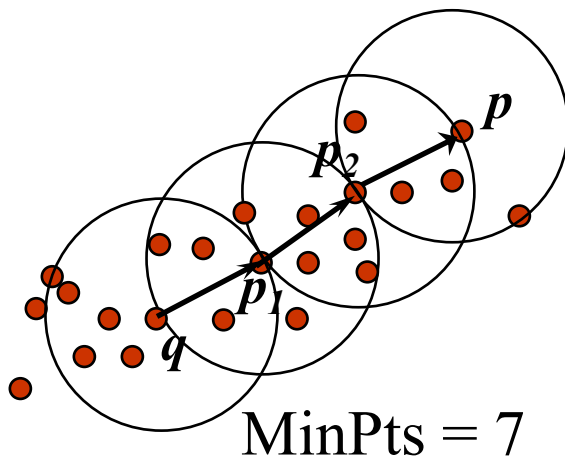


ε    ε

*q*    *p*

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 $q$

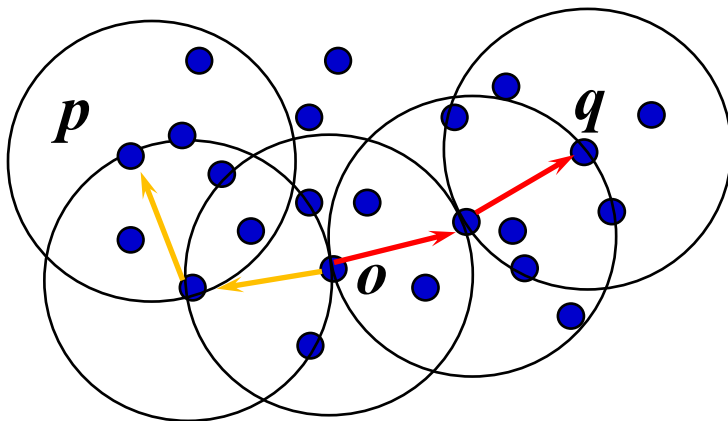■ Is $q$ not density-reachable from $p$? Yes, check whether $p_2$ is directly density-reachable from $p$

# Density-Reachability ➜ Density-Connectivity

■ **Density-Reachable is not symmetric**

  ❑ not good enough to describe clusters

■ **Density-Connected**

  ❑ A pair of points $p$ and $q$ are density-connected if they are commonly density-reachable from a point $o$.
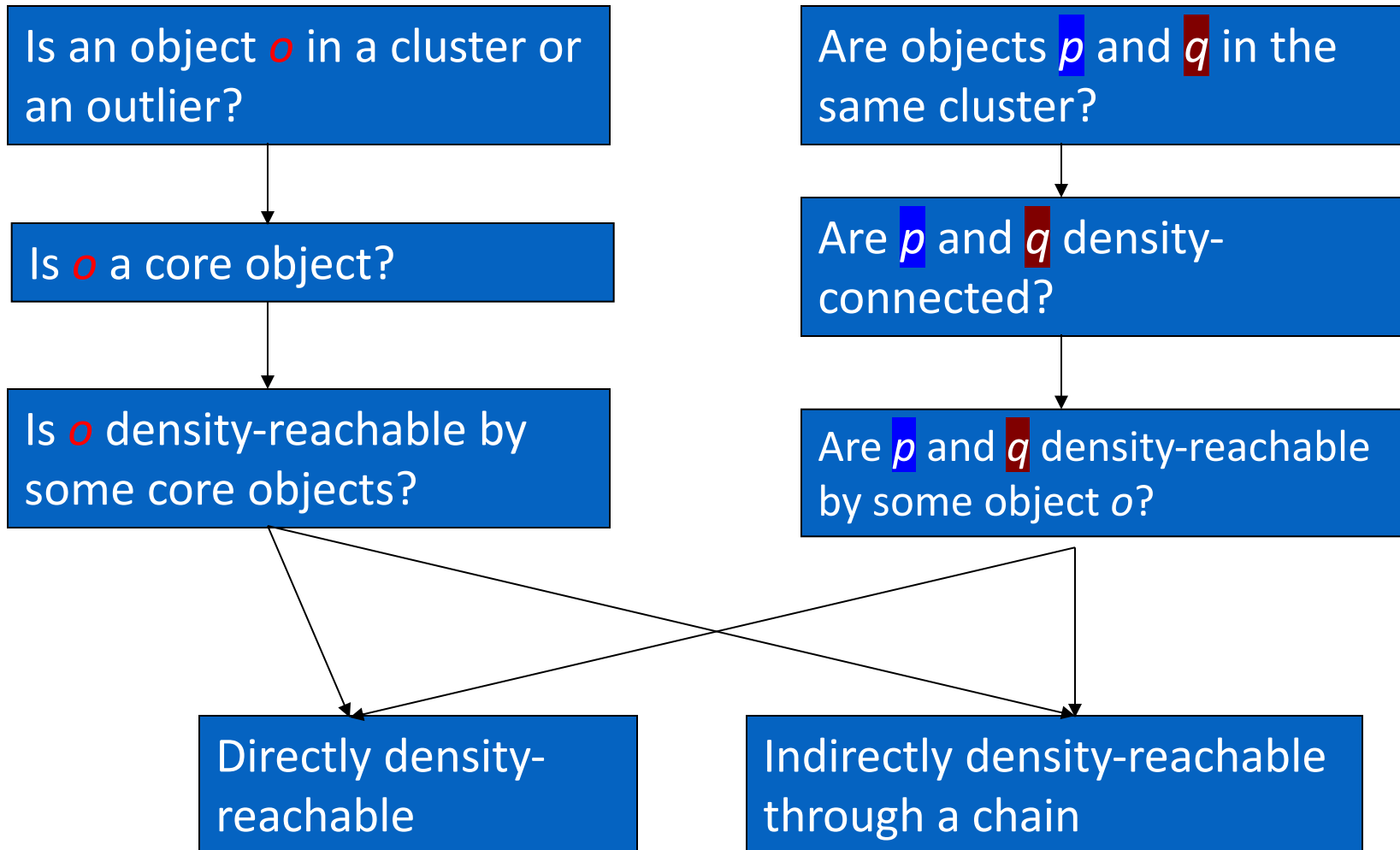


■ Density-connectivity is symmetric

# Formal Description of Cluster

- Given a data set $D$, parameter $\varepsilon$ and threshold MinPts.

- A cluster $C$ is a subset of objects (with core and border points) satisfying two criteria:
  - *Connected:* $\forall\ p,q \in C$: $p$ and $q$ are density-connected.
  - *Maximal:* $\forall\ p,q$: if $p \in C$ and if $q$ is <u>density-reachable from $p$</u>, then $q \in C$.

    *Density-reachable ---> p is a core point*
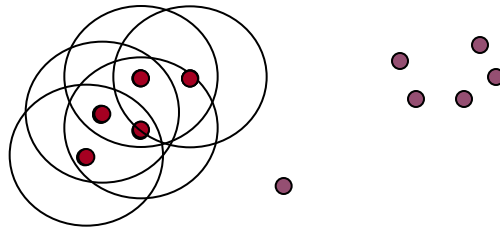
# Review of Concepts

Is an object *o* in a cluster or an outlier?

Is *o* a core object?

Is *o* density-reachable by some core objects?

Are objects *p* and *q* in the same cluster?

Are *p* and *q* density-connected?

Are *p* and *q* density-reachable by some object *o*?

Directly density-reachable

Indirectly density-reachable through a chain

# DBSCAN: The Algorithm

1. Arbitrary select a point $p$

2. Retrieve all points density-reachable from $p$ wrt *Eps* ($\varepsilon$) and *MinPts*.

3. If $p$ is a core point, a cluster is formed.

4. If $p$ is a border point, no points are density-reachable from $p$ and DBSCAN visits the next point of the database.

5. Continue the process until all of the points have been processed.

# DBSCAN Algorithm: Example
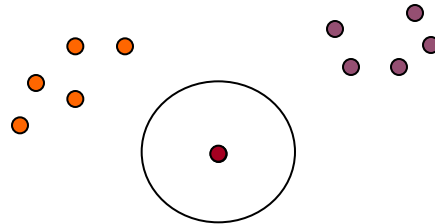
- Parameter
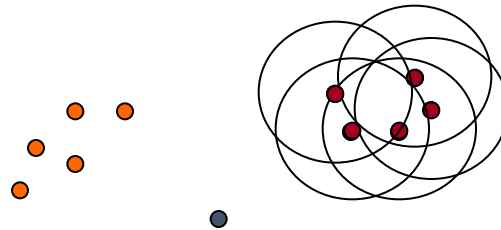  - $\varepsilon$ = 2 cm
  - *MinPts* = 3



```
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
```
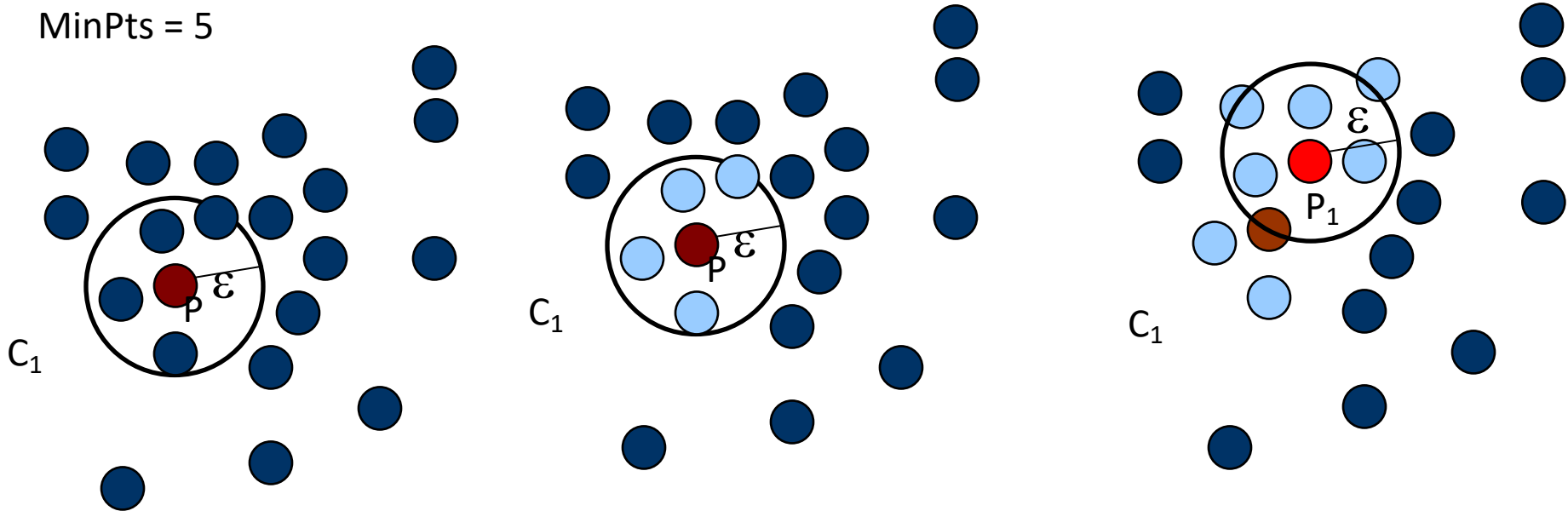
# DBSCAN Algorithm: Example

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



**for** each $o \in D$ **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
```
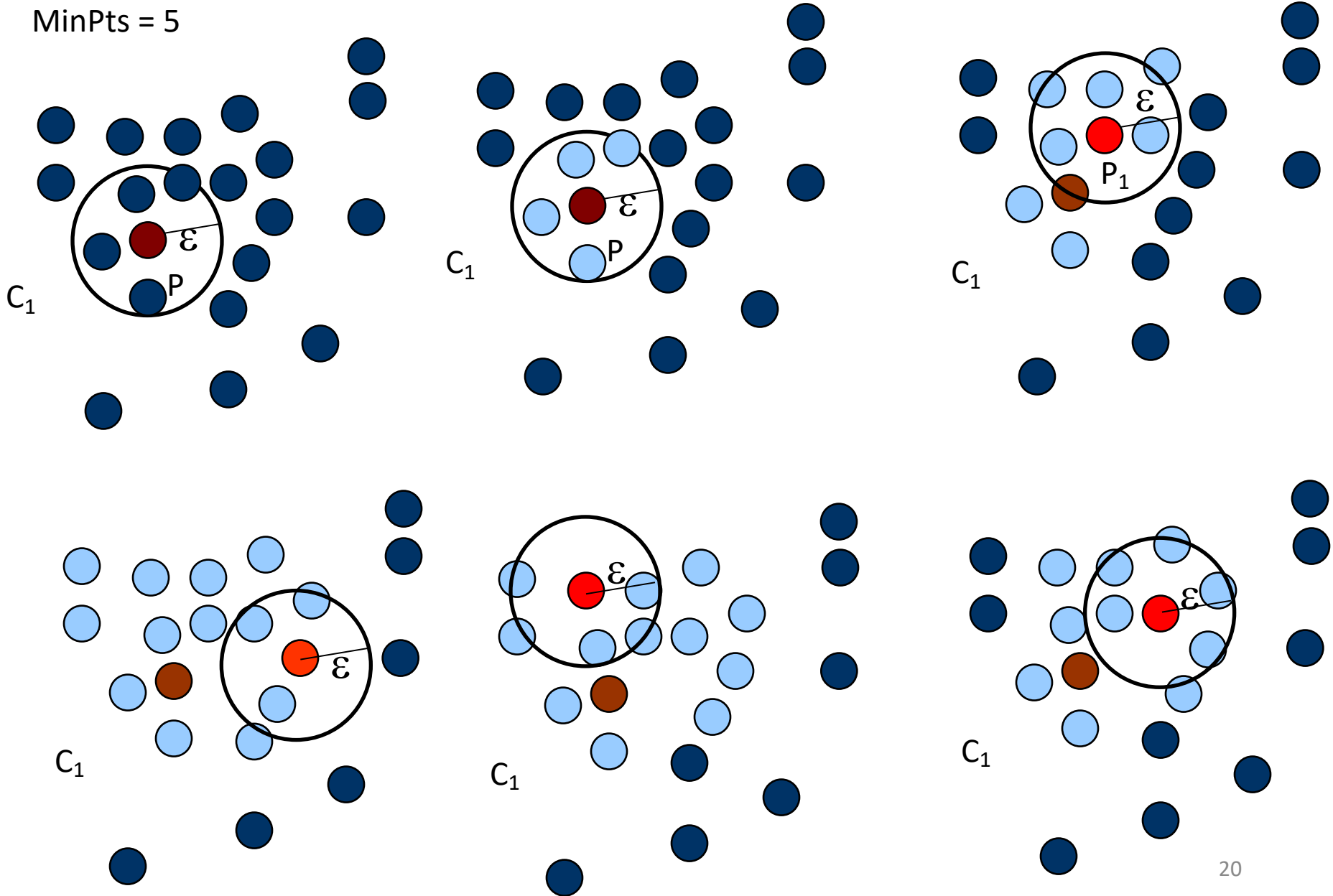
# More elaborations

MinPts = 5

$C_1$

$C_1$

$C_1$

$\varepsilon$

$\varepsilon$

$\varepsilon$

P

P

$P_1$

1. Check the $\varepsilon$-neighborhood of p;

2. If p has less than MinPts neighbors then mark p as outlier and continue with the next object

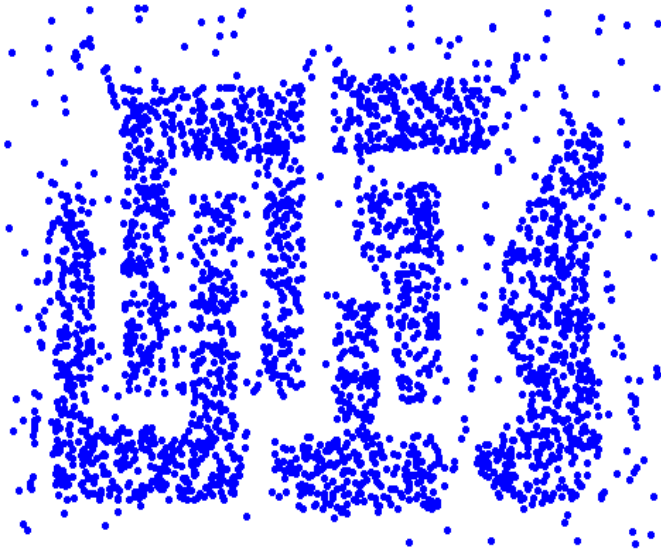3. Otherwise mark p as processed and put all the neighbors in cluster C

1. Check the unprocessed objects in C (light blue dots)

2. If no core object, return C

3. Otherwise, randomly pick up one core object $p_1$, mark $p_1$ as processed, and put all unprocessed neighbors of $p_1$ in cluster C
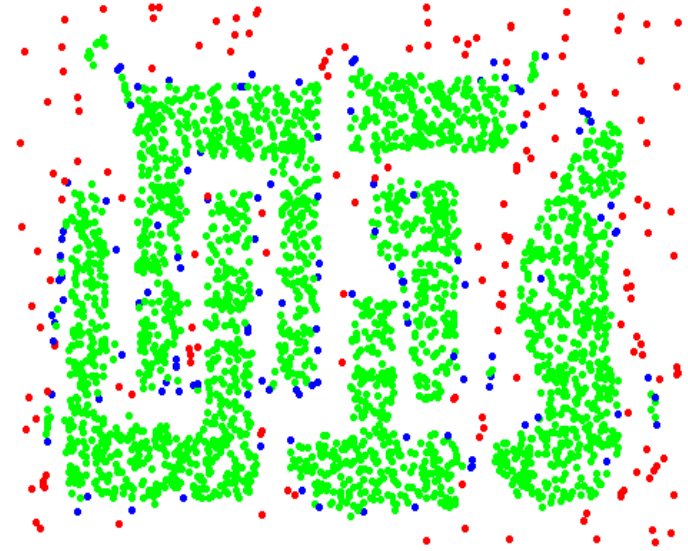
MinPts = 5

$C_1$

$P$

$\mathcal{E}$

$C_1$

$P$

$\mathcal{E}$

$C_1$

$P_1$

$\mathcal{E}$

$C_1$

$\mathcal{E}$

$C_1$

$\mathcal{E}$

$C_1$
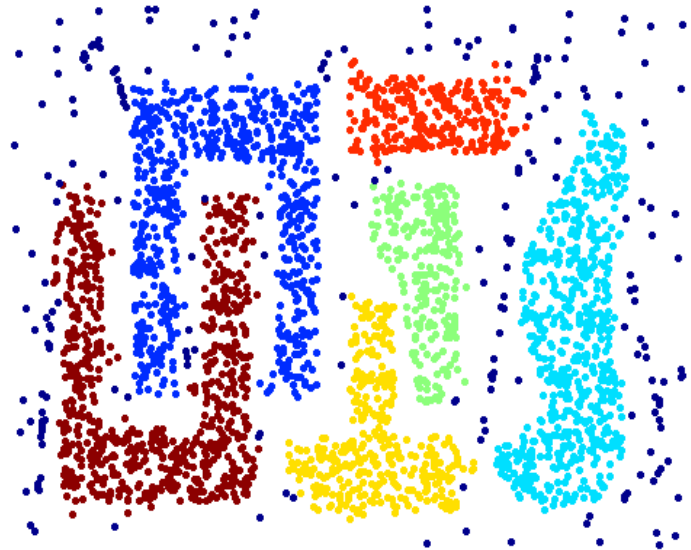
$\mathcal{E}$

# Pictorial Example



**Original Points**

**Point types: core, border and outliers**

ε = 10, MinPts = 4
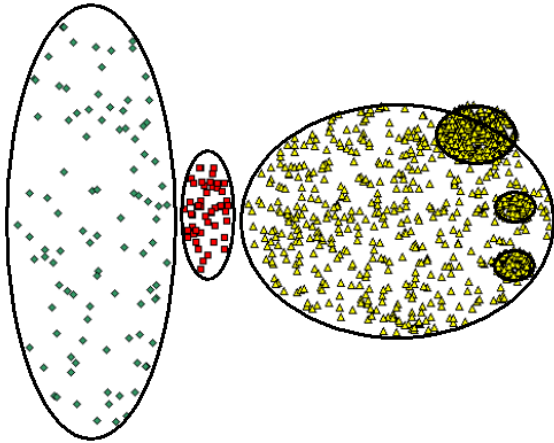
# When DBSCAN Works Well…

**Original Points**

**Clusters**
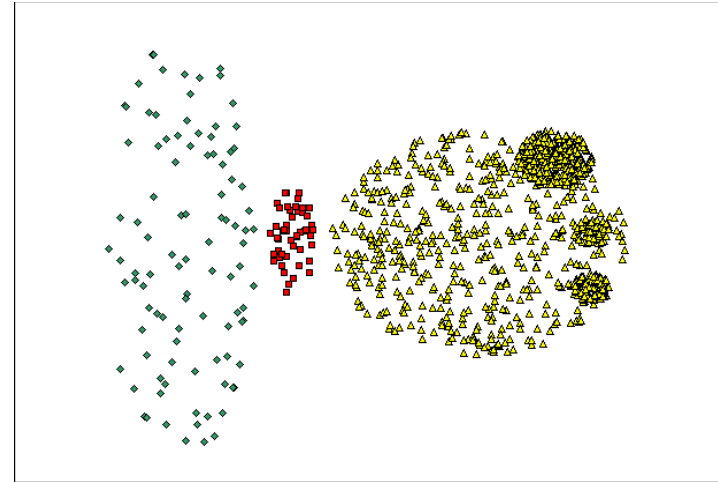
- **Resistant to Noise**

- **Can handle clusters of different shapes and sizes**
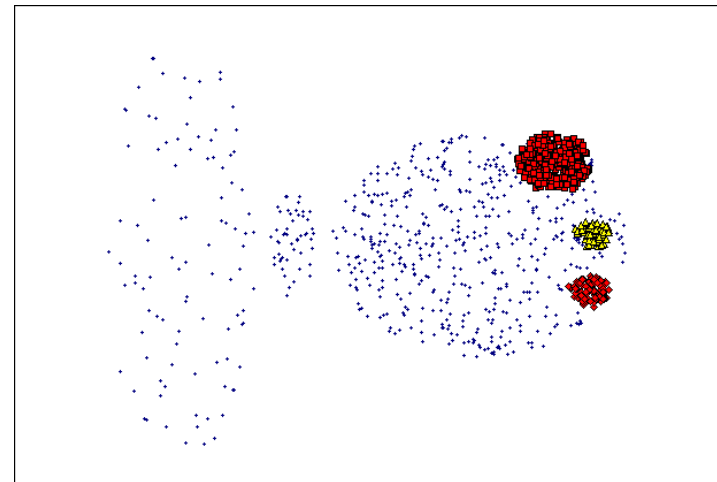
# When DBSCAN Does NOT Work Well…



(MinPts=4, Eps=9.92).



(MinPts=4, Eps=9.75)

**Original Points**

- **Cannot handle varying densities!**

- **Sensitive to parameters!**

# DBSCAN: Sensitive to Parameters

Figure 8. DBScan results for DS1 with MinPts at 4 and Eps at (a) 0.5 and (b) 0.4.
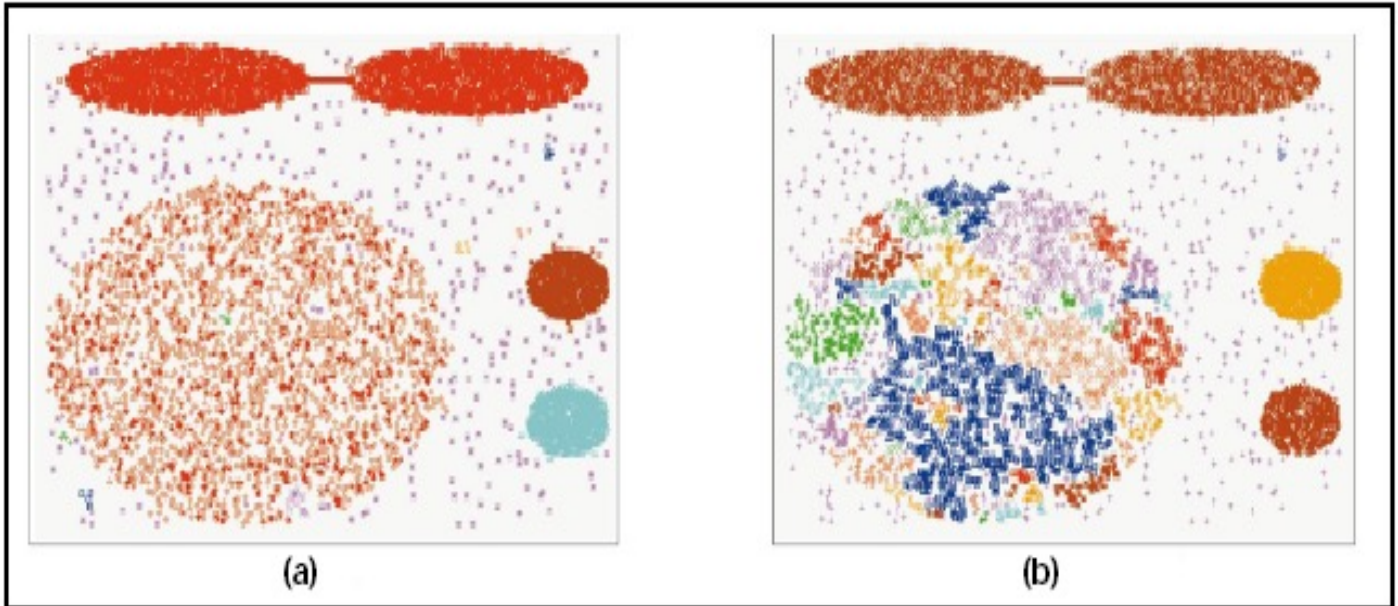
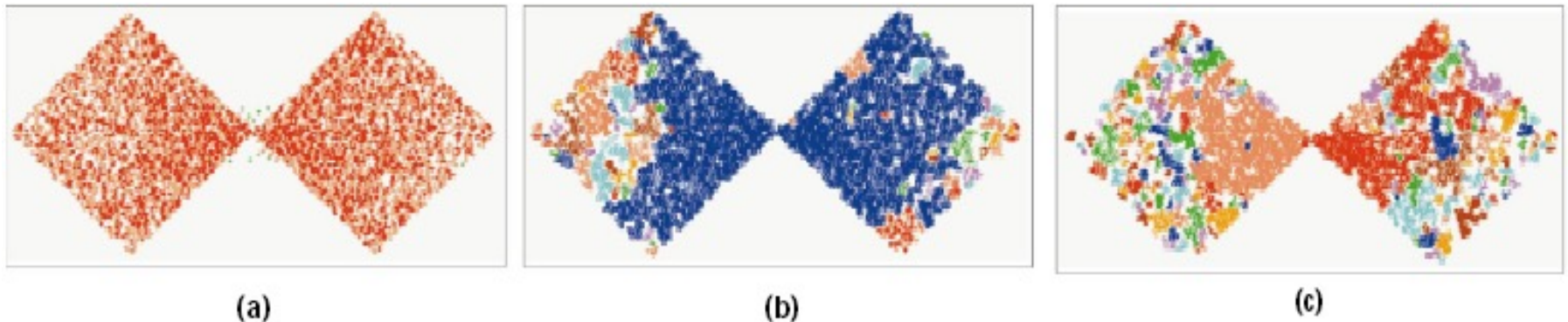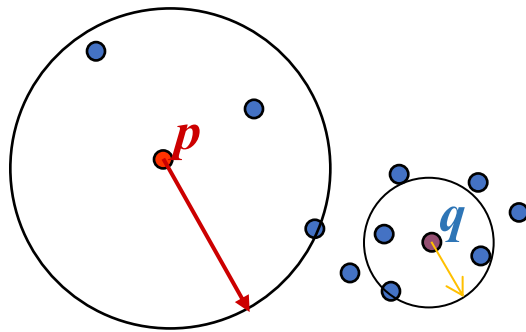Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.



(a)

(b)

(a)

(b)

(c)

# Determining the Parameters $\varepsilon$ and *MinPts*

- Cluster: Point density higher than specified by $\varepsilon$ and *MinPts*

- Idea: Use the point density of the least dense cluster in the data set as parameters – but how to determine this?

- Heuristic: look at the distances to the *k*-nearest neighbors
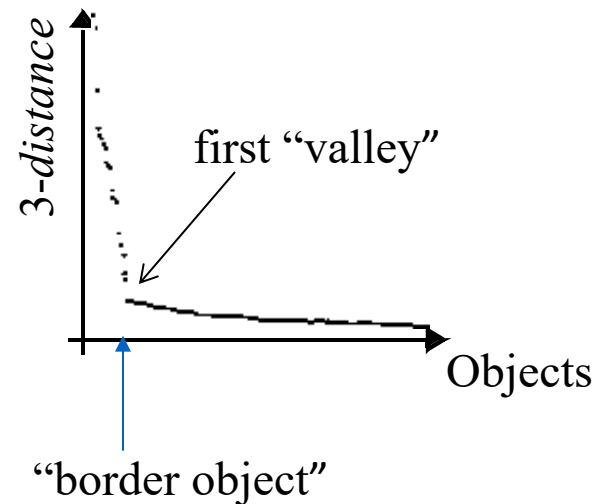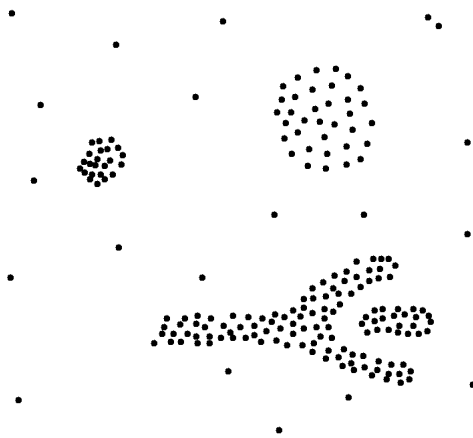


$3\text{-}distance(p):$      ⟶

$3\text{-}distance(q):$      →

- Function *k-distance*(*p*): distance from *p* to the its *k*-nearest neighbor

- *k-distance plot*: *k*-distances of ALL objects, sorted in decreasing order

# Determining the Parameters $\varepsilon$ and *MinPts*

- Example *k*-distance plot



- Heuristic method:
  - Fix a value for *MinPts* (default: $2 \times d - 1$, where d denotes dimensionality of data space)
  - User selects "border object" *o* from the *MinPts-distance* plot; $\varepsilon$ is set to *MinPts-distance*(*o*)

# Density Based Clustering: Discussion

- Advantages
  - Clusters can have arbitrary shape and size
  - Number of clusters is determined automatically
  - Can separate clusters from surrounding noise

- Disadvantages
  - Input parameters may be difficult to determine
  - In some situations very sensitive to input parameter setting

# Take-home Messages

- Cluster analysis groups objects based on their similarity and has wide applications

- Density-based clustering takes into considerations of other important concepts, i.e., (density and connectivity) vs (intra-cluster and inter-cluster similarity).

- There are still lots of research issues on cluster analysis, such as semi-supervised clustering, subspace clustering, etc.

- Yet it is always a topic of interest for emerging applications
  - Clustering in a social network graph
  - Spatial clustering of GPS data
  - Spatial clustering of farming data →



An Introduction to
**PRECISION FARMING**

# Acknowledgement

- Slides/Materials of
    - https://www.cse.buffalo.edu/faculty/azhang/
- Photos from Internet