# Chapter 7. Cluster Analysis

**Dong-Kyu Chae**

**PI of the Data Intelligence Lab @HYU**
**Department of Computer Science & Data Science**
**Hanyang University**

Data
Intelligence
LAB

# DBSCAN Algorithm

**ALGORITHM 1**: Pseudocode of Original Sequential DBSCAN Algorithm

**Input**: *DB*: Database
**Input**: $\varepsilon$: Radius
**Input**: *minPts*: Density threshold
**Input**: *dist*: Distance function
**Data**: *label*: Point labels, initially *undefined*

1  **foreach** *point p* in *database DB* **do**     // Iterate over every point
2     **if** *label(p)* ≠ *undefined* **then continue**     // Skip processed points
3     Neighbors $N \leftarrow$ RANGEQUERY(*DB, dist, p, $\varepsilon$*)     // Find initial neighbors
4     **if** $|N| < minPts$ **then**   → not core point 를 의미     // Non-core points are noise
5         *label(p)* ← Noise
6         **continue**
7     $c \leftarrow$ next cluster label   $C_1)C_2,\cdots$     // Start a new cluster
8     *label(p)* ← c
        exclude
9     Seed set $S \leftarrow N \setminus \{p\}$     // Expand neighborhood
10     **foreach** *q* in *S* **do**
11         **if** *label(q)* = *Noise* **then** *label(q)* ← c   → miss given noise → border
12         **if** *label(q)* ≠ *undefined* **then continue**
13         Neighbors $N \leftarrow$ RANGEQUERY(*DB, dist, q, $\varepsilon$*)
14         *label(q)* ← c
        → borderpoint initially undefined → border
15         **if** $|N| < minPts$ **then continue**     // Core-point check
16         $S \leftarrow S \cup N$   expand the seed set

# DBSCAN Algorithm: Example

❑ **Parameter**

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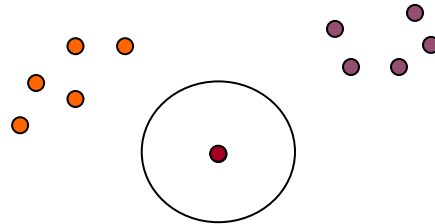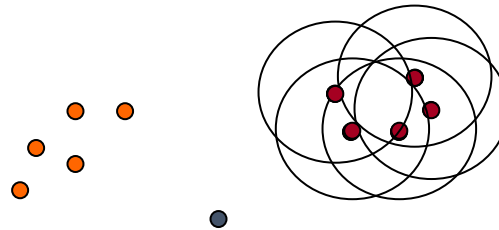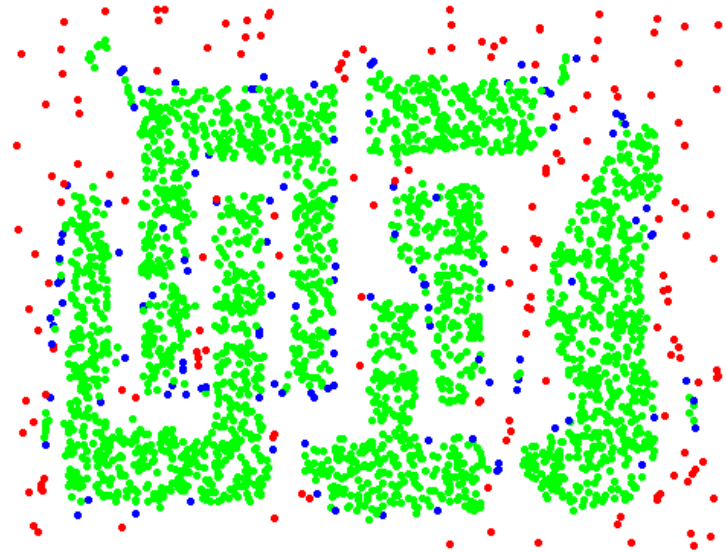


```
for each o ∈ D do
    if o is not yet marked then
        if o is a core-object then
            collect all objects density-reachable from o
            and assign them to a new cluster.
        else
            assign o to NOISE
```

# DBSCAN Algorithm: Example

❏ **Parameter**

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



> **for** each $o \in D$ **do**
>     **if** $o$ is not yet marked **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 marked **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

# Point Types Marked by DBSCAN



**Original Points**

**Point types: core** <sub>green</sub>**, border and outliers**

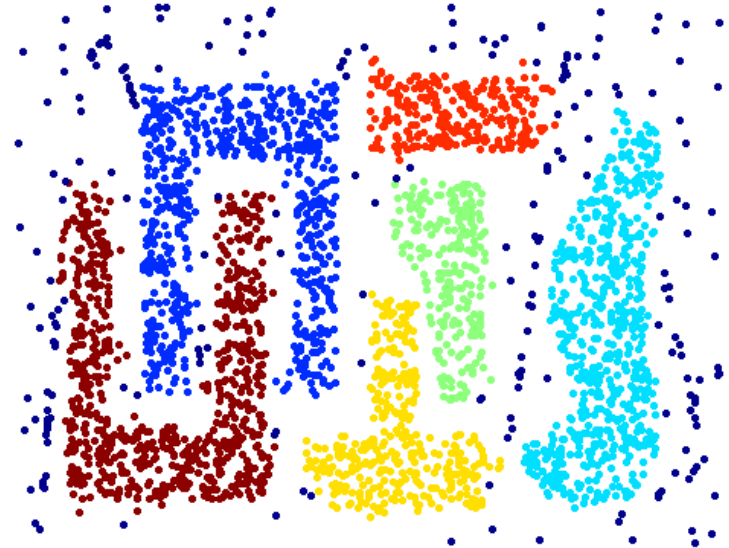$\varepsilon = 10$, **MinPts = 4**

# When DBSCAN Works Well
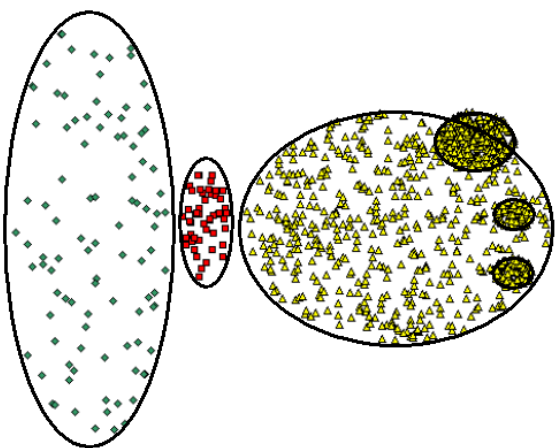


**Original Points**                    **Clusters**
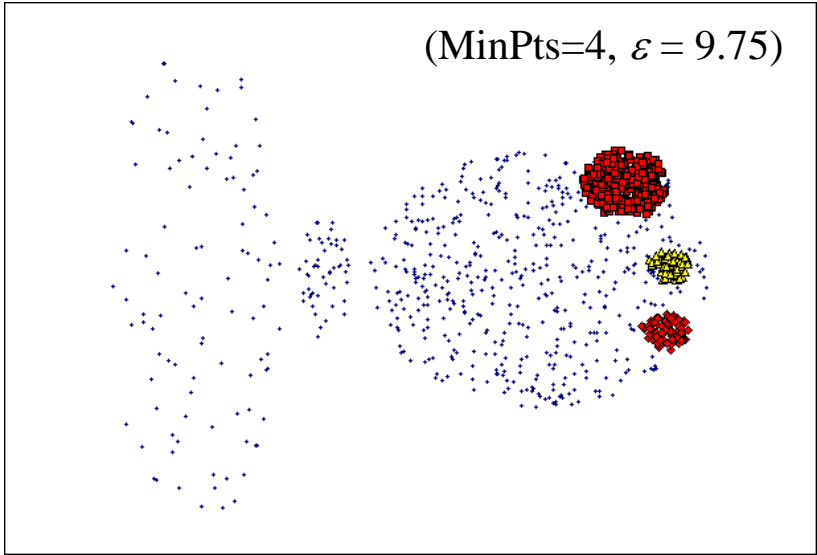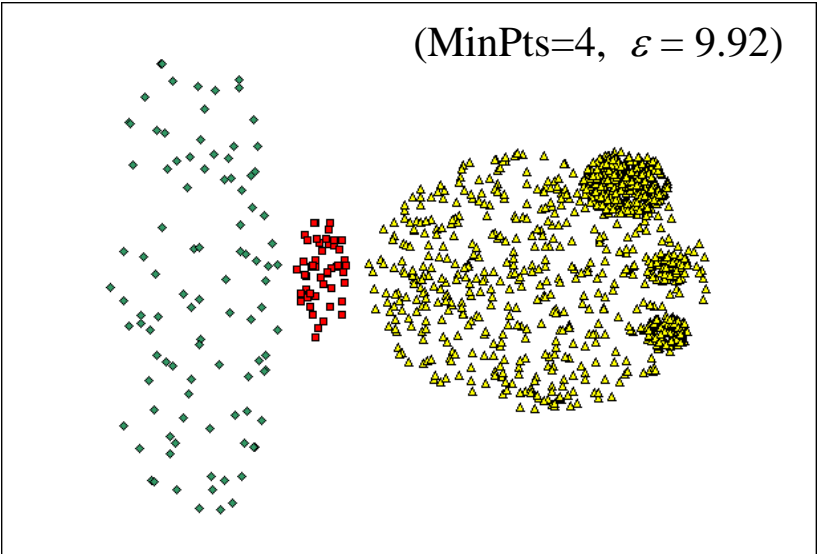
- **Resistant to Noise**

- **Can handle clusters of different shapes and sizes**

# DBSCAN: Sensitive to Parameters



(MinPts=4, $\varepsilon = 9.92$)

(MinPts=4, $\varepsilon = 9.75$)

**Original Points**

# DBSCAN: Sensitive to Parameters

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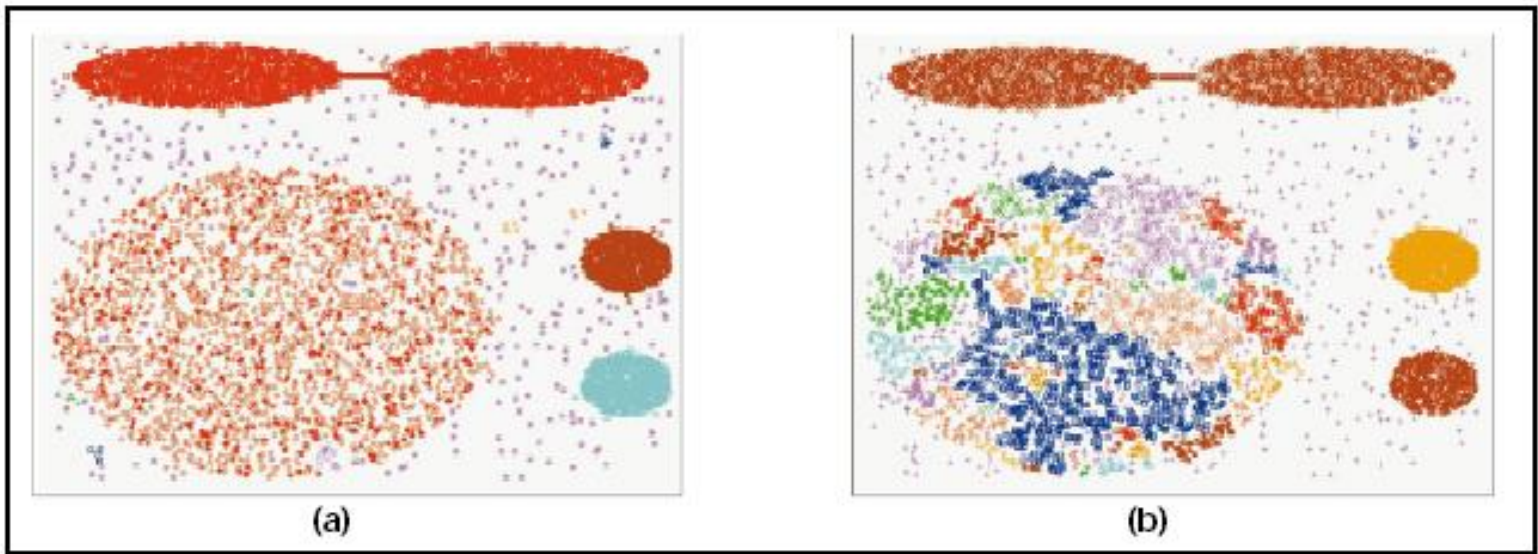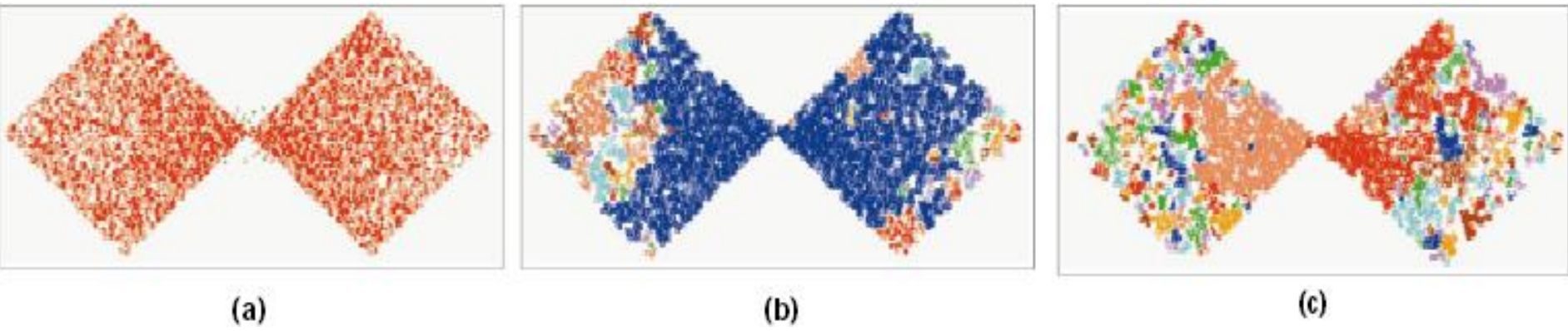
Figure 9. DBScan results for DS2 with MinPts at 4 and Eps at (a) 5.0, (b) 3.5, and (c) 3.0.

# DBSCAN Summary

❑ **Advantages**

   ❑ Clusters can have arbitrary shape and size

   ❑ Number of clusters is determined automatically

   ❑ Can separate clusters from noise and outliers

❑ **Disadvantages**

   ❑ Input parameters may be difficult to determine

   ❑ In some situations very sensitive to input parameter setting

❑ **OPTICS**

   ❑ Based on DBSCAN

   ❑ Does not produce clusters explicitly

   ❑ Rather generate **an ordering of data objects** representing density-based clustering structure

# OPTICS: Some Extension from DBSCAN

❑ **OPTICS: <u>O</u>rdering <u>P</u>oints <u>T</u>o <u>I</u>dentify the <u>C</u>lustering <u>S</u>tructure**

    ❑ It aims to answer:  "*how to choose proper $\varepsilon$ value?*"

    ❑ Produces a linear **order** of objects such that **spatially closest points become neighbors in the ordering**

    ❑ This ordering can produce a **graphical information** equivalent to density-based clustering structure corresponding to a broad range of parameter settings ($\varepsilon$)

    ❑ OPTICS can be seen as a **visualization technique** for clustering, rather than a clustering solution

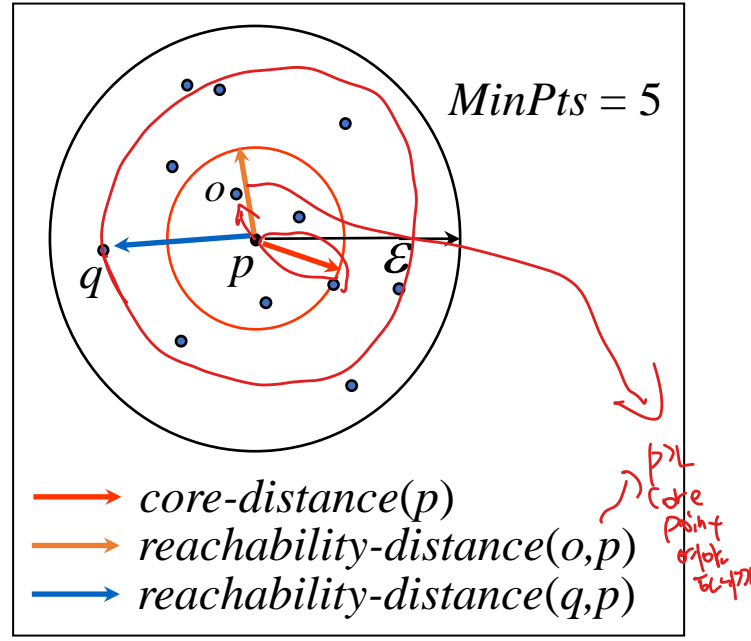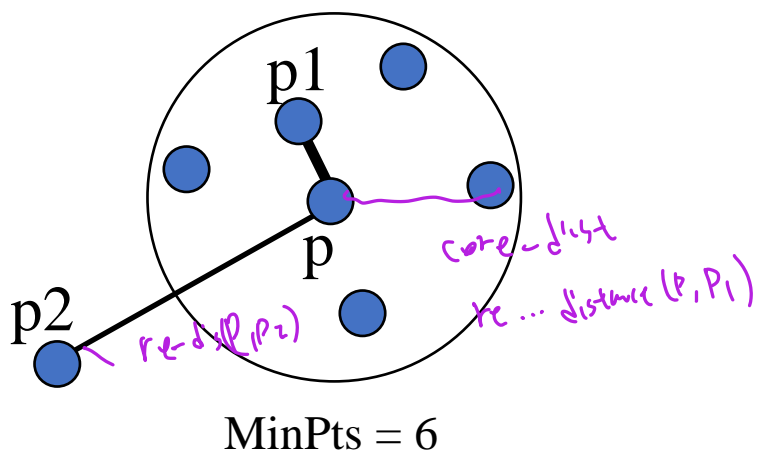    ❑ Good for both automatic and interactive cluster analysis, including finding intrinsic clustering structure

# OPTICS: Some Extension from DBSCAN

- $\text{core-dist}_{\varepsilon,MinPts}(p) = \begin{cases} \text{UNDEFINED} & \text{if } |N_\varepsilon(p)| < MinPts \\ MinPts\text{-th smallest distance in } N_\varepsilon(p) & \text{otherwise} \end{cases}$

  **"smallest distance that makes _p_ a core"**

- $\text{reachability-dist}_{\varepsilon,MinPts}(o,p) = \begin{cases} \text{UNDEFINED} & \text{if } |N_\varepsilon(p)| < MinPts \\ \max(\text{core-dist}_{\varepsilon,MinPts}(p), \text{dist}(p,o)) & \text{otherwise} \end{cases}$

  **"smallest distance that makes _o_**
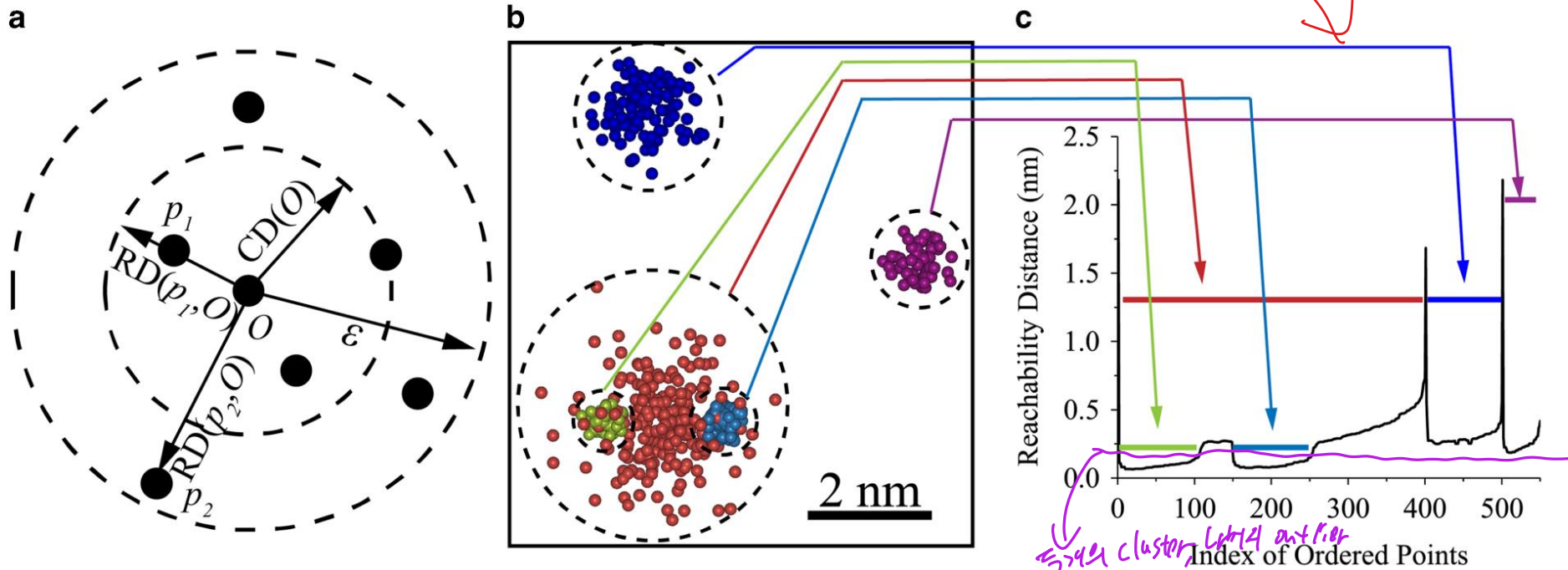
  **directly density-reachable** from _p_"



MinPts = 6

MinPts = 5

core-distance(p)
reachability-distance(o,p)
reachability-distance(q,p)

# OPTICS: Some Extension from DBSCAN

❑**Idea**

*detail ship*

❑ Order data points by ***shortest reachability distance***

❑ The clusters show up as **valleys** in the **reachability plot**. The deeper the valley, the denser the cluster (having shorter reachability distances).

❑ Extracting clusters from this plot can be done manually by selecting a threshold on the y-axis

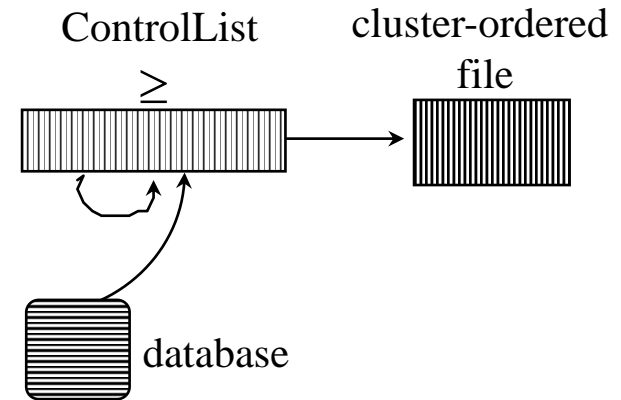❑ Needs an initial ε*, but it does not require a careful decision.

a

b

c

2 nm

Reachability Distance (nm)

Index of Ordered Points

등수가 cluster Lebel4 outlier

# The Pseudocode of OPTICS (skip)

❏ **Output: ordered data points**

ControlList   cluster-ordered file

$\geq$

database

**foreach** $o \in$ Database
  // initially, $o$.processed = false for all objects $o$
  **if** $o$.processed = false;
    insert ($o$, *"undefined"*) into *ControlList;*
  **while** *ControlList* is not empty
     select first element ($o$, *r-dist*) from *ControlList*;
     retrieve $N_\varepsilon(o)$ and determine *c_dist= core-distance*($o$);
     set $o$.processed = true;
     write ($o$, *r_dist*, *c_dist*) to file;
     **if** $o$ is a core object at any distance $\leq \varepsilon$
      **foreach** $p \in N_\varepsilon(o)$ not yet processed;
        determine $r\_dist_p = $ *reachability-distance*($p$, $o$);
        **if** ($p$, _) $\notin$ *ControlList*
          insert ($p$, $r\_dist_p$) in *ControlList;*
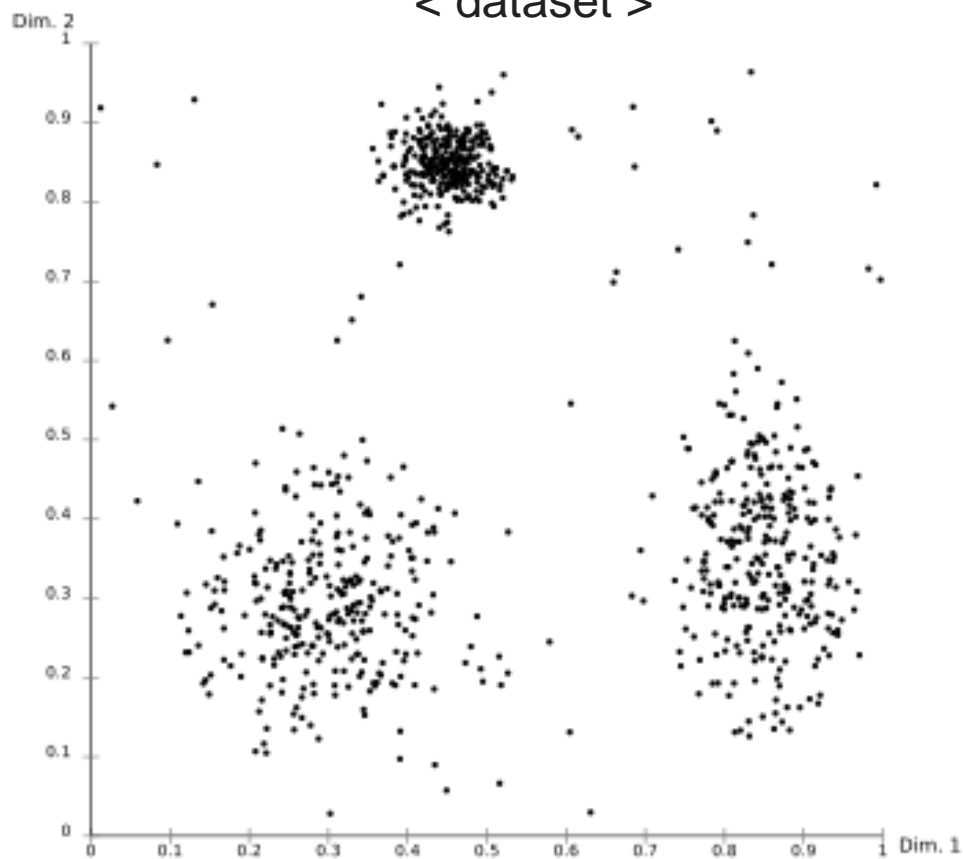        **else if** ($p$, *old_r_dist*) $\in$ *ControlList* **and** $r\_dist_p < old\_r\_dist$
          update ($p$, $r\_dist_p$) in *ControlList;*

# OPTICS: The Reachability Plot



**Reachability-distance**

neighboring objects stay close to each other in a linear sequence.

undefined

$\varepsilon^*$

$\varepsilon$

$\varepsilon$

**Cluster-order of the objects**

< dataset >

< 1 Nearest-neighbor tree >

Dim. 2

Dim. 1
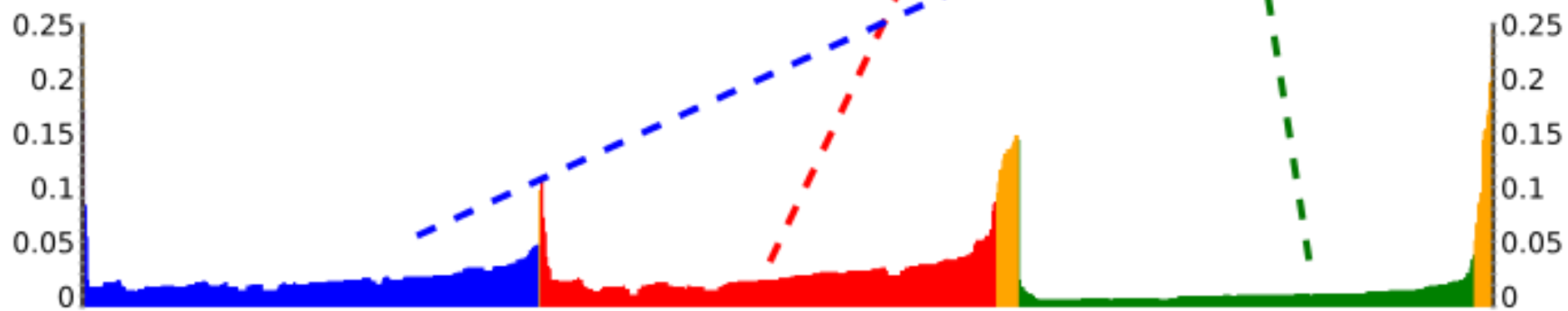
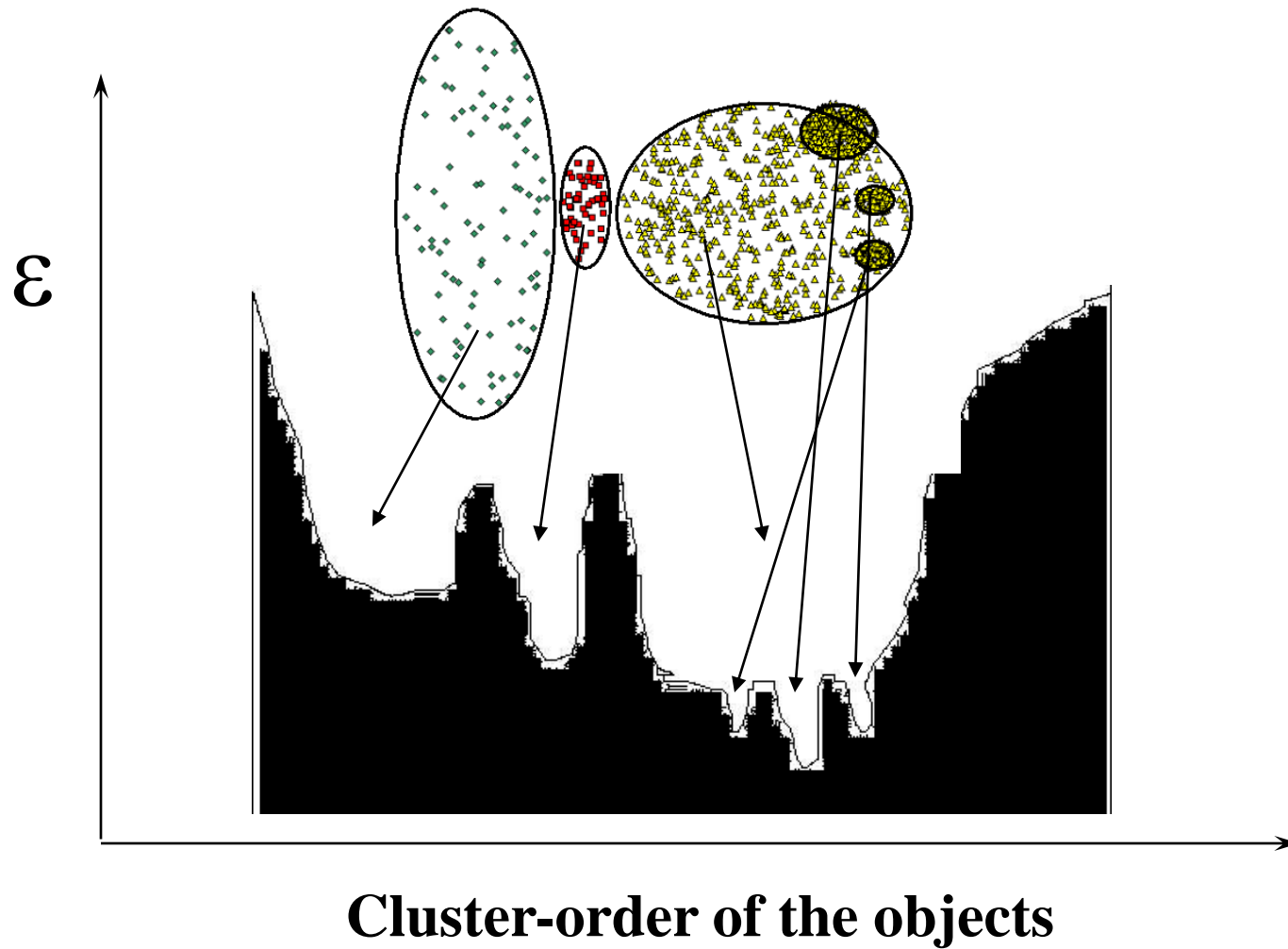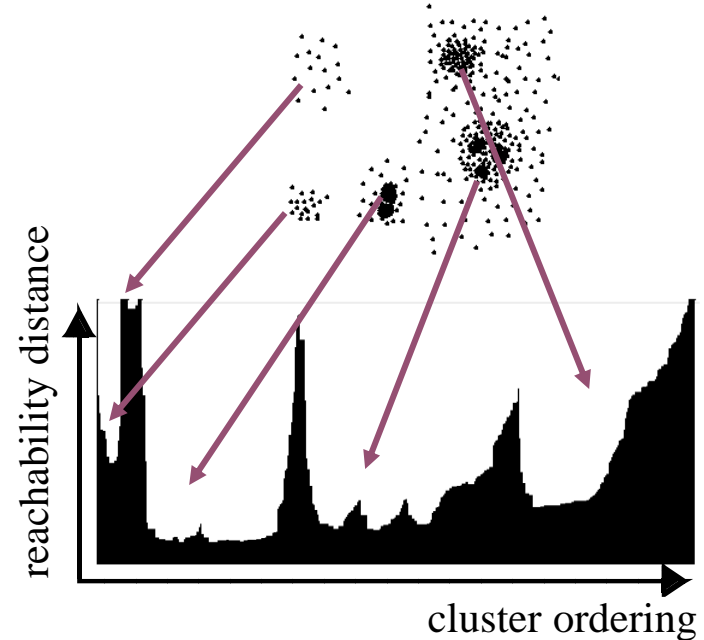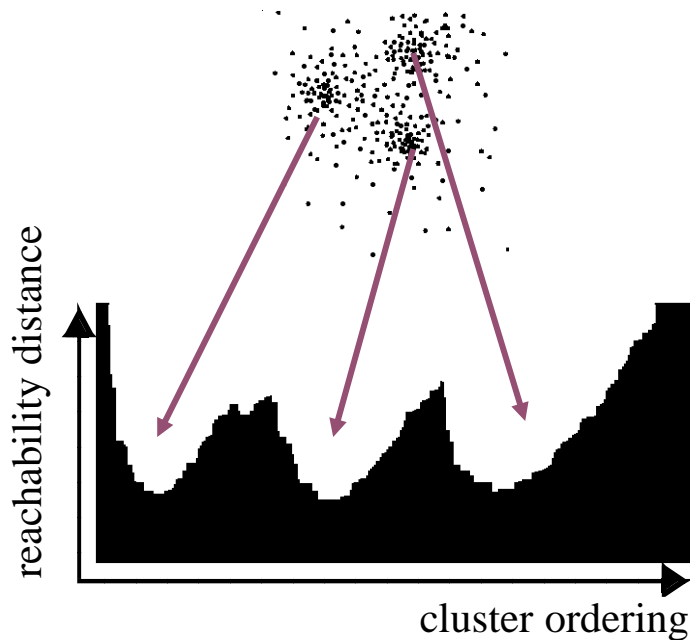distance

order of objects

$\varepsilon$

**Cluster-order of the objects**

# OPTICS: The Reachability Plot

❑ **Represents the density-based clustering structure**

❑ **Easy to analyze**

❑ **Independent of the dimension of the data**

# DBSCAN  VS  OPTICS

| | DBSCAN | OPTICS |
|---|---|---|
| **Density** | Boolean value (high/low) *border or outlier* | Numerical value (core distance) |
| **Density-connected** | Boolean value (yes/no) | Numerical value (reachability distance) |
| **Searching strategy** | random | greedy |

34

# Thank You