# Cluster Analysis:
# Density-Based Methods

CS5483 Data Warehousing and Data Mining

# How to identify the following clusters?



$d_1 < d_2$

- Remedy?
- Why they should be clusters?
  They are **d_____** regions.

- Why centroid-based method fails?
  - Both clusters have the same c_____.
  - Bias towards s_____ cluster.
- Why single-linkage fails?
  - The c_____ distance of the two clusters is no larger than those between two points in the same cluster.
  - Chaining phenomenon
- Why complete-linkage method fails?
  - The f_____ distance of the two clusters is no larger than those between two points in one cluster.
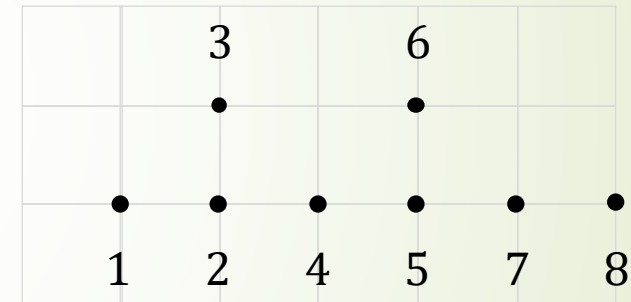  - Bias towards s_____ cluster.

# Identify pillars of dense regions

- $\varepsilon$-neighbourhood of $\boldsymbol{p} \in D$ is the region:
  $$N_\varepsilon(\boldsymbol{p}) := \{\boldsymbol{q} \in \mathbb{R}^d | \text{dist}(\boldsymbol{p}, \boldsymbol{q}) \leq \varepsilon\}$$
  - Within r_____ $\varepsilon > 0$
  - from the c_____ $\boldsymbol{p}$.
- C_____ **point**: $\boldsymbol{p} \in D$ such that
  $$|D \cap N_\varepsilon(\boldsymbol{p})| \geq \text{MinPts}$$
- How to find clusters?

$\varepsilon = 1$ MinPts $= 4$

```
            3       6
            •       •


    •   •   •   •   •   •
    1   2   4   5   7   8
```
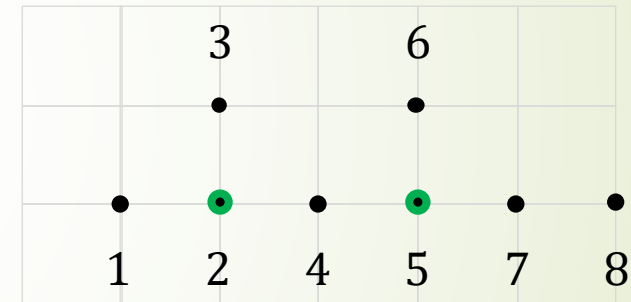
Core points:_____

# Density-reachability

- *q* is **directly density-reachable** from *p*, denoted as $p \rightarrow q$, if

  1. *p* is a core point and

  2. $q \in N_\varepsilon(p)$

- *q* is **density-reachable** from *p* if there is a **path**
$$p \rightarrow \cdots \rightarrow q.$$

  - All points in the path must be core points except *q*.

  - *q* is called a **border** point if it is not a core point.

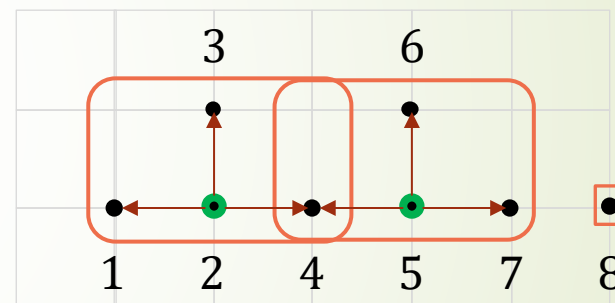$\varepsilon = 1$     $\text{MinPts} = 4$



- $p_1 \rightarrow p_2$ ? _____

# Density-connectedness

- $q$ is **density-connected** with $p$, denoted as $p \sim q$, if $p, q$ are reachable from a c_____ core point.
- Density-connected components as clusters? Y/N because _____

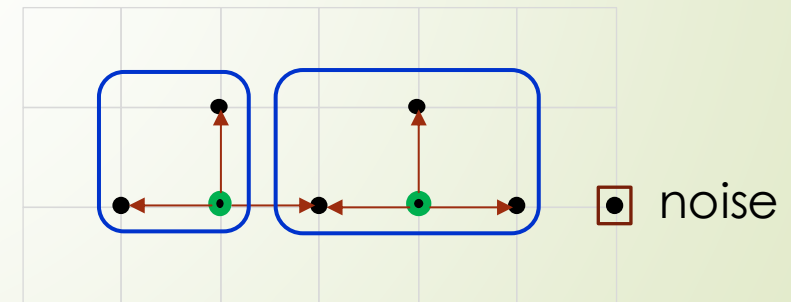$\varepsilon = 1$     MinPts $= 4$
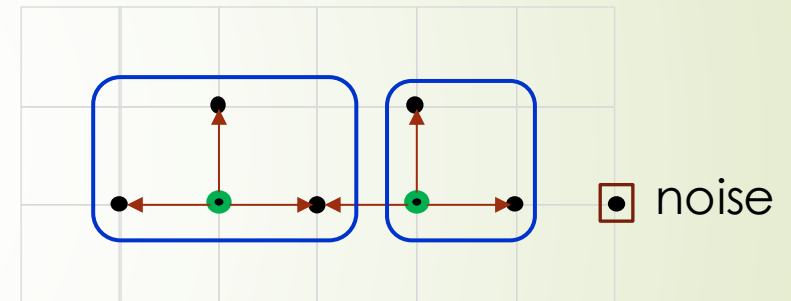


- $p_1 \sim p_2$ ? _____

Border point: _____

# DBSCAN
## **D**ensity-**B**ased **S**patial **C**lustering of **A**pplications with **N**oise

1. Identify core points and their density reachable points.

2. Return density-connected components of core points as clusters.

3. Assign border points to one of the cluster it is density-connected to.

4. Label the remaining points as noise.
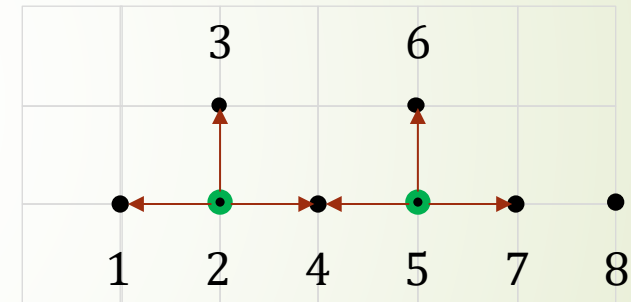
➡ Clustering solution is not **u**_____.

$\varepsilon = 1$      $\text{MinPts} = 4$



▣ noise



▣ noise

# Uniqueness of clusters for core points

- DBSCAN give unique clusters for core points because density-connectedness is an equivalence relation on core points:

  - r_____: $p \sim p$

  - s_____: $p \sim q \Leftrightarrow q \sim p$

  - t_____: $p \sim r, r \sim q \Rightarrow p \sim q$

- However, transitivity can fail on border points:

  - E.g., _____
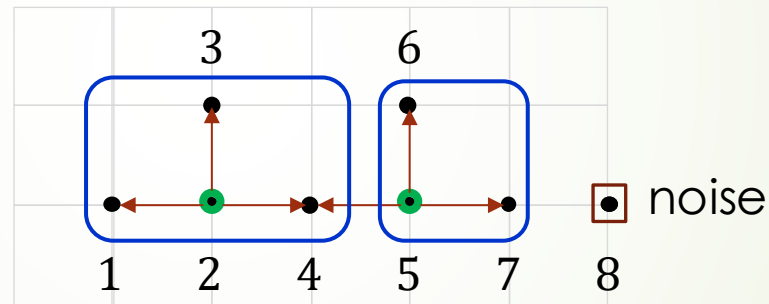
  - see the last slide for the mistake in [Han11].

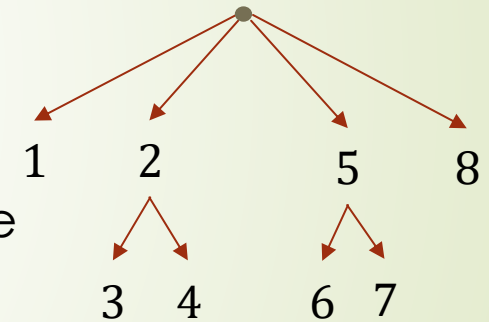$\varepsilon = 1 \qquad \mathrm{MinPts} = 4$

# Implementation

- Repeatedly start a breadth first search (BFS) on a new point to add density-reachable points (not already assigned to another cluster) to a cluster.

- Label remaining points as noise.

- Complexity: _____ (or $O(n \log n)$ with spatial indexing.)
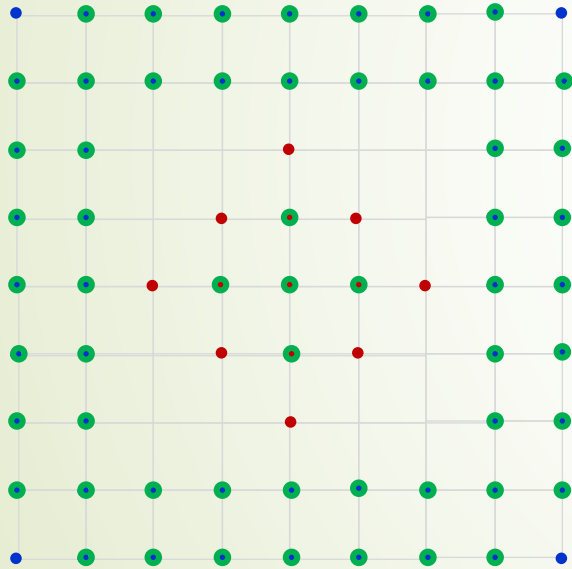
$\varepsilon = 1$     MinPts $= 4$

Nodes visited by BFS
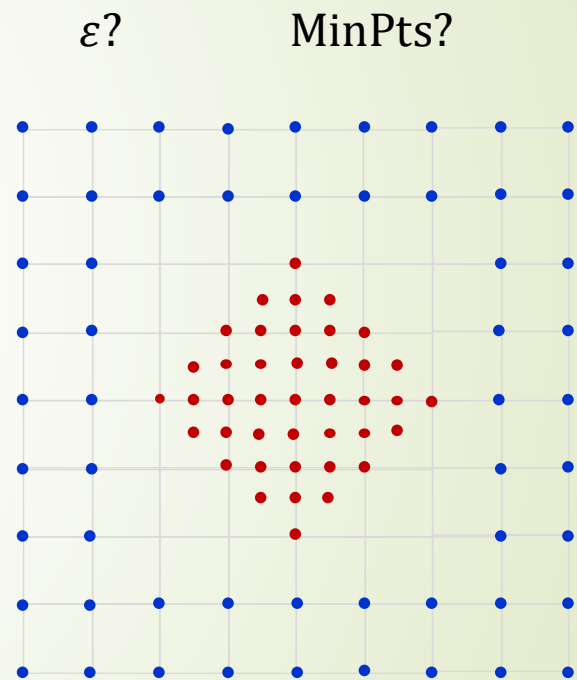
# Clusters can be non-convex

$\varepsilon = 1$     MinPts $= 4$



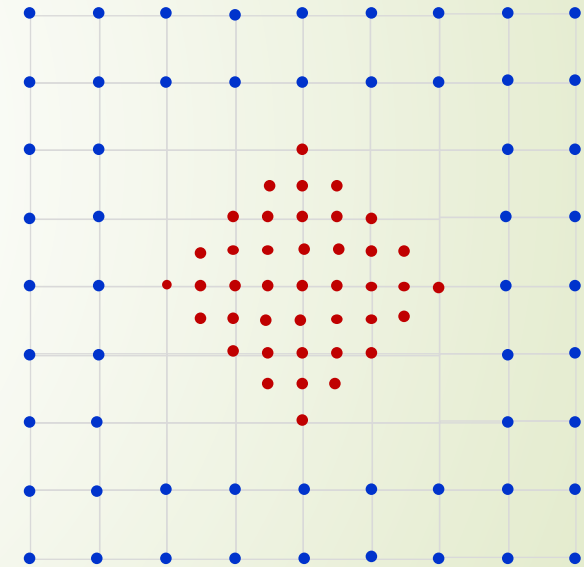➡ Identify the core points, clusters, ambiguous border points and the noise.

# Limitations

- How to choose MinPts and $\varepsilon$?
  - For the outer ring, want $\varepsilon \geq$ ___.
  - But the corner point of the inner diamond will be _____.
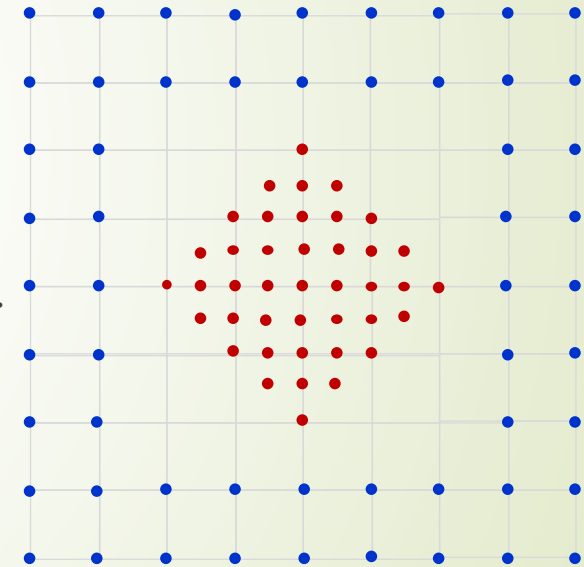- What about clusters with different densities?

$\varepsilon$?          MinPts?

# Density-based method

- Can DBSCAN recover the two clusters?
  - How to choose the parameters $\varepsilon$ and MinPts?
- How to handle clusters with different densities?
  - Varying $\varepsilon$ to obtain a h_____ of clusters.
  - How to do this efficiently?
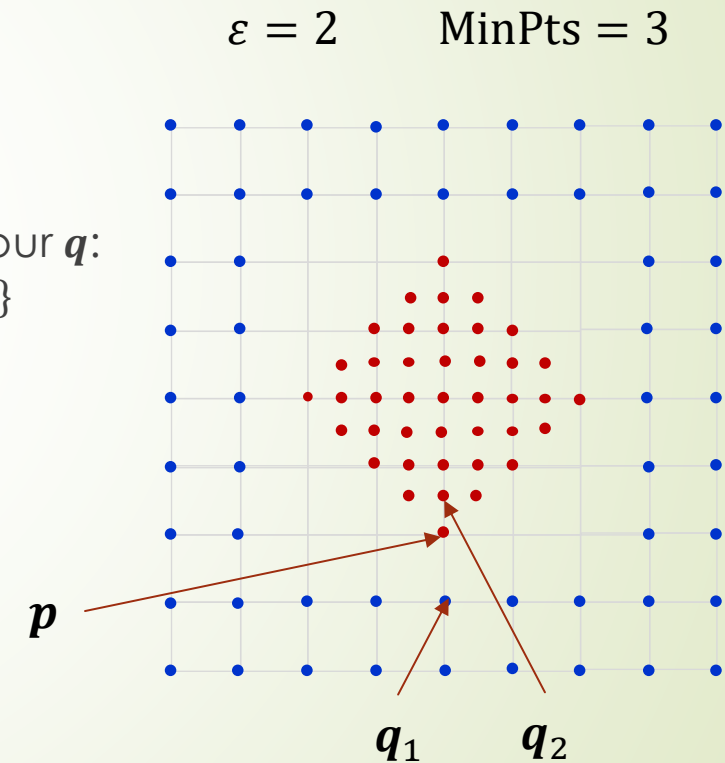
# OPTICS

- **O**rdering **P**oints to **I**dentify the **C**lustering **S**tructure.

  - Idea: Apply DBSCAN but visit c_____ points first.

- Choice of parameters:

  - MinPts normally chosen as dimension + 1. Why?
    [Optional] See Carathéodory's theorem.

  - Choose the worst-case (largest) radius $\varepsilon$ for the neighborhood.

- For the example, we can choose

  - MinPts =_____

  - $\varepsilon \geq$_____, where all the points are density-connected.

# Core and reachability distances

- When a core point $\boldsymbol{p}$ is reached,

  - Calculate the true density:
    $$\text{core}-\text{distance}(\boldsymbol{p}) \coloneqq \min\{0 \leq \varepsilon' \leq \varepsilon || D \cap N_{\varepsilon'}(\boldsymbol{p})| \geq \text{MinPts}\}$$

  - Calculate the distance to its density-reachable neighbour $\boldsymbol{q}$:
    $$\text{reachability}-\text{distance}(\boldsymbol{q}) \coloneqq \max\{\text{dist}(\boldsymbol{p}, \boldsymbol{q}), \text{core}-\text{distance}(\boldsymbol{p})\}$$

- Visit neighbors with smaller reachability distance first.
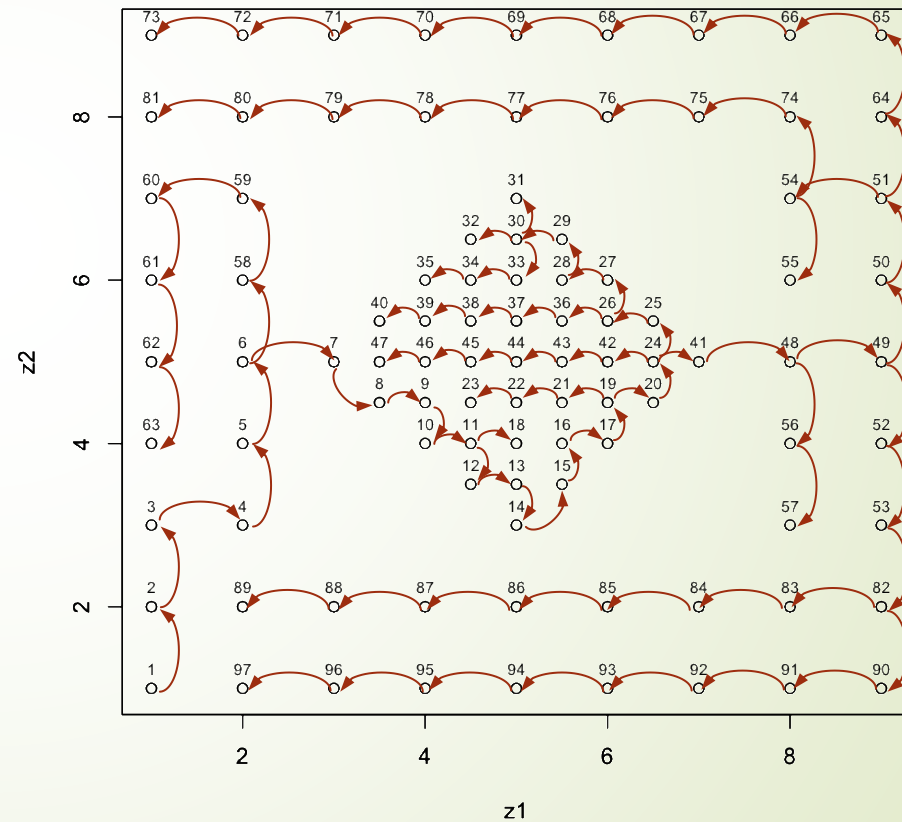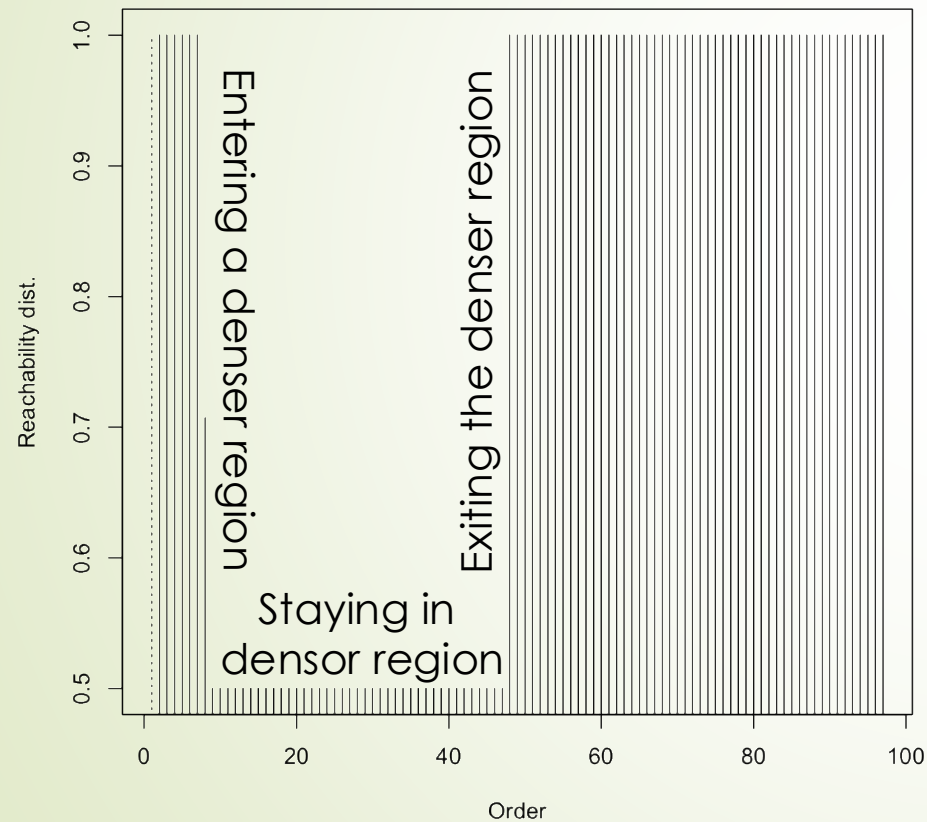
- Why use reachability distance instead of distance?

_____

$\text{core}-\text{distance}(\boldsymbol{p}) = \underline{\qquad\qquad}$
$\text{reachability}-\text{distance}(\mathbf{p}, \boldsymbol{q}_1) = \underline{\qquad\quad}$
$\text{reachability}-\text{distance}(\mathbf{p}, \boldsymbol{q}_2) = \underline{\qquad\quad}$

$\varepsilon = 2 \qquad \text{MinPts} = 3$

$\boldsymbol{p}$
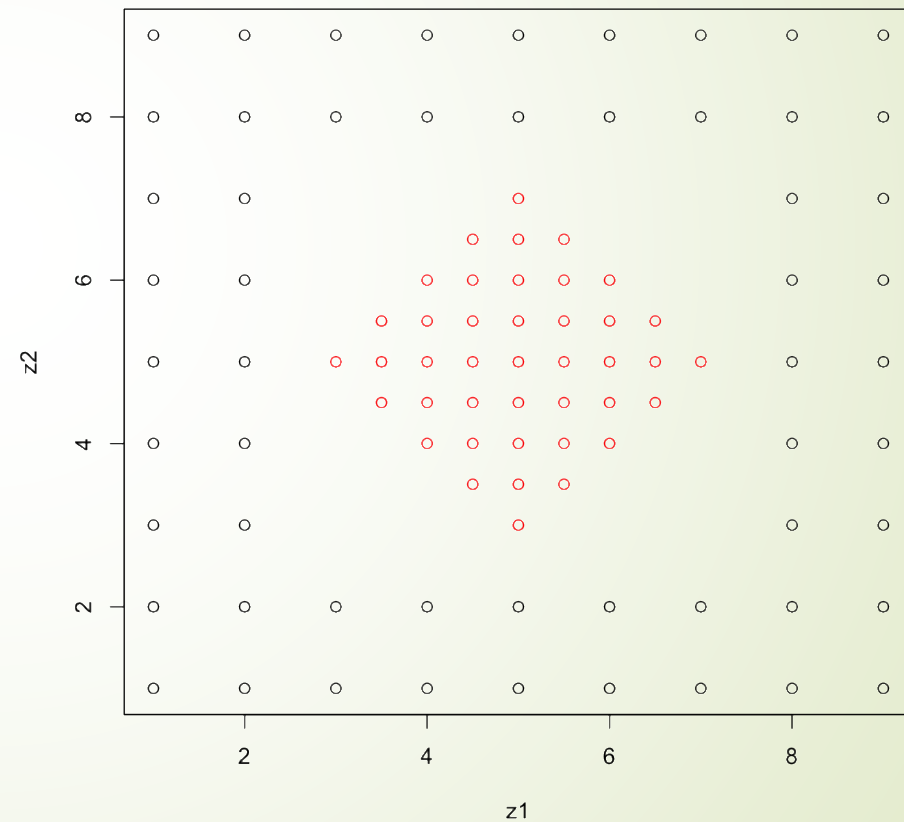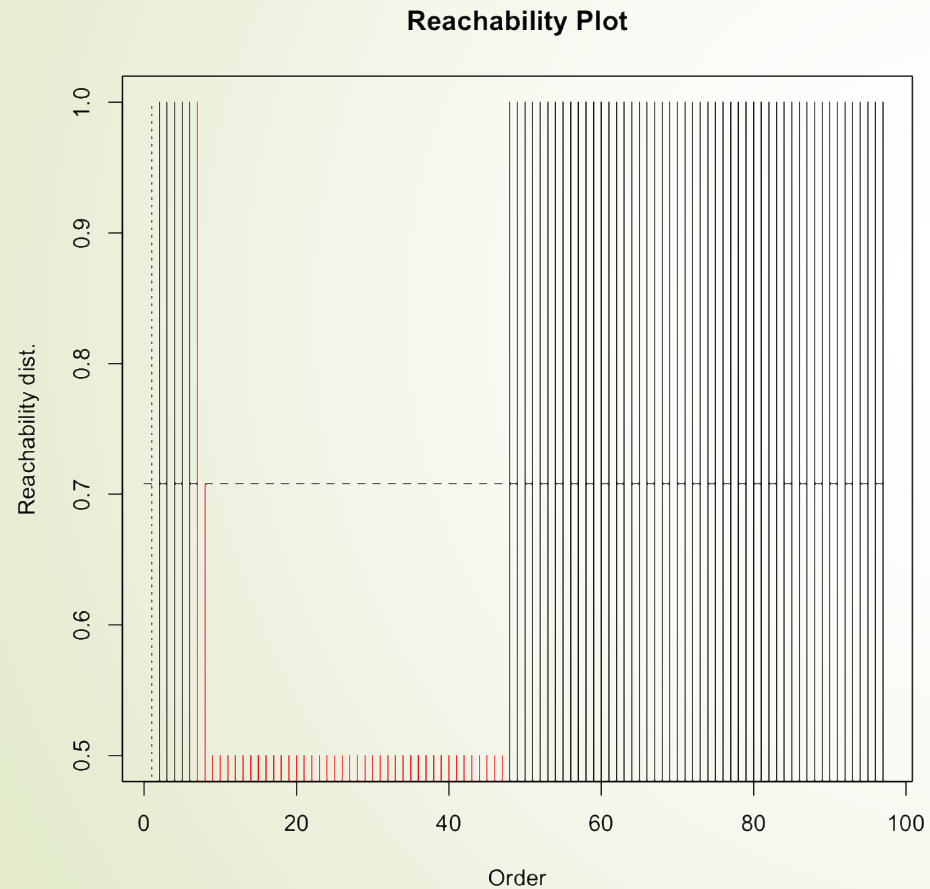
$\boldsymbol{q}_1 \qquad \boldsymbol{q}_2$

# Reachability plot

 ➧ Plot the reachability distances of the sequences of visited points.



Reachability Plot

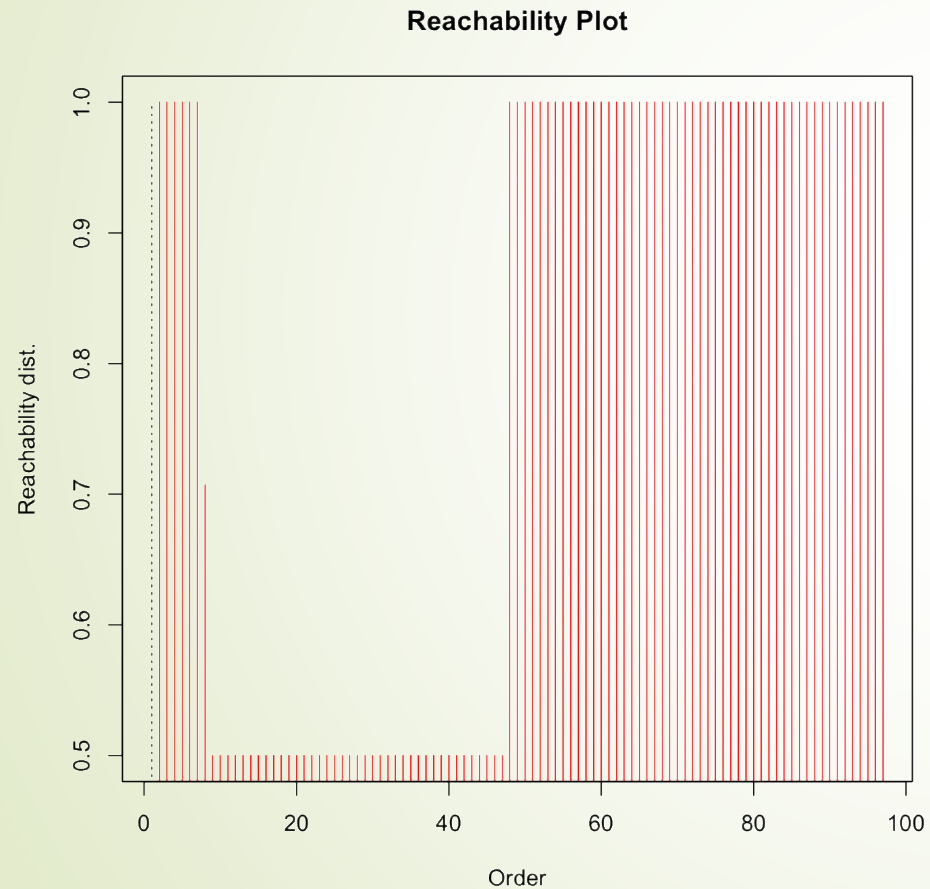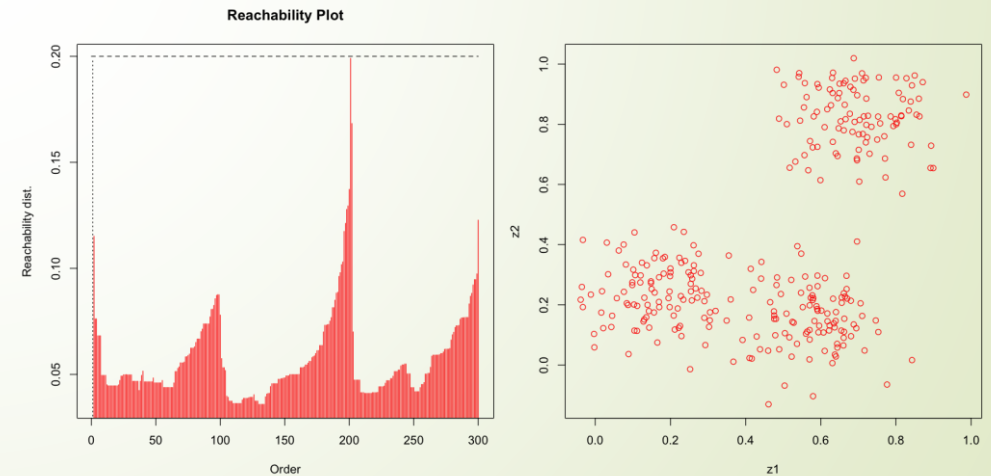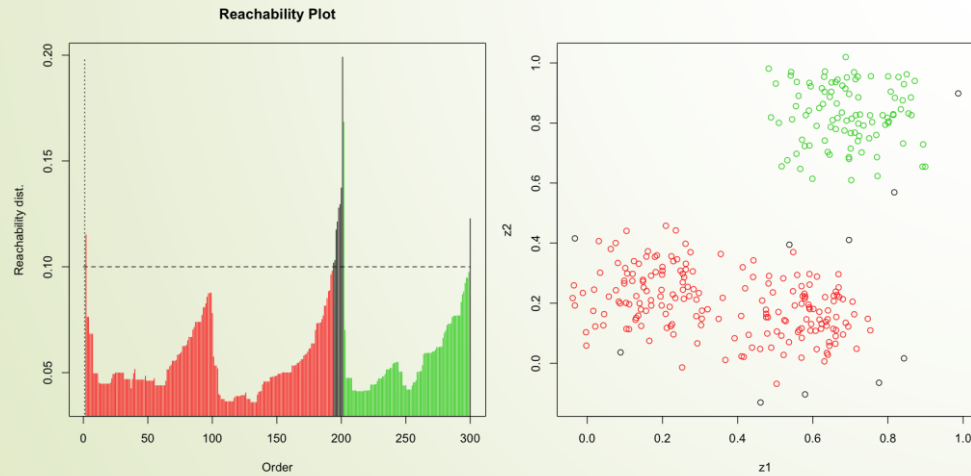A cluster at threshold $\varepsilon$ is discovered as the valley of points with reachability distance below $\varepsilon$ except for the first point.



Reachability Plot

- With the threshold of reachability raised to __ or above, the valley covers all the points.
- Note that the outer ring is not identified as a separate cluster from the inner diamond.



Reachability Plot

# An example with multiple clusters

# References

- 10.4.1 DBSCAN: Density-Based Clustering Based on Connected Regions with High Density
  - Errata on p.472: Density-connectedness is NOT an equivalence relation unless restricted to only core points.
- McInnes, Leland, and John Healy. "Accelerated Hierarchical Density Based Clustering." *Data Mining Workshops (ICDMW), 2017 IEEE International Conference on.* IEEE, 2017.