# 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

# Contents

# CHAMELEON

❑ **Main idea**

  ❑ Two clusters can be merged only if the interconnectivity and closeness (proximity) between two clusters are high

  ▪ **Relative** to the internal interconnectivity of the clusters and internal closeness of items within the clusters

# 1. Draw a $k$-nearest neighbor graph (KNN graph) first

  ❑ **Node**: object, **edge**: k-nearest neighbor's link, **weight**: similarity
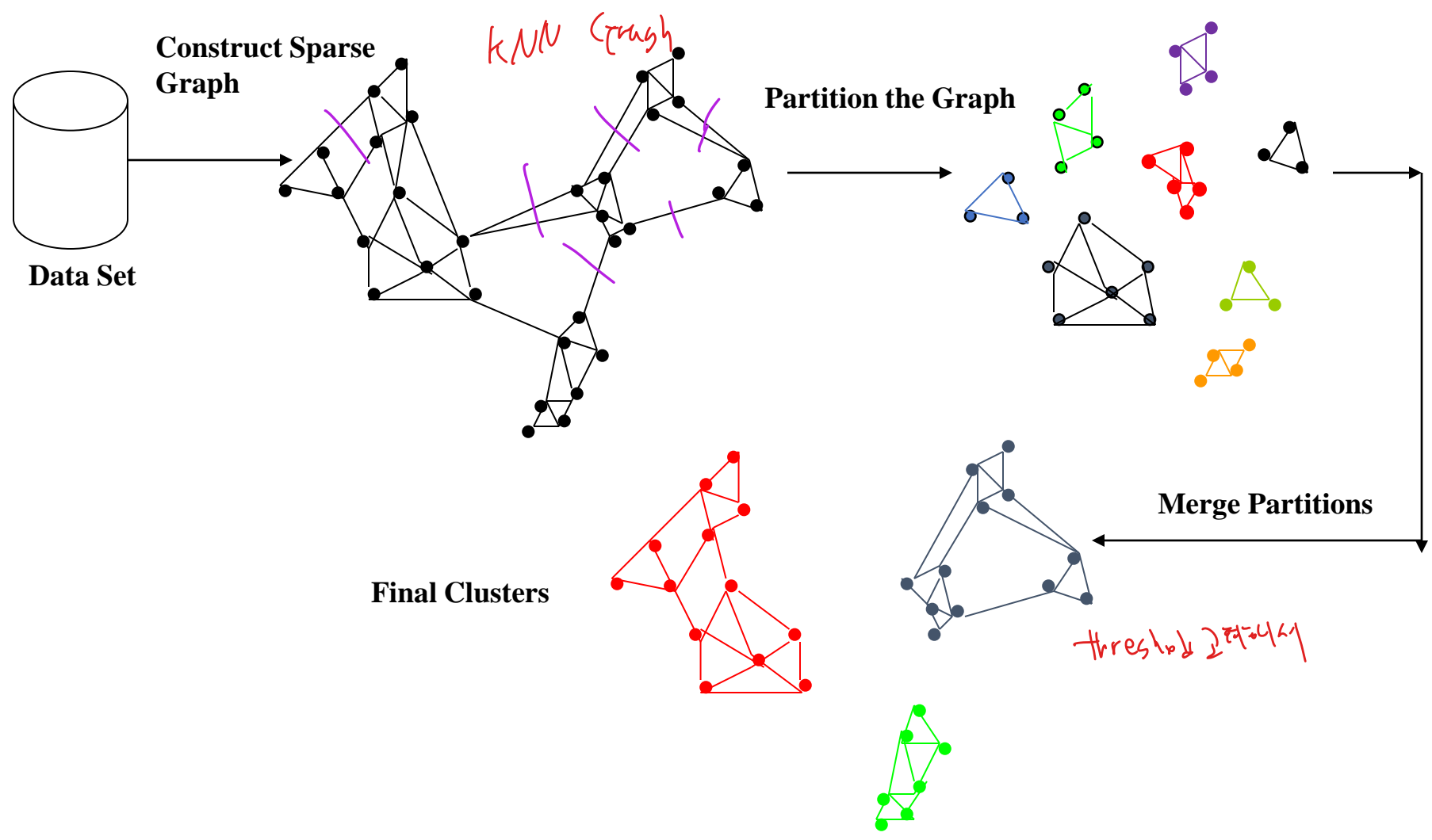
# 2. Partition: Use a graph partitioning algorithm

  ❑ Divide the KNN graph into a large number of relatively small sub-clusters

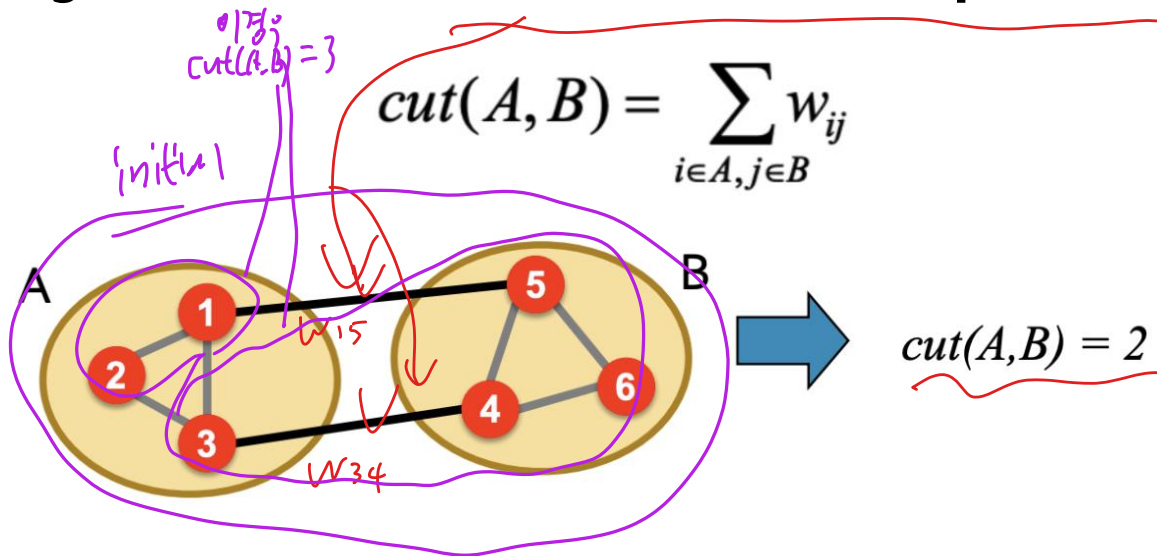# 3. Merge: Use an agglomerative hierarchical clustering algorithm

  ❑ Iteratively find clusters by repeatedly combining these sub-clusters

# Overall Framework of CHAMELEON

Construct Sparse Graph

KNN Graph

Data Set

Partition the Graph

Merge Partitions

Final Clusters

threshold 값 이상이면 합침

# CHAMELEON: Partitioning

❑ **Partition the KNN graph such that the edge cut is minimized.**

   ❑ The edge-cut of a partition is **the sum of the weights** of edges whose **vertices lie in different partitions**.



$$cut(A,B) = \sum_{i \in A, j \in B} w_{ij}$$

*(handwritten annotations: cut(A,b) = 3, initial, W15, W34, cut(A,B) = 2)*

   ❑ hMeTiS library (**METIS**) is used

      ▪ Tries to split a graph into two subgraphs of nearly equal sizes

*(handwritten: not only edge cut is minimized but also balance number of data points inside partition)*

# CHAMELEON: Merging

❏ **Merging the partitions**

   ❏ This step computes the cluster similarity based on the **relative inter-connectivity** and **relative closeness** of the clusters.

❏ **Relative inter-connectivity**

$$RI(C_i, C_j) = \frac{|EC_{\{C_i, C_j\}}|}{\frac{1}{2}(|EC_{C_i}| + |EC_{C_j}|)},$$

*(handwritten annotations: "criteria", "set of edges", "Standard", "Ci", "3", "2", "Cj", Korean notes, "merge")*

   ❏ $EC_{\{Ci, Cj\}}$ = edges that connect Ci and Cj.
   ❏ $EC_{Ci}$ = edges that partition the cluster into roughly equal parts.

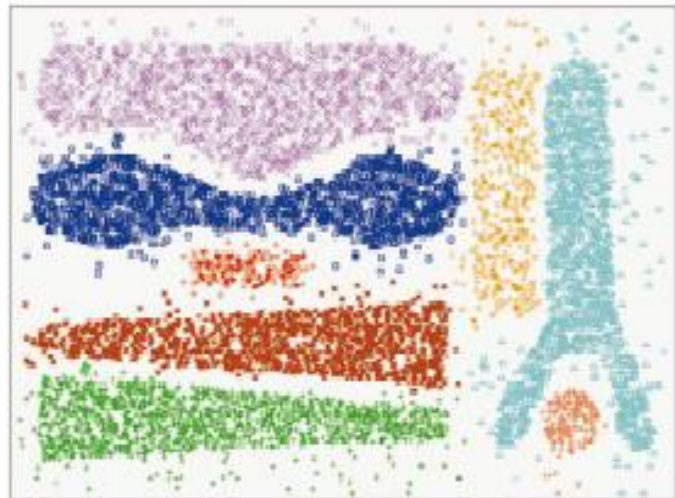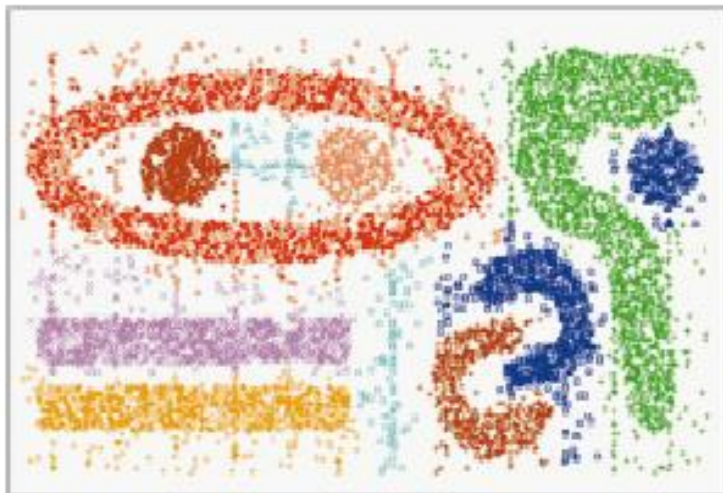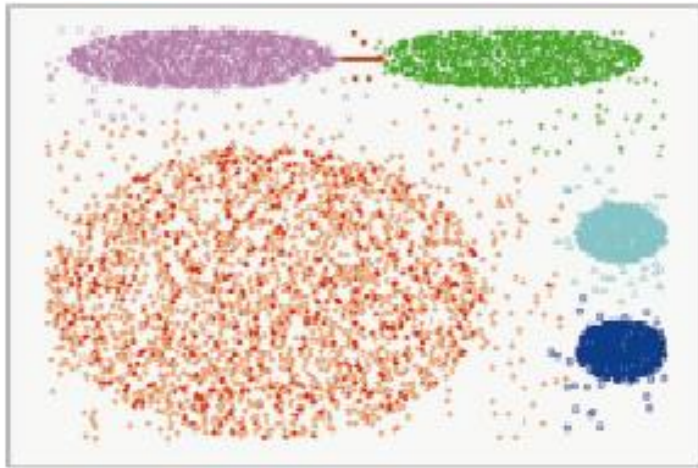# CHAMELEON: Merging

❑ **Relative closeness**

$$RC(C_i, C_j) = \frac{\bar{S}_{EC_{\{C_i, C_j\}}}}{\frac{|C_i|}{|C_i|+|C_j|}\bar{S}_{EC_{C_i}} + \frac{|C_j|}{|C_i|+|C_j|}\bar{S}_{EC_{C_j}}},$$

❑ $\bar{S}_{EC_{\{C_i, C_j\}}}$ = average weight of the edges from Ci to Cj

❑ $\bar{S}_{EC_{C_i}}$ = average of the weights of the edges in the cluster.

❑ **Merging**

❑ So far, we have got **Relative Inter-Connectivity** and **Relative Closeness**

❑ Using them:

$$RI(C_i, C_j) * RC(C_i, C_j)^\alpha$$

▪ where alpha controls the importance of RC

# CHAMELEON (Clustering Complex Objects)

# Contents

# Why Density-based Clustering?



**DATASET**

**K-MEANS**

**Hierarchical Clustering**

**DBSCAN**

Controid에
의거에 부하기 저서
가깝지도
다른 cluster
에 있다

handle noise
→ not the member of cluster

does not depend on distance

different shapes 가능

# Density-Based Clustering Methods

❑ **Clustering based on density (local cluster criterion), such as density-connected points, rather than just a distance**

❑ **Major features:**

  ❑ Discover clusters of arbitrary shape

  ❑ Handle noise

  ❑ One scan, thus being efficient

  ❑ Need density parameters as termination condition

  *programmer가 optimal한 hyper parameter 정해야함*

❑ **Several interesting studies:**

  ❑ **DBSCAN**

  ❑ **OPTICS**

  ❑ CLIQUE

# Density-Based Clustering: Hyper-Parameters

❑ **Two parameters:**

    ❑ $\varepsilon$: radius for the neighborhood of any point p:

        $N_\varepsilon(p) := \{$any $q$ in dataset $D\,|\,dist(p, q) \leq \varepsilon\}$

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

    ❑ *MinPts*: minimum number of points in the given neighborhood

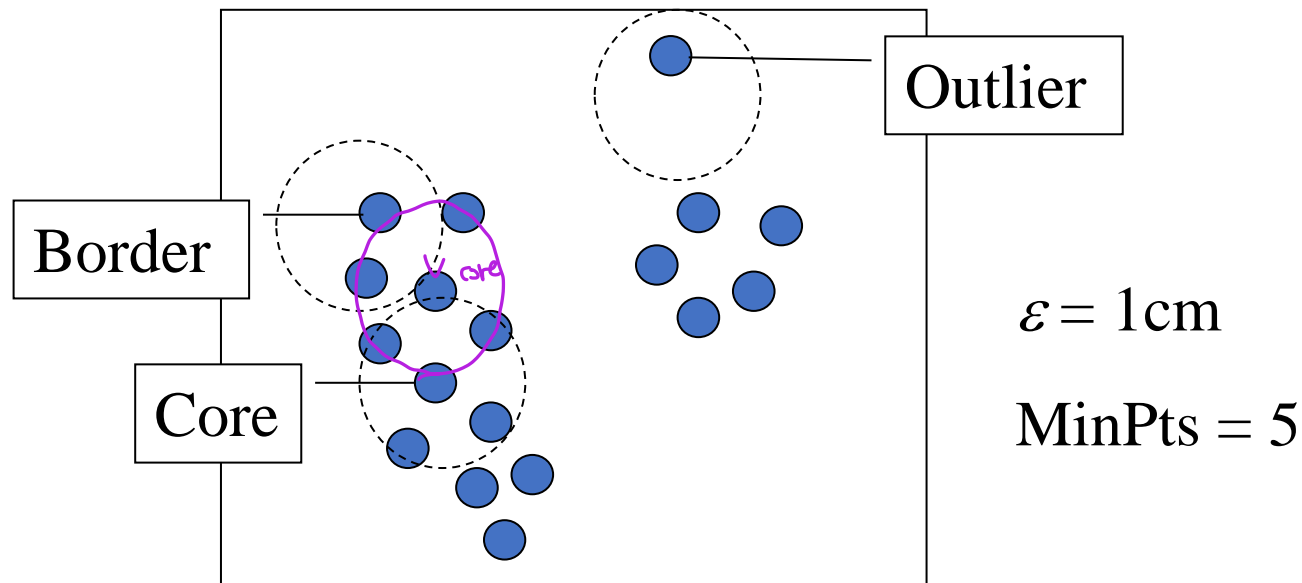        ▪ "High density" : ε-Neighborhood of an object contains at least *MinPts* of objects



ε-Neighborhood of p

ε-Neighborhood of q

ε  ε

q  p

| Density of p is "high" |
| Density of q is "low" |

**(MinPts = 4)**

# Density-Based Clustering: Types of Points

❑ **A point is a core point** if it has points more than *MinPts* within $\varepsilon$

❑ **A border point** has fewer than MinPts within $\varepsilon$, but is in the neighborhood of a core point

❑ **Outlier** is any point that is not a core point nor a border point. It is thus a noise, or an outlier.

Outlier

Border

Core

$\varepsilon = 1\text{cm}$

$\text{MinPts} = 5$

# Directly Density-Reachable

❑ **Directly density-reachable**: **A point $q$ is directly density-reachable from a point $p$ if $p$ is a core object and $q$ is in $p$'s ε-neighborhood.**
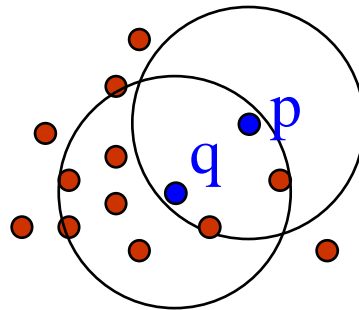


MinPts = 4

- $q$ is directly density-reachable from $p$
- $p$ is **NOT** directly density-reachable from $q$
- Density-reachability is **asymmetric**.

# Directly Density-Reachable

❑ **Directly density-reachable: A point $q$ is directly density-reachable from a point $p$ if $p$ is a core object and $q$ is in $p$'s ε-neighborhood.**

    ❑ *$p$ belongs to $N_\varepsilon(q)$*

    ❑ *$p$ is directly density-reachable from $q$*

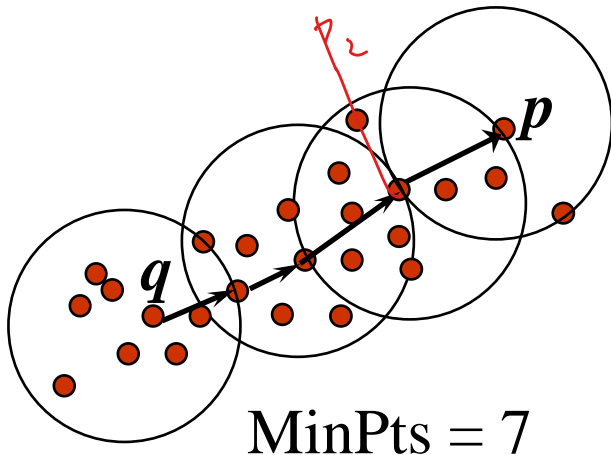    ❑ *$q$ is NOT directly density-reachable from $p$*



MinPts = 5

# Density-Reachability

- **Density-Reachable (directly and indirectly):**

  - A point p is directly density-reachable from p2;

  - p2 is directly density-reachable from p1;

  - p1 is directly density-reachable from q;

  - p←p2←p1←q form a chain.

  - Then, $p$ is (**indirectly**) **density-reachable** from $q$



MinPts = 7

- A point $p$ is density-reachable from a point $q$ if **there is a chain of points** $p_1, \ldots, p_n$, $p_1 = q$, $p_n = p$ such that $p_{i+1}$ is directly density-reachable from $p_i$

- Is $q$ not density-reachable from $p$? => NO

  $p$ is not core

# Thank You