



KUBIG

Data Science and Machine Learning

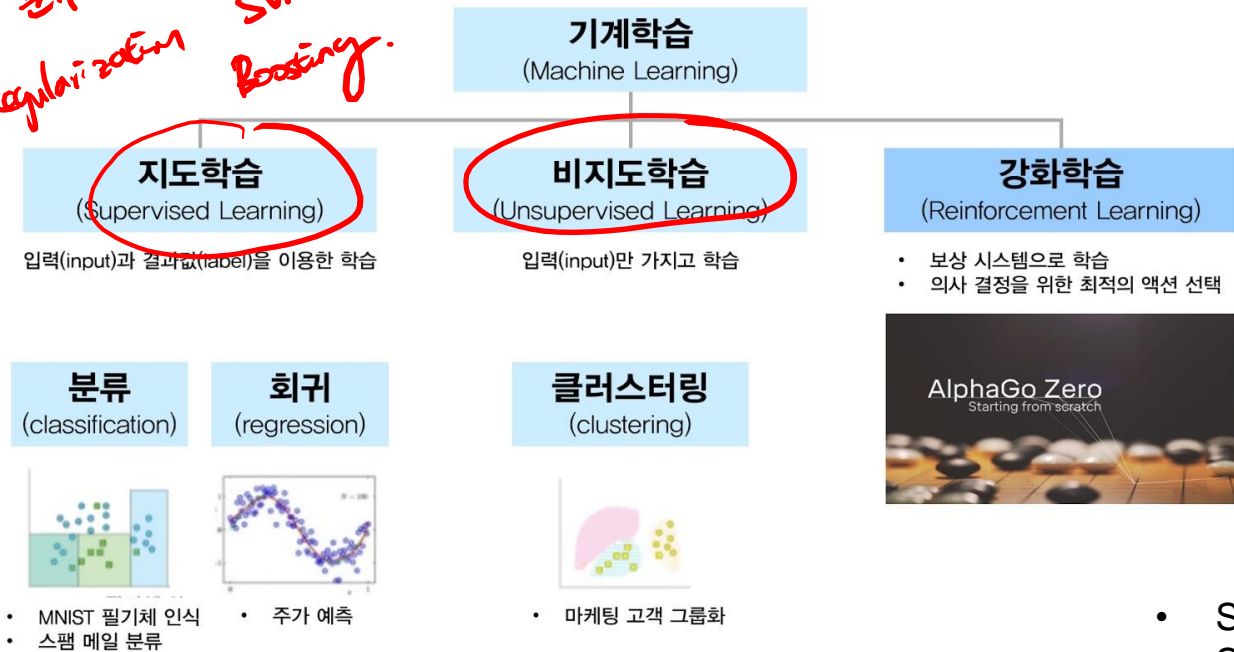
Week 7. Unsupervised Learning



Unsupervised Learning

Machine Learning

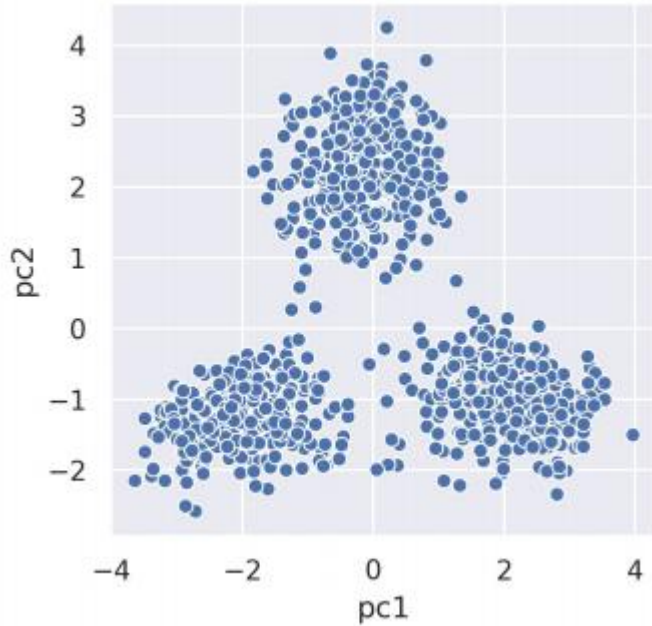
회귀 분석 회귀
Regularization SUM
Boosting.



PCA t-SNE

- Semi-Supervised Learning
- Self-Supervised Learning

Clustering

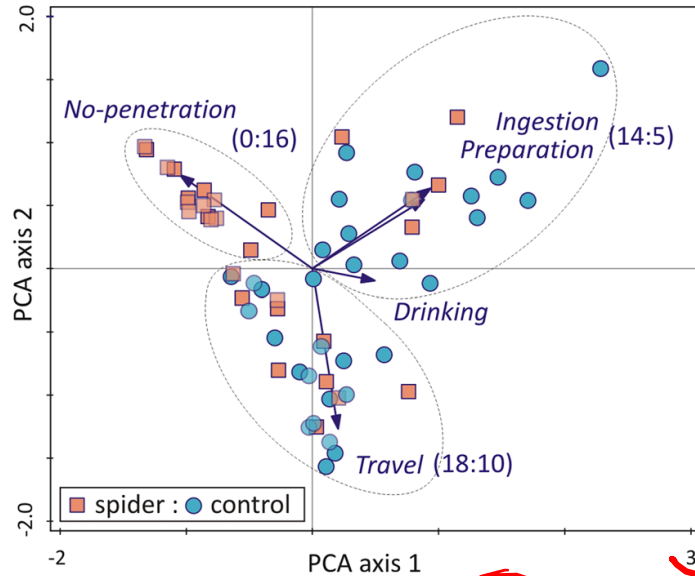


- Don't have labels for each point
- Want to infer pattern without labels

Dimensionality Reduction

Curse of Dimensionality

$(x_1, \dots, x_n) \rightarrow (z_1, z_2)$



• feature selection

• $(x_1, x_2, x_3, x_4) \rightarrow (z_1, z_2)$
 $2x_1 + x_2$ $x_3^2 + x_4$

- High dimensional data leads to high computational cost to perform learning
- Reduce correlation and complexity in data while preserving most of the relevant information in the data

PCA
SVD
Eigenvector

t-SNE

Hyper Dimensionality Reduction

Loss Function

Clustering

- ① Prototype-based Clustering
- ② Hierarchical Clustering
- ③ Density-Based Clustering.

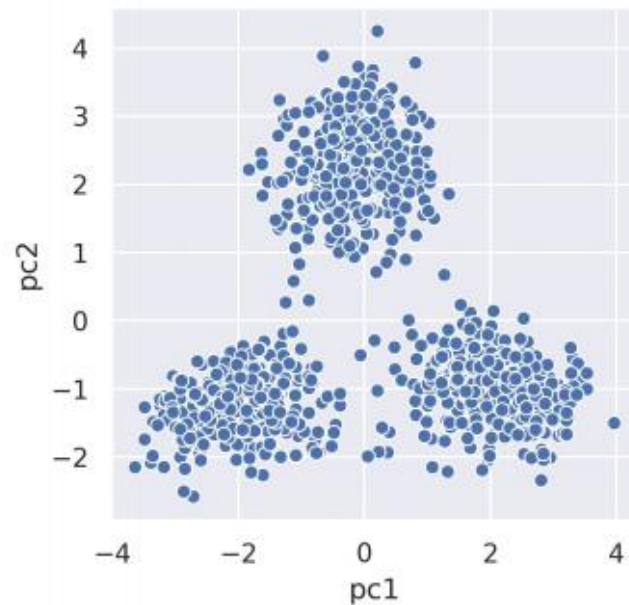
① 2nd performance metric



② Distance Metric

Validity Index

- What is the best cluster for this data?



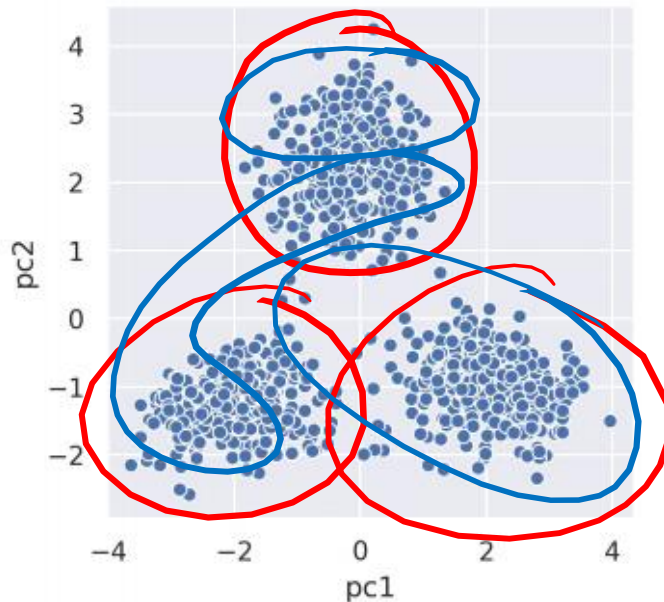
Validity Index

유효성 지수 / 성능 지수

Supervised Accuracy
Recall

- What is the best cluster for this data?

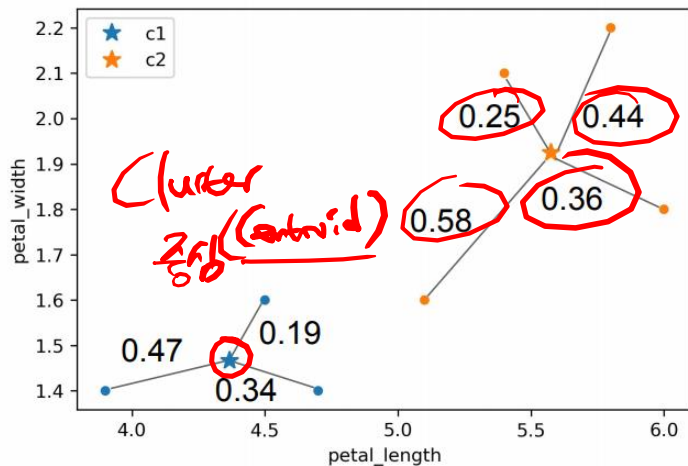
- High Intra-cluster similarity
- Low Inter-cluster similarity



Validity Index

Two common loss functions:

- ~~SSE~~ • Inertia: Sum of squared distances from each data point to its center
- Distortion: Weighted sum of squared distances from each data point to its center



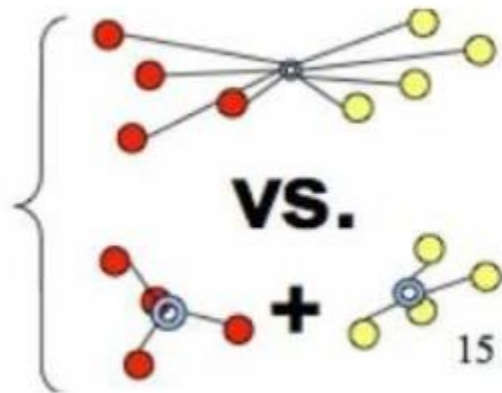
Example:

- ~~SSE~~ • Inertia: $0.47^2 + 0.19^2 + 0.34^2 + 0.25^2 + 0.58^2 + 0.36^2 + 0.44^2$
- Distortion: $(0.47^2 + 0.19^2 + 0.34^2)/3 + (0.25^2 + 0.58^2 + 0.36^2 + 0.44^2)/4$

Lower is Better

Validity Index

- Ward's Method



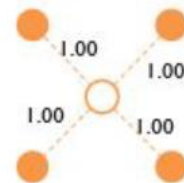
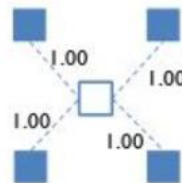
$$d(i+j, k) = \frac{\|\mu_{i+j} - \mu_k\|^2}{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$ESS = \sum_{k=1}^K \sum_{x_i \in C_k} \sum_{j=1}^n (x_{ij} - \bar{x}_{kj})^2$$

* ESS : error sum of squares
 ** k : number of clusters (1 ~ K)
 x_i : elements of cluster C_k
 j : number of variables (1 ~ n)

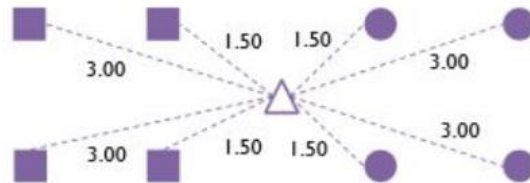
■ SSE before merge: $1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 1^2 + 1^2 = 8$

1042227



■ SSE after merge: $4 \times 1.5^2 + 4 \times 3^2 = 45$

114227



■ Ward distance: $45 - 8 = 37$

Validity Index

external index

internal index

- External Index (Using Reference Model)

Data $D = \{x_1, x_2, \dots, x_m\}$

Clusters $C = \{C_1, C_2, \dots, C_k\}$

Reference Clusters $C^* = \{C_1^*, C_2^*, \dots, C_s^*\}$

Cluster Labels: λ, λ^*

$$\begin{aligned} a &= |SS|, & SS &= \{(x_i, x_j) \mid \lambda_i = \lambda_j, \lambda_i^* = \lambda_j^*, i < j\} \\ b &= |SD|, & SD &= \{(x_i, x_j) \mid \lambda_i = \lambda_j, \lambda_i^* \neq \lambda_j^*, i < j\} \\ c &= |DS|, & DS &= \{(x_i, x_j) \mid \lambda_i \neq \lambda_j, \lambda_i^* = \lambda_j^*, i < j\} \\ d &= |DD|, & DD &= \{(x_i, x_j) \mid \lambda_i \neq \lambda_j, \lambda_i^* \neq \lambda_j^*, i < j\} \end{aligned}$$

★ Jaccard Coefficient

$$JC = \frac{a}{a + b + c}$$

CV
(Iou)

Fowlkes and Mallows Index

$$FMI = \sqrt{\frac{a}{a+b} * \frac{a}{a+c}}$$

Rand Index

$$RI = \frac{2(a+d)}{m(m-1)}$$

Validity Index

Model Score

Data $D = \{x_1, x_2, \dots, x_m\}$

Clusters $C = \{C_1, C_2, \dots, C_k\}$

- Internal Index

$$avg(C) = \frac{2}{|C|(|C| - 1)} \sum_{1 \leq i < j \leq |C|} dist(x_i, x_j)$$

$$diam(C) = \max_{1 \leq i < j \leq |C|} dist(x_i, x_j)$$

$$d_{min}(C_i, C_j) = \min_{x_i \in C_i, x_j \in C_j} dist(x_i, x_j)$$

$$d_{cen}(C_i, C_j) = dist(\mu_i, \mu_j)$$

$$\mu = \frac{1}{|C|} \sum_{1 \leq i \leq |C|} x_i$$

Daives-Bouldin Index

$$DBI = \frac{1}{k} \sum_{i=1}^k \max_{j \neq i} \left(\frac{avg(C_i) + avg(C_j)}{d_{cen}(C_i, C_j)} \right)$$

Dunn Index

$$DI = \frac{1}{|C|} \min_{1 \leq i \leq k} \left\{ \min_{j \neq i} \left(\frac{d_{min}(C_i, C_j)}{\max_{1 \leq l \leq k} diam(C_l)} \right) \right\}$$

Distance Measure

거리
측정

거리
(kg)

2차원 Euclidean

- Distance measure

7 (mm)

p=2

$$d(\mathbf{u}, \mathbf{v}) = (\sum |u_i - v_i|^2)^{\frac{1}{2}} = ||\mathbf{u} - \mathbf{v}||_2$$

Euclidean (L2 norm)

p=1

$$d(\mathbf{u}, \mathbf{v}) = \sum |u_i - v_i| = ||\mathbf{u} - \mathbf{v}||_1$$

Manhattan (L1 norm)

$$d(\mathbf{u}, \mathbf{v}) = (\sum |u_i - v_i|^p)^{\frac{1}{p}} = ||\mathbf{u} - \mathbf{v}||_p$$

Minkowski (Lp norm)

$$d(\mathbf{u}, \mathbf{v}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

Mahalanobis Distance

거리 측정

1-2 cor.



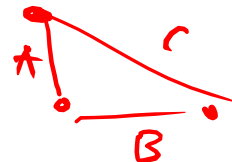
Distance Measure

노출, 불안, 보통, 만족, 매우 만족.
 1 2 3 4 5
 1 1

- Distance measure

특성변수	비유사성 측도	
연속형	제곱 유클리디안 거리 (squared Euclidean distance)	$d(x_i, x_{i'}) = \sum_{j=1}^m (x_{ij} - x_{i'j})^2 = (x_i - x_{i'})^T (x_i - x_{i'})$
	1차 유클리디안 거리 (L1 Euclidean distance)	$d(x_i, x_{i'}) = \sum_{j=1}^m x_{ij} - x_{i'j} $
	마할라노비스 거리 (Mahalanobis distance)	$d(x_i, x_{i'}) = (x_i - x_{i'})^T \Sigma^{-1} (x_i - x_{i'})$ (Σ : 특성변수들의 분산-공분산 행렬)
순서형 (ordinal)	특성변수를 $\frac{k-1}{M}$ ($k = 1, \dots, M$; 순서의 크기)로 변환 후 연속형 비유사성 측도 적용	
범주형 (categorical)	두 관측치가 같은 범주에 속하면 0, 아니면 1로 값 부여	

연속 숫자
(숫자)



distance

VDM

dummy variable

Prototype – based Clustering

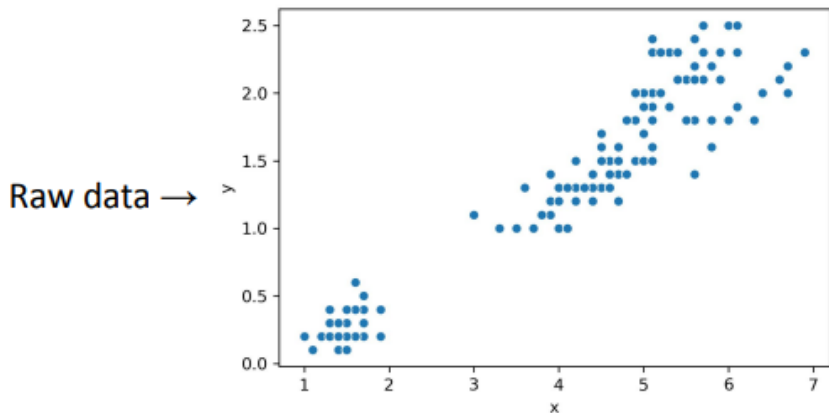
Centroid

K-Means Clustering

Most popular clustering approach: K-Means

$K = \text{클러스터의 개수}$

- Pick an arbitrary k , and randomly place k “centers”, each a different color
- Repeat until convergence:
 - Color points according to the closest center
 - Move center for each color to center of points with that color

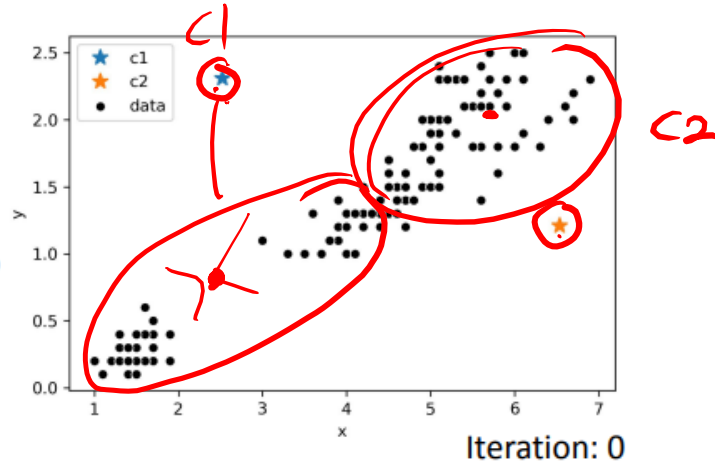


K-Means Clustering

Most popular clustering approach: K-Means.

- Pick an arbitrary k , and **randomly place k “centers”**, each a different color
- Repeat until convergence:
 - Color points according to the closest center
 - Move center for each color to center of points with that color

Initial random
placement of two
centers

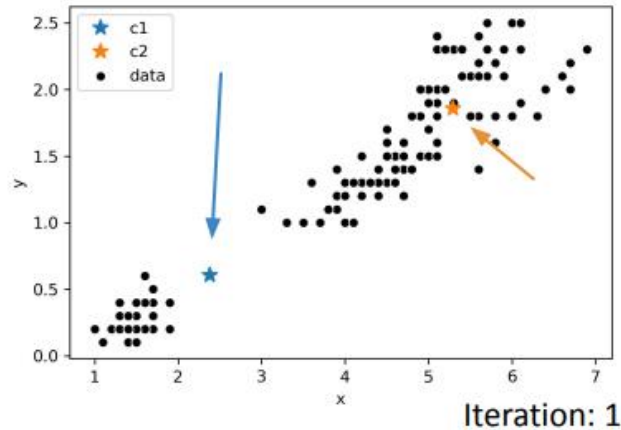


K-Means Clustering

Most popular clustering approach: K-Means

- Pick an arbitrary k , and randomly place k “centers”, each a different color
- Repeat until convergence:
 - Color points according to the closest center
 - **Move center for each color to center of points with that color**

Centers moved to
their new homes

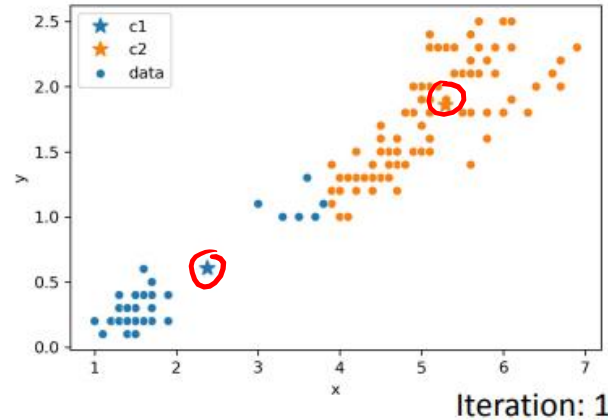


K-Means Clustering

Most popular clustering approach: K-Means

- Pick an arbitrary k , and randomly place k “centers”, each a different color
- Repeat until convergence:
 - **Color points according to the closest center**
 - Move center for each color to center of points with that color

Data colored by
closest center (in
new position)

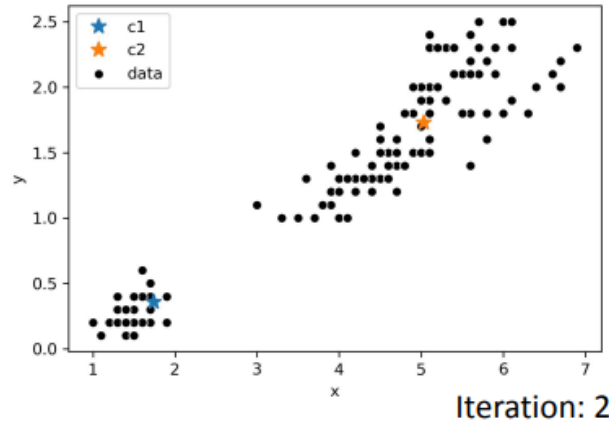


K-Means Clustering

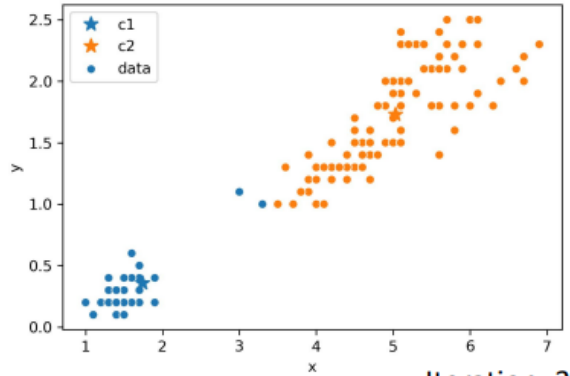
Most popular clustering approach: K-Means

- Pick an arbitrary k , and randomly place k “centers”, each a different color
- Repeat until convergence:
 - Color points according to the closest center
 - **Move center for each color to center of points with that color**

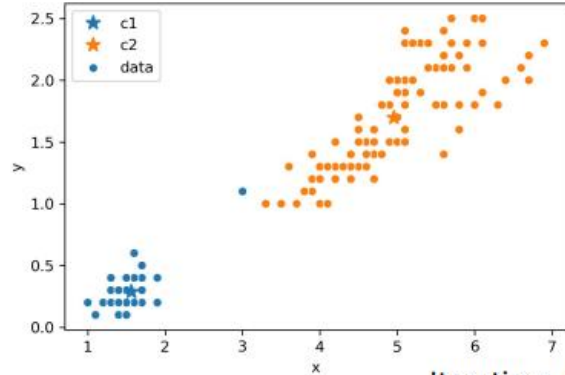
Centers moved to
new position



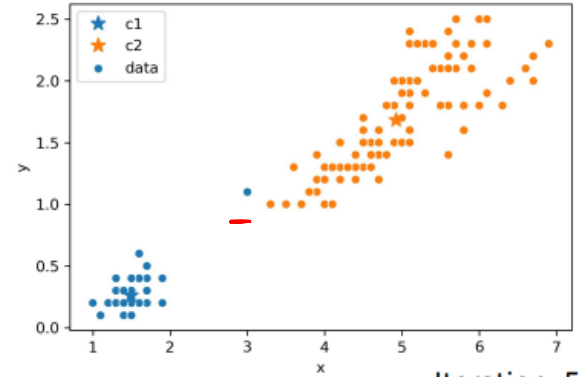
K-Means Clustering



Iteration: 3



Iteration: 4



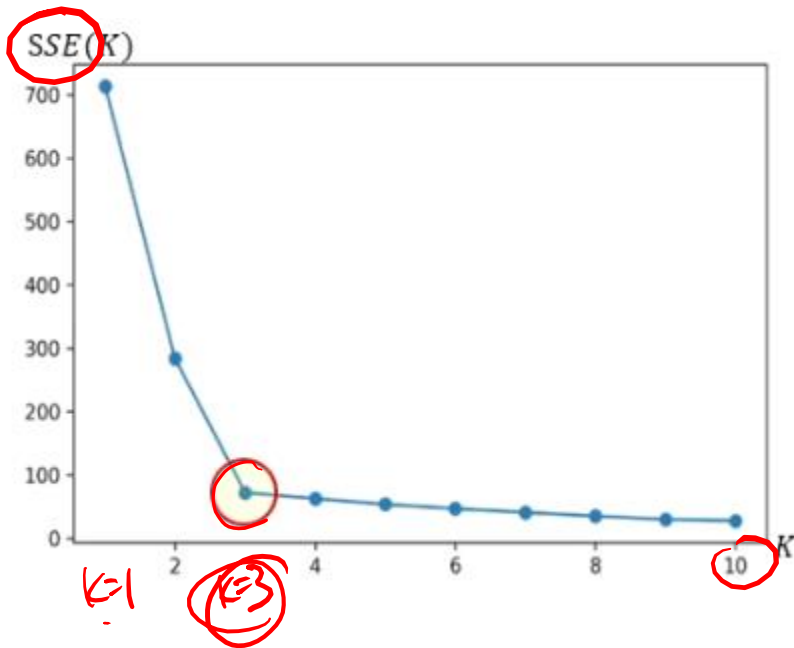
Iteration: 5

K-Means Clustering

$$S_l^2 = \sum_{i \in l} \sum_{j=1}^d ((x_{ij} - \mu_{lj})^2$$

$$SSE(K) = \sum_{l=1}^K S_l^2$$

K-Means Clustering



▶ K 가 커질수록 $SSE(K)$ 감소

▶ $K=3$ 에서 $SSE(K)$ 의 감소속도(기울기)가 현저하게 줄어듦

▶ $K=3$ 선택

→ "Elbow method"

K-Means Clustering

- Downfalls of K-Means Clustering
 - Output is inconsistent
 - Performance depends on hyperparameter K

c

K-Means ++

1. N개의 표본으로 구성된 학습데이터로부터 1개의 임의표본을 뽑는다. 이를 초기 중심값 μ_1 으로 한다.
2. 나머지 (n-1) 개의 자료에 대해 μ_1 으로부터의 유클리디안 거리를 구하고 이 거리에 비례한 확률을 (n-1)개의 자료 각각에 부여한 후, 이 확률에 비례하여 1개의 임의표본을 뽑는다. 이를 두번째 초기 중심값 μ_2 로 한다.
3. (n-2)개의 자료에 대해 $\min(d(x_i, \mu_1), d(x_i, \mu_2))$ 를 구해 이 최소거리에 비례한 확률을 (n-2)개 자료 각각에 부여한 후, 이 확률에 비례하여 1개의 임의표본을 뽑아 이를 두번째 초기 중심값 μ_3 로 한다.
4. 앞의 방법을 계속 반복

K-Means ++

1. N개의 표본으로 구성된 학습데이터로부터 1개의 임의표본을 뽑는다. 이를 초기 중심값 μ_1 으로 한다.
2. 나머지 (n-1) 개의 자료에 대해 μ_1 으로부터의 유클리디안 거리를 구하고 이 거리에 비례한 확률을 (n-1)개의 자료 각각에 부여한 후, 이 확률에 비례하여 1개의 임의표본을 뽑는다. 이를 두번째 초기 중심값 μ_2 로 한다.
3. (n-2)개의 자료에 대해 $\min(d(x_i, \mu_1), d(x_i, \mu_2))$ 를 구해 이 최소거리에 비례한 확률을 (n-2)개 자료 각각에 부여한 후, 이 확률에 비례하여 1개의 임의표본을 뽑아 이를 두번째 초기 중심값 μ_3 로 한다.
4. 앞의 방법을 계속 반복

Other Prototype Based Clustering

이이 클러스터가 있다고 가정

- Learning Vector Quantization (LVQ)
- Mixture-of-Gaussian Clustering

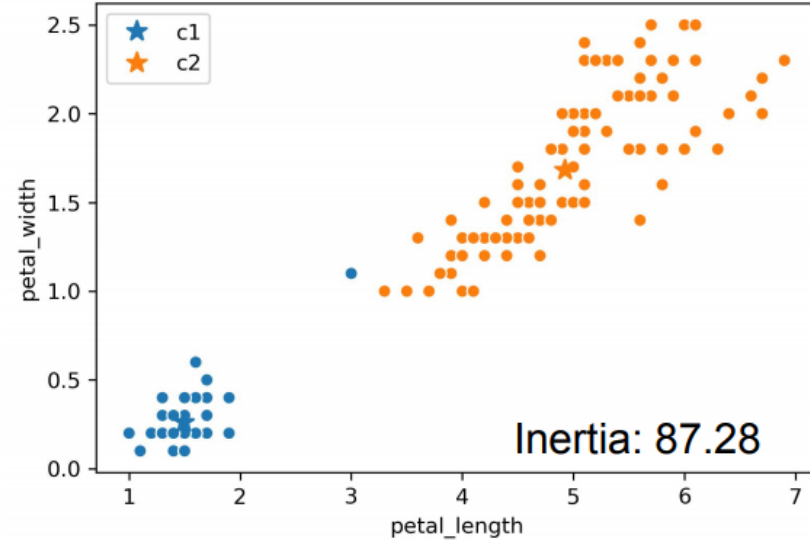
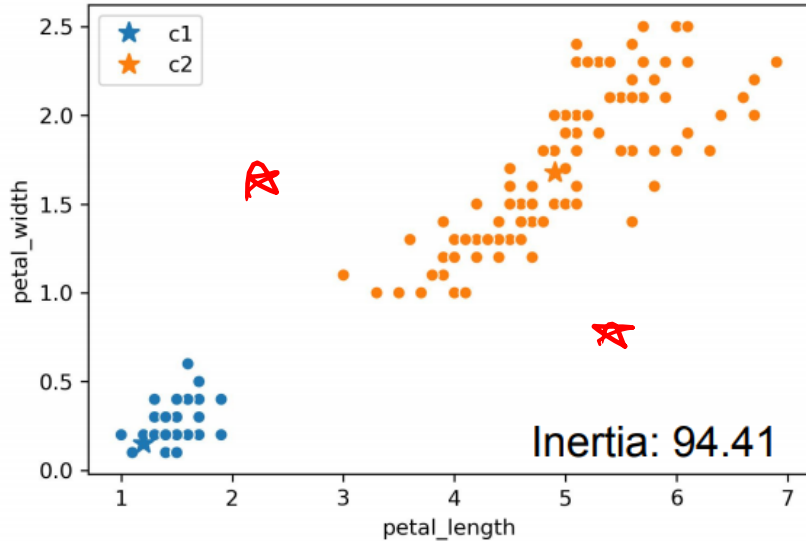
- Centroid
- 본토

(본토) \leftarrow EM 알고리즘

Hierarchical Clustering

다변량 통계분석

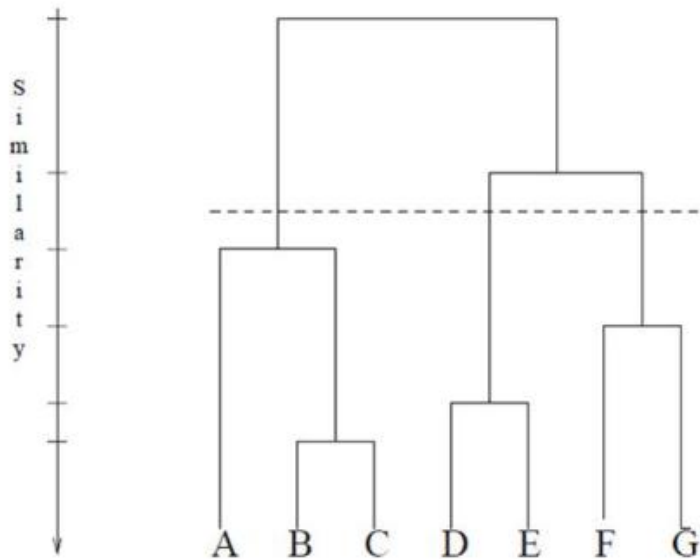
Problems of Prototype – Based Clustering



Hierarchical Clustering

진돗개, 셰퍼드, 요크셔테리어, 푸들, 몰소, 짚소

인들레



- Direction

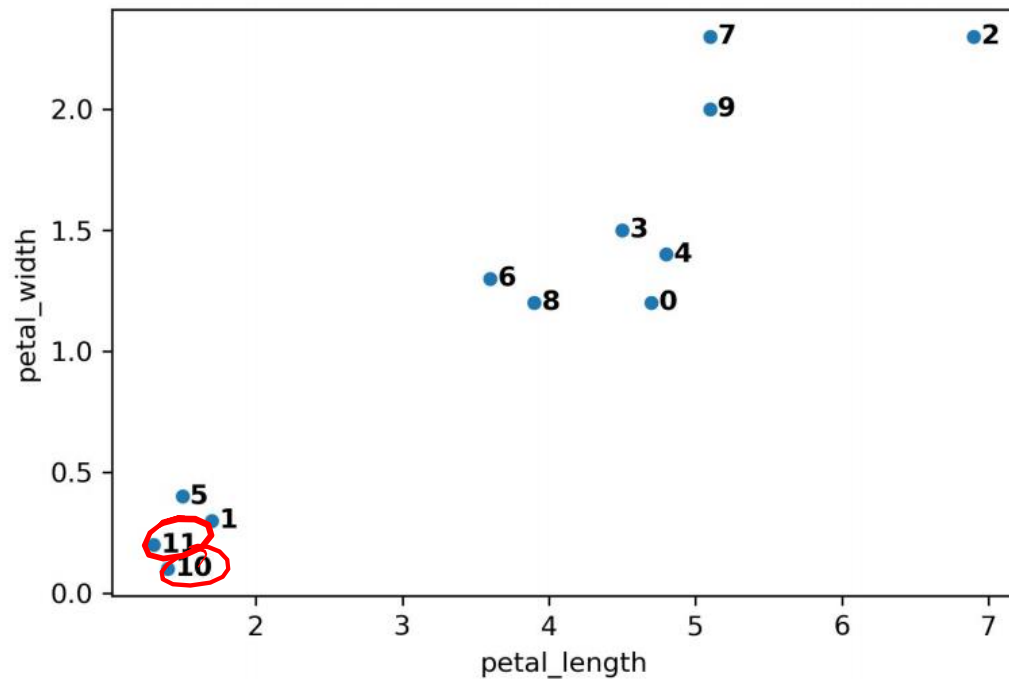
Agglomerative Clustering (Bottom Up Approach)

- Divisive Clustering

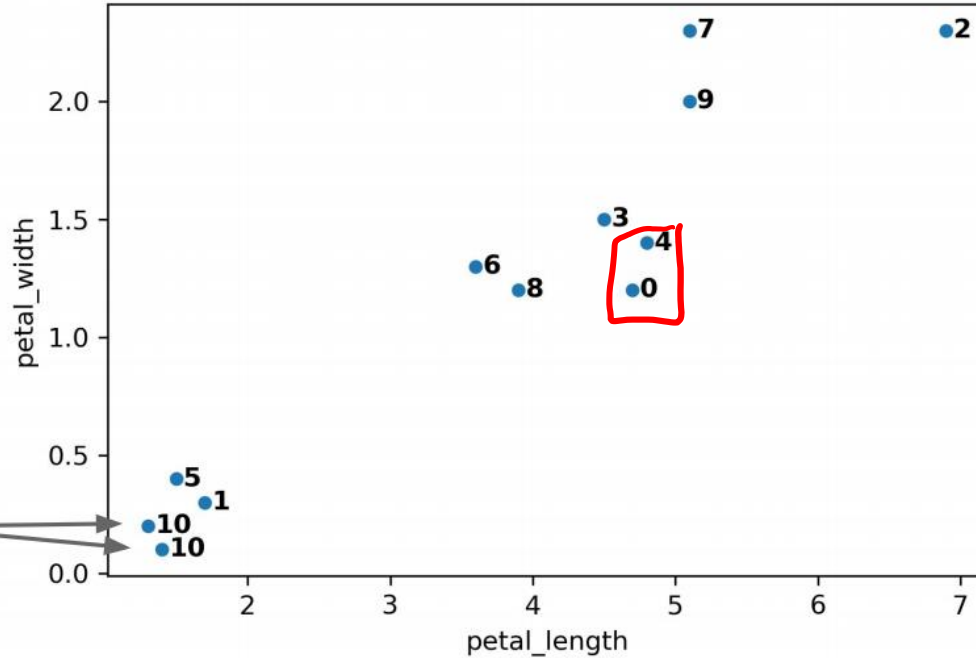
- Distance

- Complete-Linkage Clustering
 - Single-Linkage Clustering
 - Average-Linkage Clustering

Agglomerative Clustering

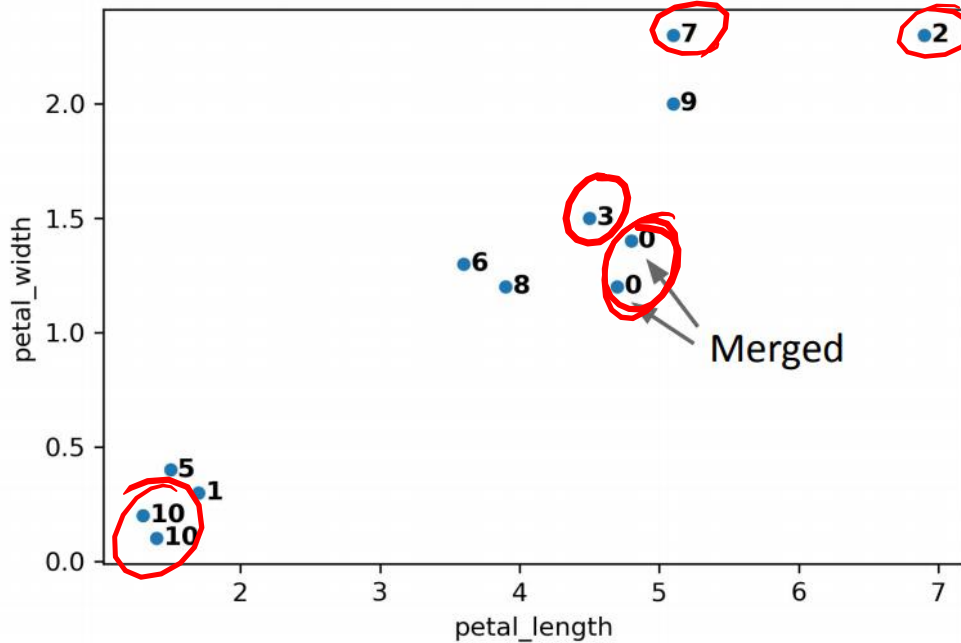


Agglomerative Clustering

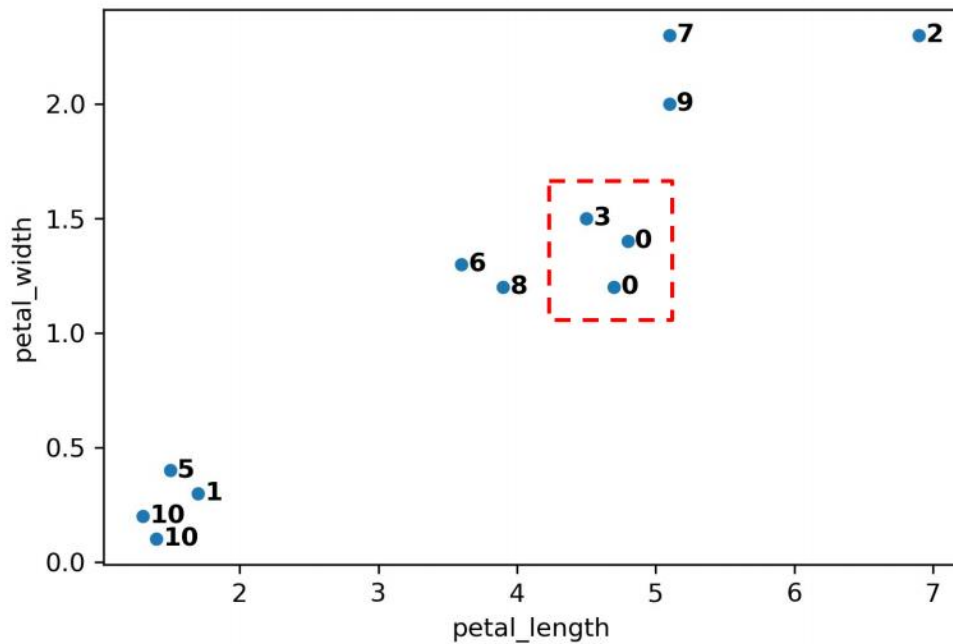


Two points in
the same cluster

Agglomerative Clustering



Agglomerative Clustering



Agglomerative Clustering

Cluster A: $\mathbf{x}_1, \dots, \mathbf{x}_{n_A}$
Clusters B : $\mathbf{z}_1, \dots, \mathbf{z}_{n_B}$

- Distance between Clusters

Complete – Linkage Clustering

$$\max\{d(\mathbf{x}_i, \mathbf{z}_j); i = 1, \dots, n_A, j = 1, \dots, n_B\}$$

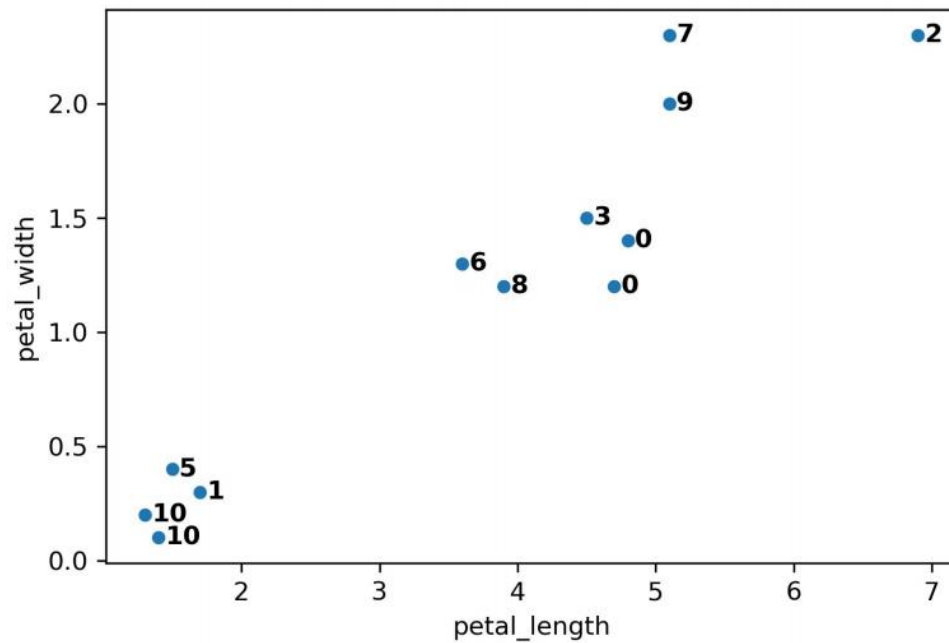
Single – Linkage Clustering

$$\min\{d(\mathbf{x}_i, \mathbf{z}_j); i = 1, \dots, n_A, j = 1, \dots, n_B\}$$

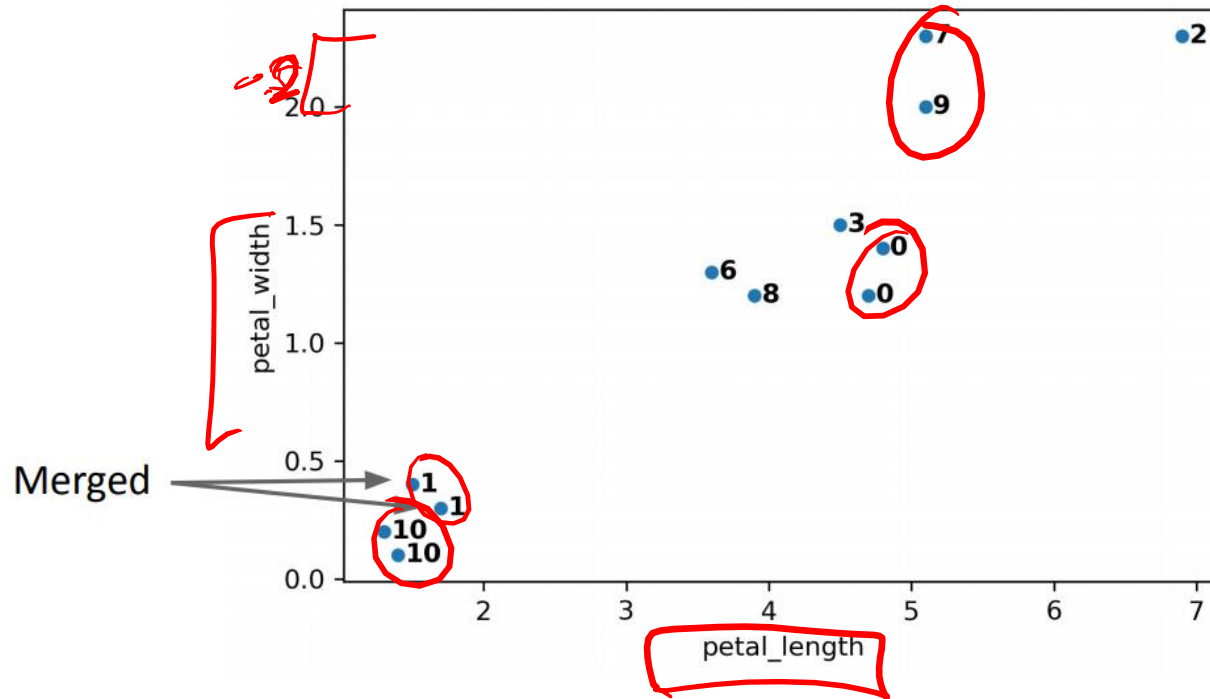
Average – Linkage Clustering

$$\frac{1}{n_A n_B} \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} d(\mathbf{x}_i, \mathbf{z}_j)$$

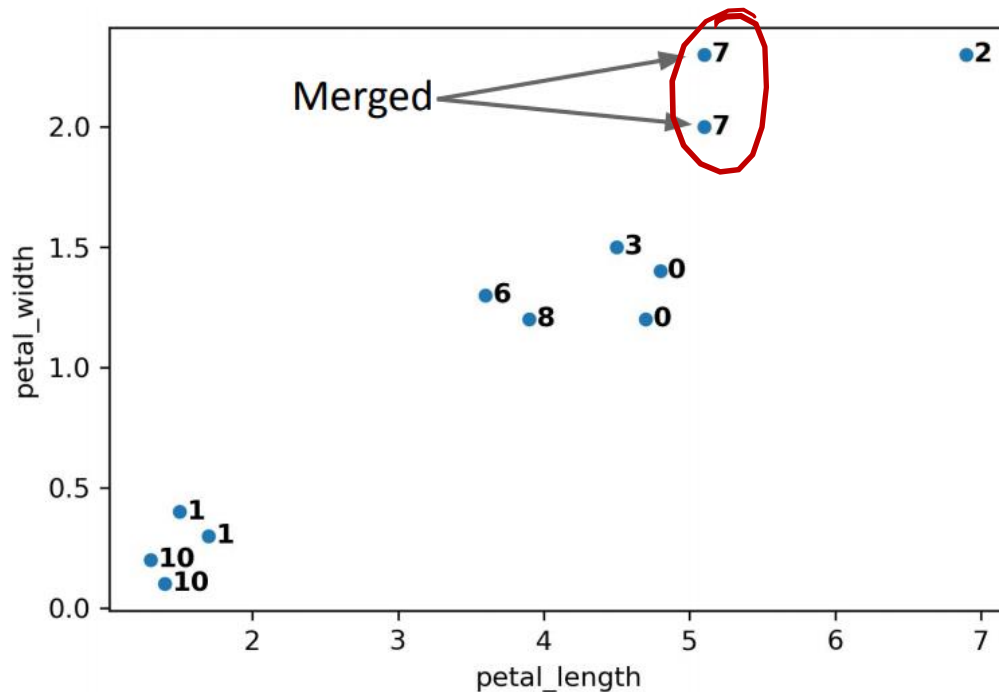
Agglomerative Clustering



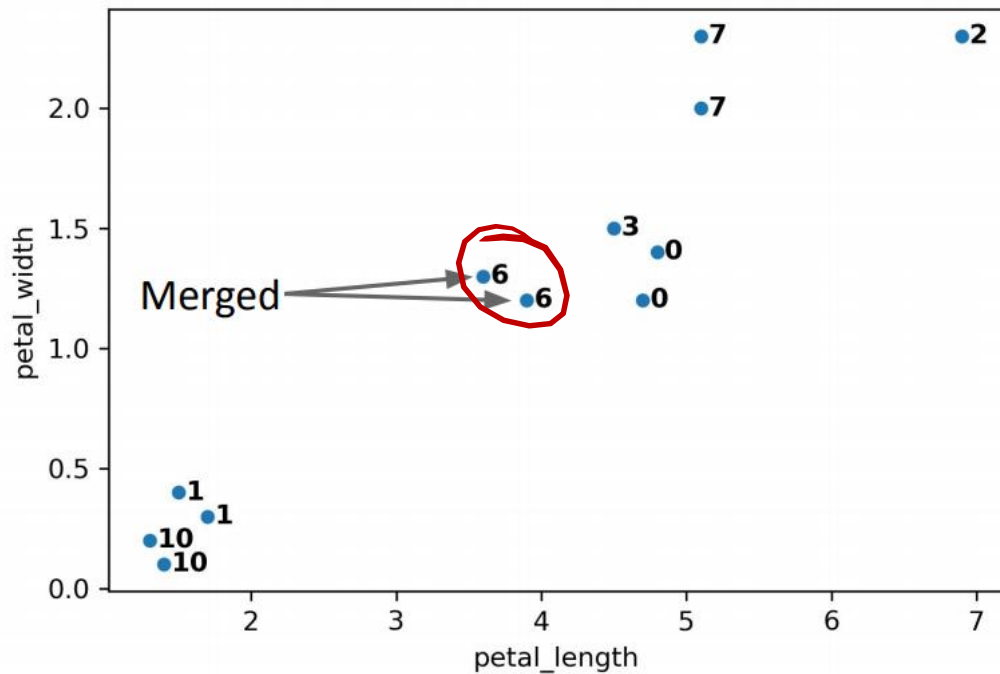
Agglomerative Clustering



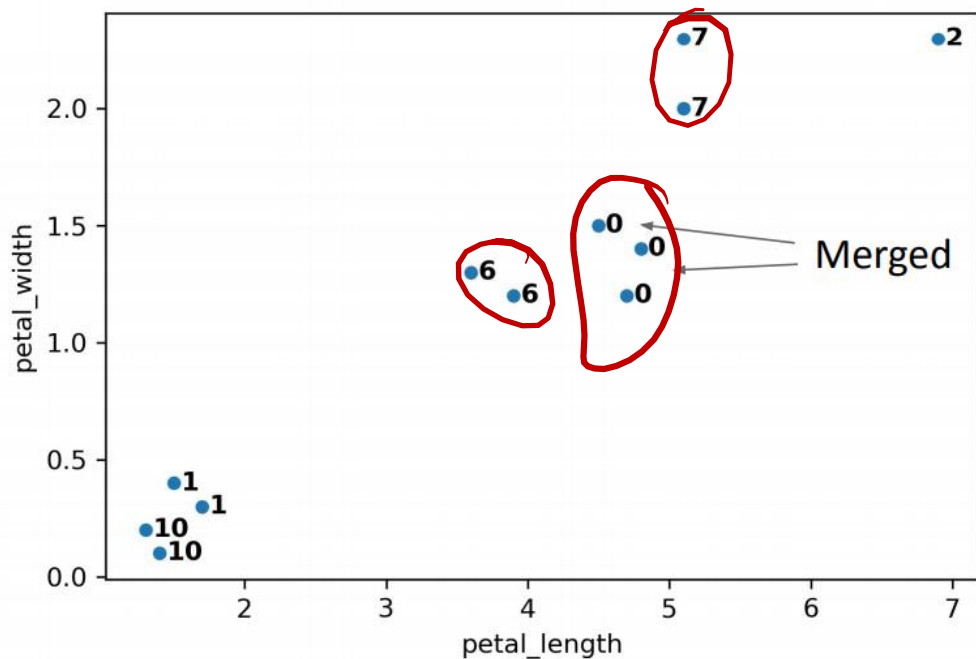
Agglomerative Clustering



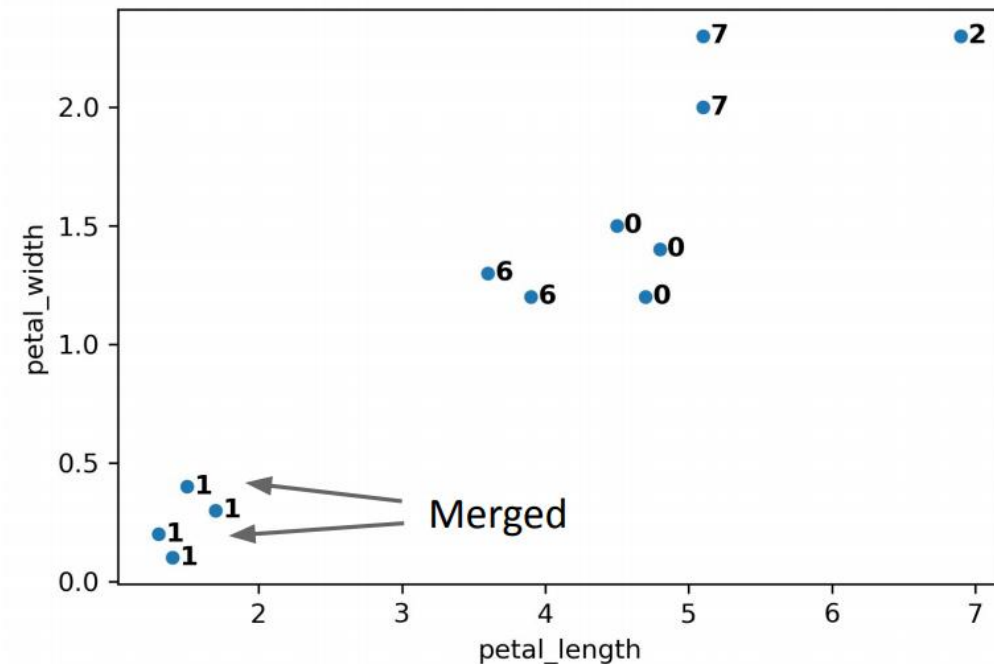
Agglomerative Clustering



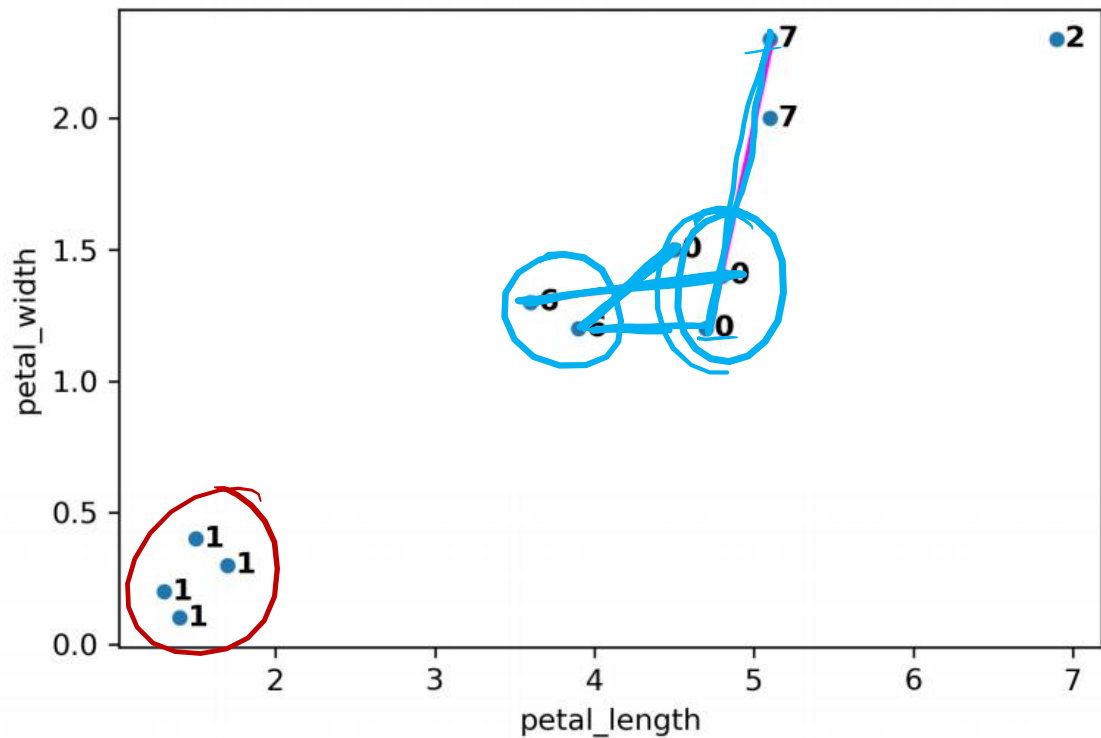
Agglomerative Clustering



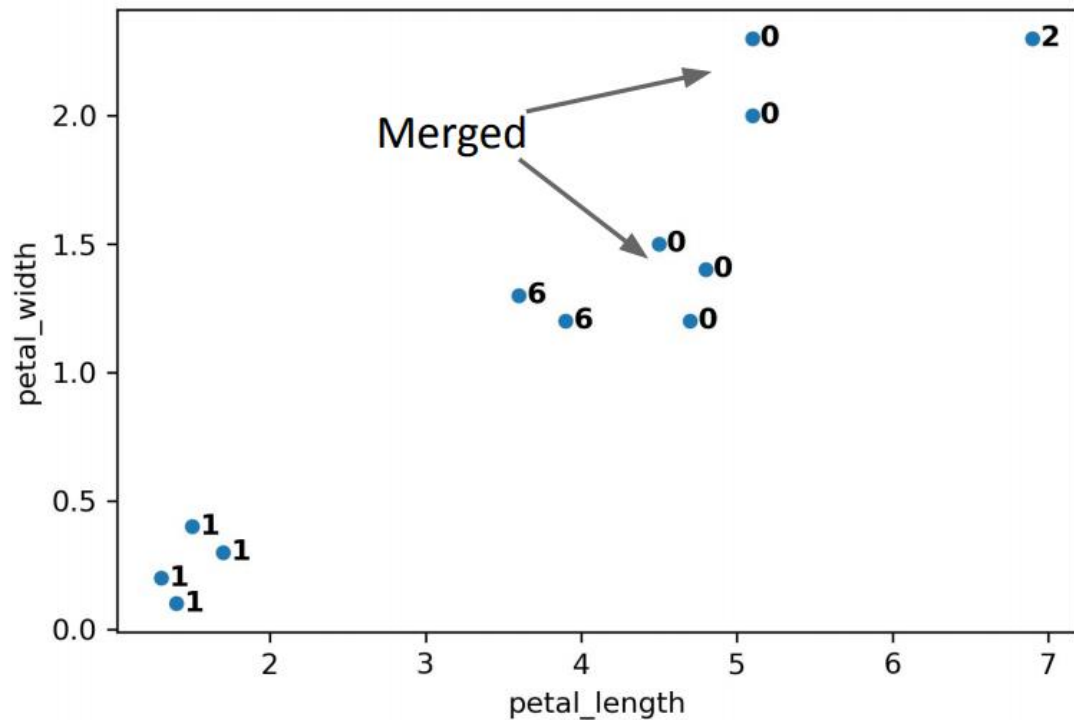
Agglomerative Clustering



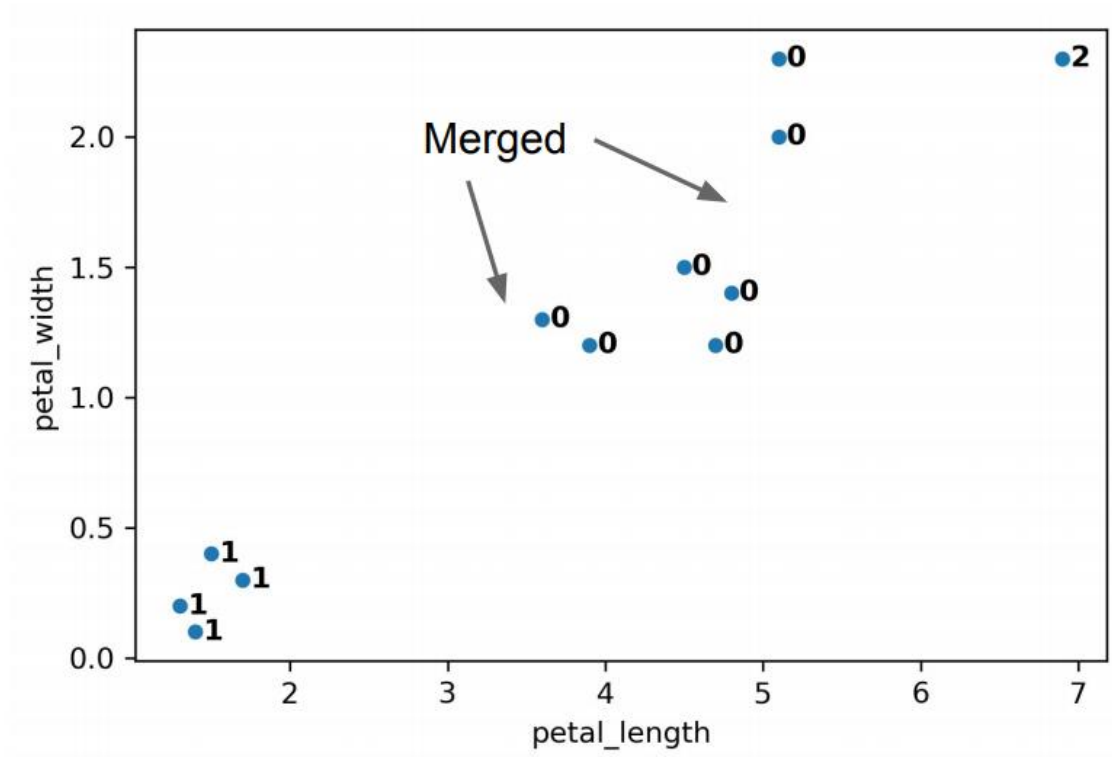
Agglomerative Clustering



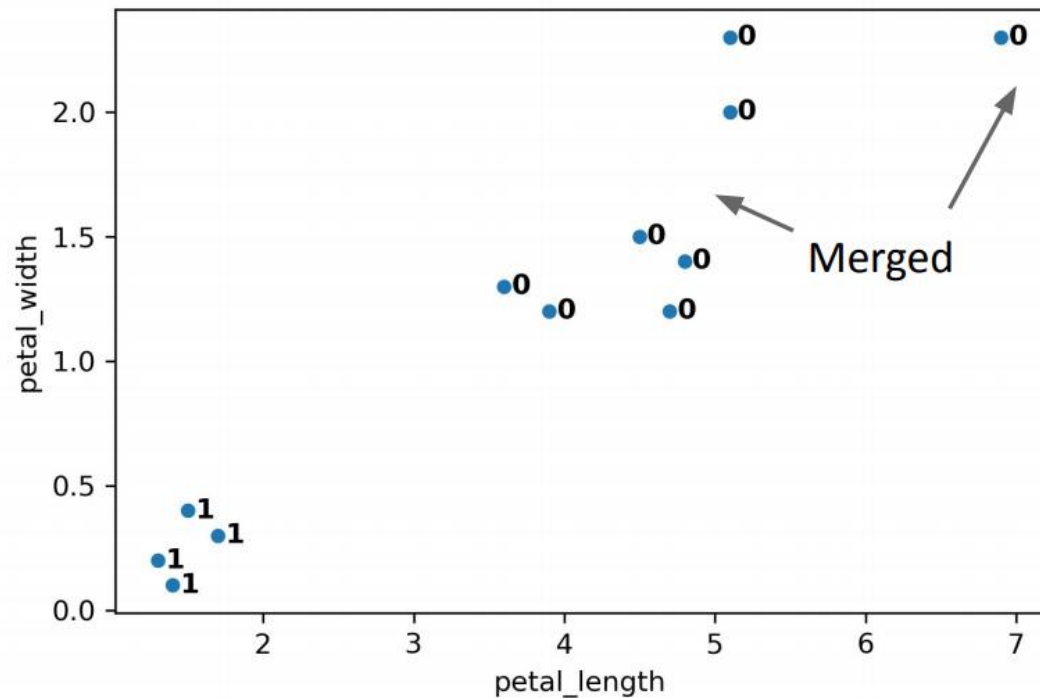
Agglomerative Clustering



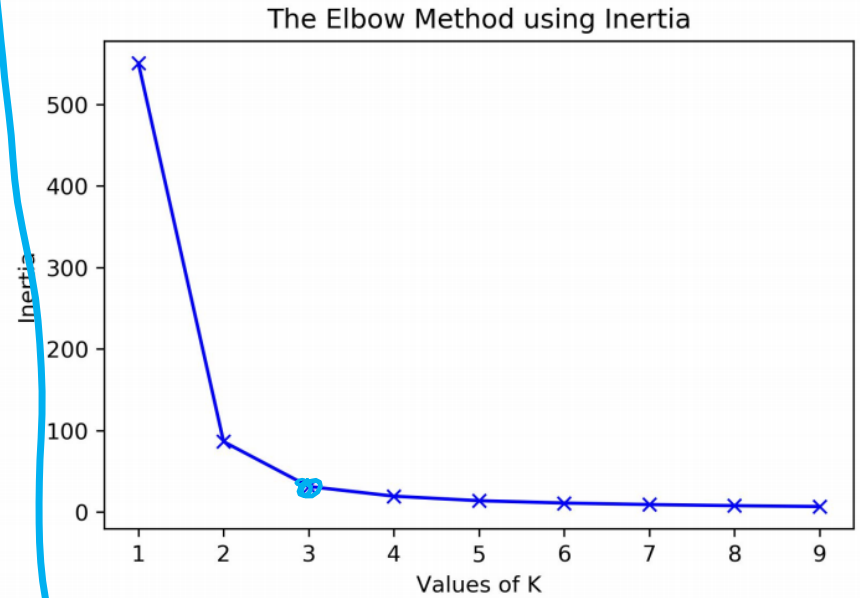
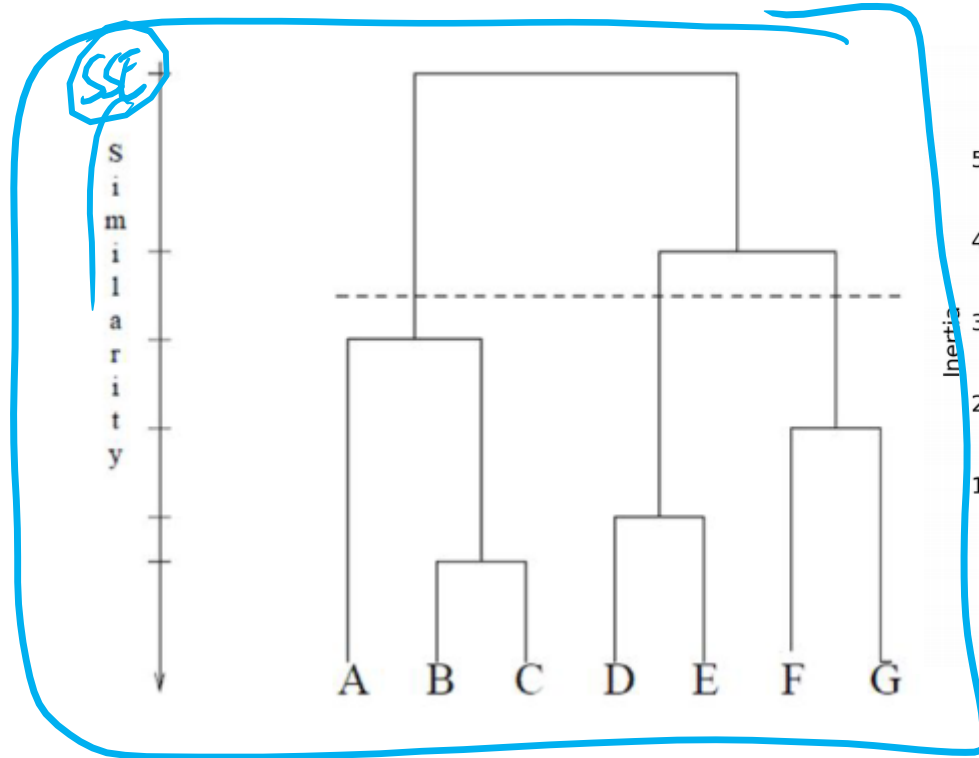
Agglomerative Clustering



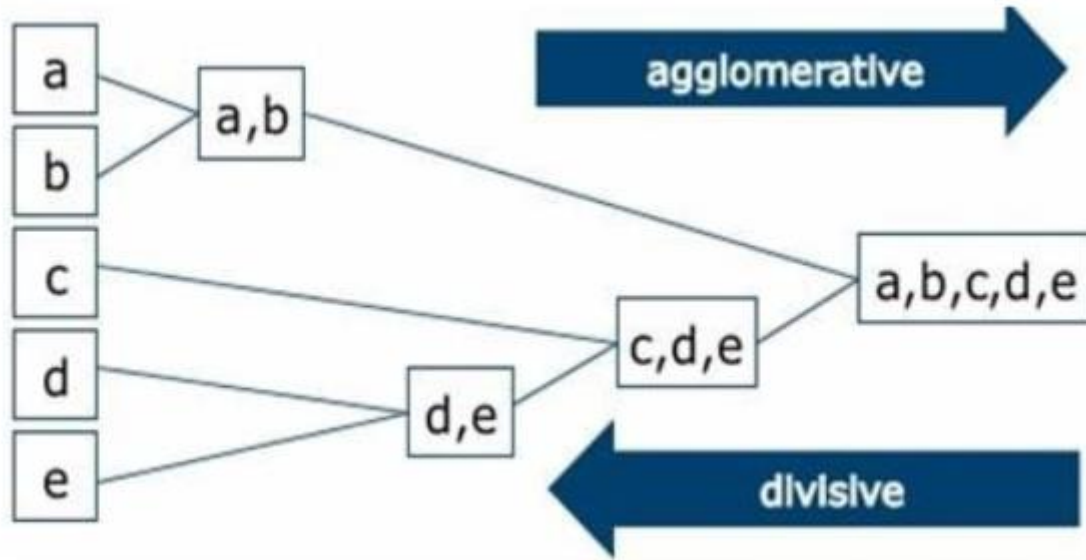
Agglomerative Clustering



Agglomerative Clustering



Agglomerative Clustering and Divisive Clustering



Downfalls of Hierarchical Clustering

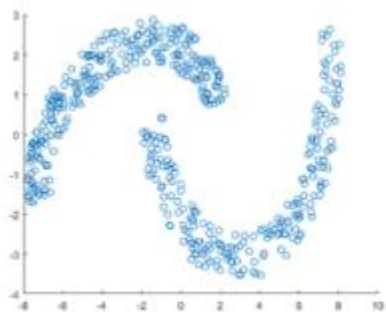
★ Noise와 Outlier에 민감하다 (Single Linkage)

- 한번 두 cluster를 묶으면 그 결정을 되돌릴 수 없다.

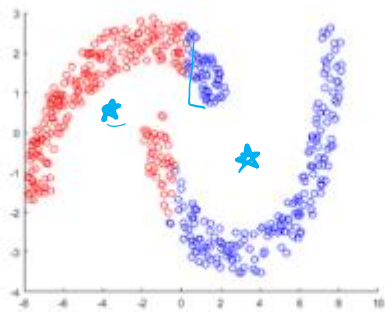


Density - based Clustering

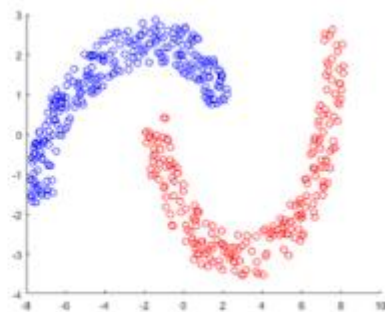
Limitations of Distance-based Clustering



(a) 원본 데이터



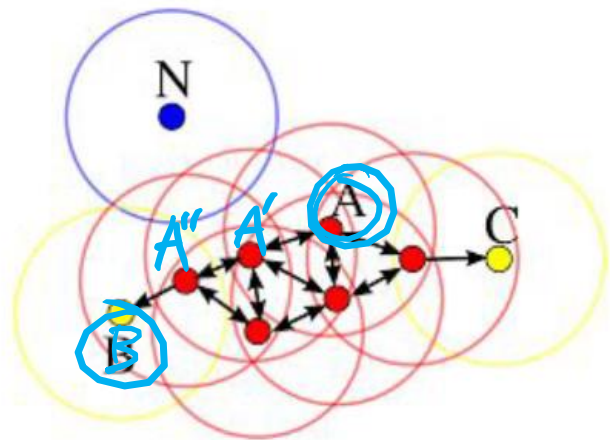
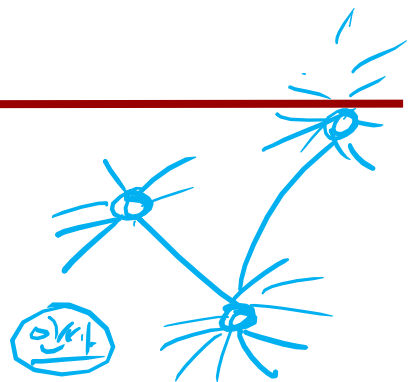
(b) k-means clustering의 결과



(c) DBSCAN의 결과

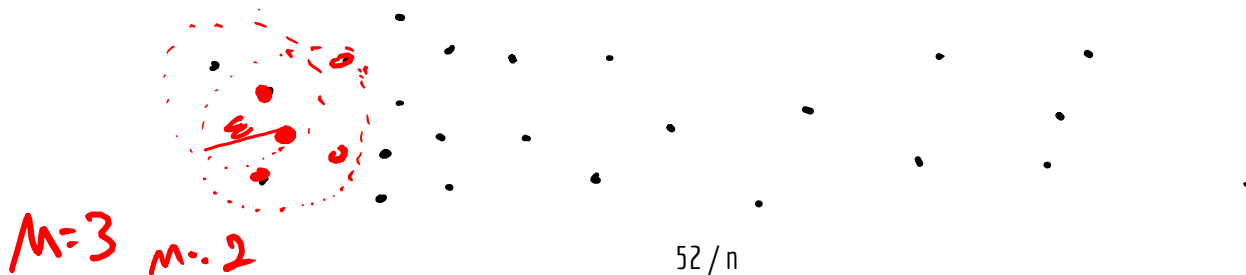
DBSCAN: Density-Based Spatial Clustering of Approximations with Noise

- Cluster: 접근 가능 관계로 유도해 낸 최대의 밀도 연결 샘플 집합
- Hyperparameter
 - M: Minimum # of samples in a neighborhood for core point
 - ϵ - Nearest Neighbors x와의 거리가 ϵ 보다 작은 샘플들

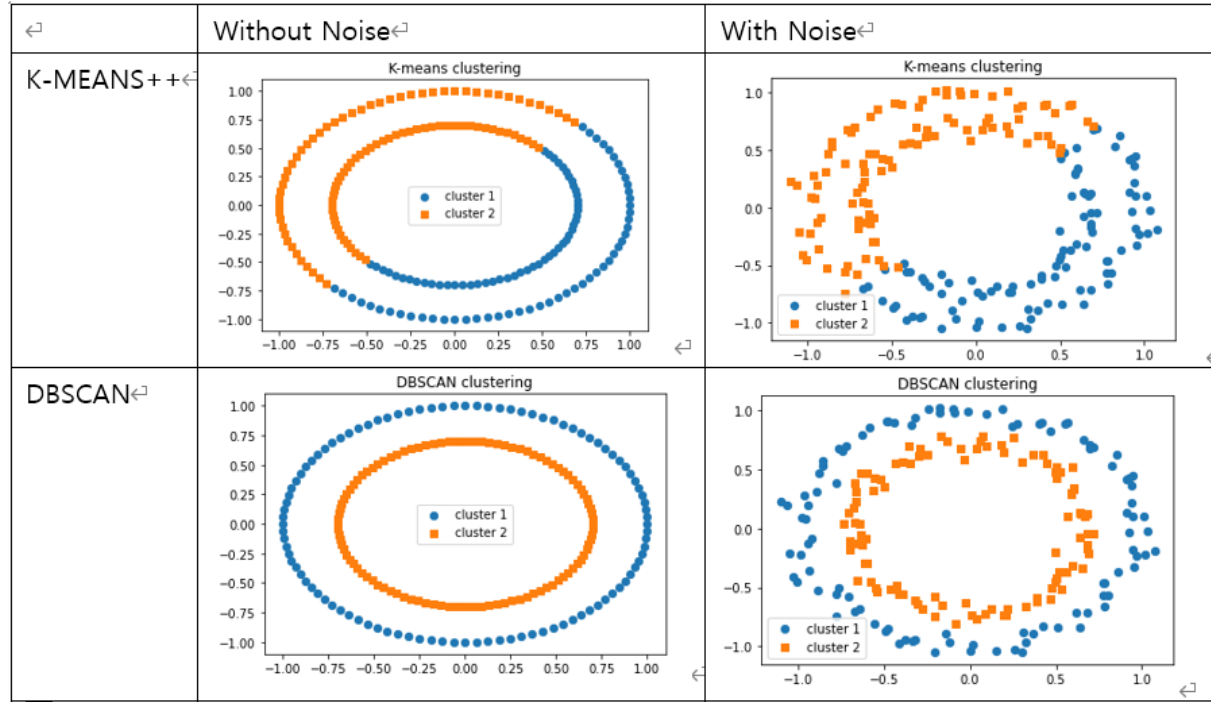


DBSCAN: Density-Based Spatial Clustering of Approximations with Noise

1. 임의 관측치를 선택하고, 군집 1 부여
2. 이 관측치의 $\epsilon - NN$ 을 구하고, $\epsilon - NN < M$ 이면 Noise나 Outlier로 표시
3. $\epsilon - NN \geq M$ 이면 모두에게 군집 1을 할당
4. 새로 군집 1이 부여된 NN에 대하여 그들의 $\epsilon - NN$ 이 M 보다 크면 이들 $\epsilon - NN$ 도 군집 1을 할당한다.
5. 군집 1의 어느 관측치도 M개 이상의 $\epsilon - NN$ 가 존재하지 않을 때까지 반복한다.
6. 군집 2를 형성하기 위해 (2) ~ (5) 반복
7. 모든 관측치가 군집 소속으로 분류되거나 잡음으로 분류될 때까지 반복



DBSCAN: Density-Based Spatial Clustering of Approximations with Noise



Other ^{Density} ~~Prototype~~ Based Clustering

- HDBSCAN: Hierarchical DBSCAN)
- OPTICS
- DENCLUE

DBSCAN & HDBSCAN