데이터 과학 기초

05

# 군집분석

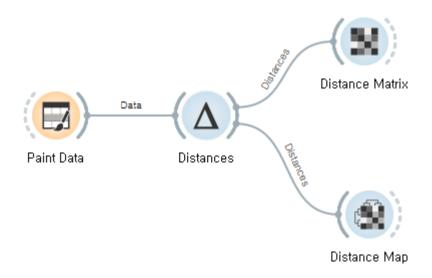
경북대학교 배준현 교수 (joonion@knu.ac.kr)

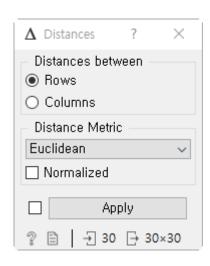


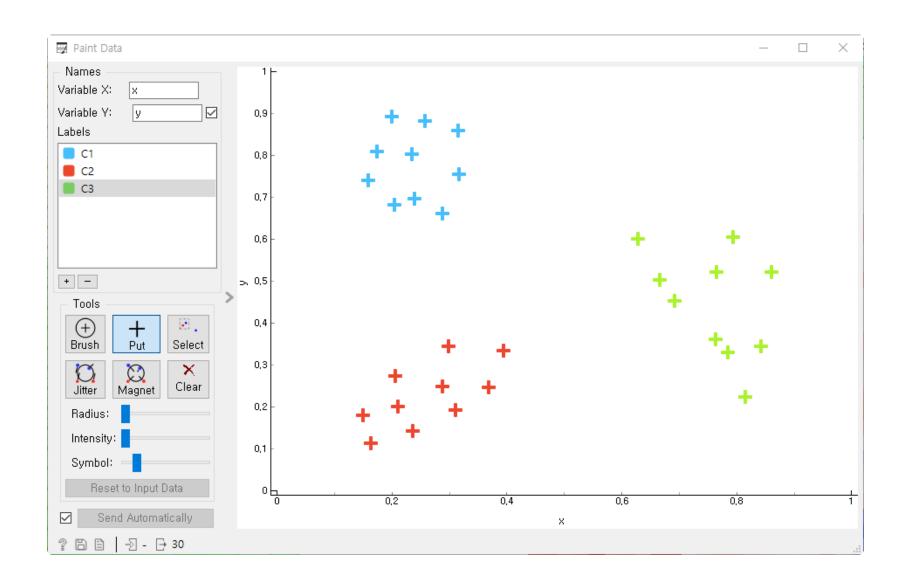
#### ○ O5. 군집 분석

- 군집화: Clustering
  - 데이터를 서로 비슷하거나, 서로 가까운 것들끼리 묶어서 나누는 것
  - 비지도 학습: 정답(라벨)이 없음. 유사도와 거리로 판단
    - 유사도(*similarity*): 두 데이터가 얼마나 가까운가를 나타내는 척도
    - 거리(*distance*): 두 데이터 사이의 거리
    - -s = 1 d, s = similarity, d = distance

Orange: Distance



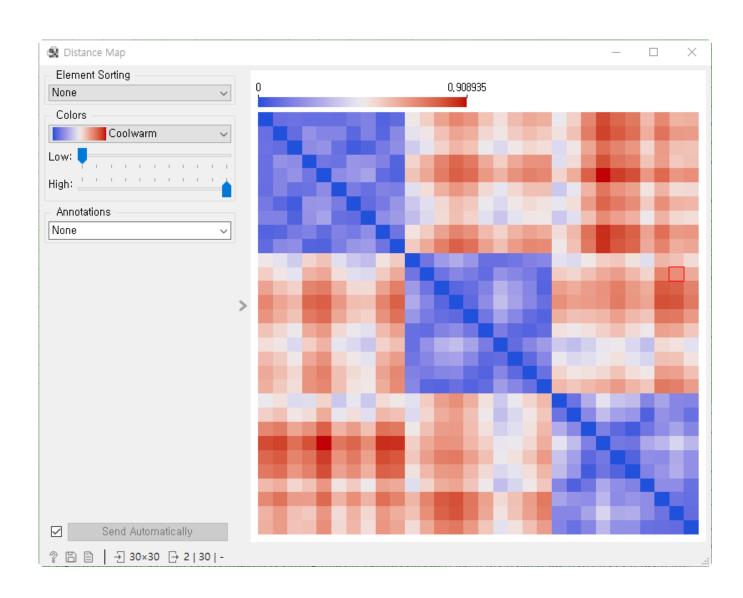






#### 🔊 05. 군집 분석







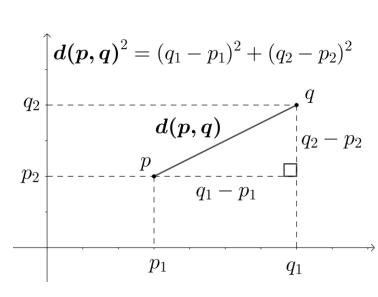
- 유클리드 거리: *Euclidean* Distance
  - 두 점 사이의 거리를 계산하는 가장 일반적인 방법
  - 유클리드 공간에서 두 점 사이의 거리

- 2차원: 
$$d = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2}$$

- 3차원: 
$$d = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + (p_3 - q_3)^2}$$

- 
$$n$$
차원:  $P = (p_1, p_2, \dots, p_n), Q = (q_1, q_2, \dots, q_n)$ 

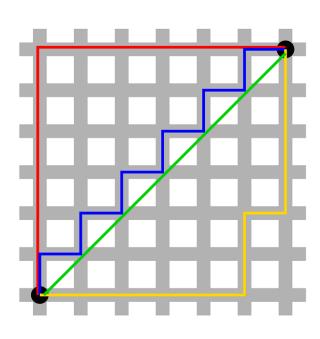
• 
$$d = \sqrt{\sum_{i=1}^{n} (p_i - q_i)^2}$$





- 맨하탄 거리: *Manhattan* Distance
  - 택시 거리라고도 함: 맨하탄에서 택시타고 가는 거리
  - 두 점 사이의 데카르트(Cartesian) 좌표계에서의 거리차의 절대값의 총합

$$- d = \sqrt{\sum_{i=1}^{n} |p_i - q_i|}$$



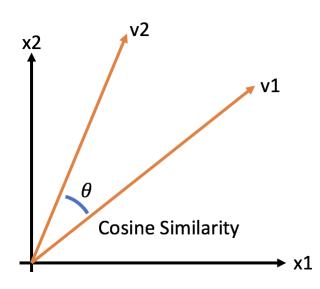


- 코사인 유사도: *Cosine* Similarity
  - 내적공간에서 두 벡터의 방향이 이루는 각의 코사인 값: [0, 1]

• 
$$S = \frac{P \cdot Q}{\|P\| \times \|Q\|} = \frac{\sum p_i \times q_i}{\sqrt{\sum p_i^2} \times \sqrt{\sum q_i^2}}$$

• 코사인 거리: Cosine Distance

$$-d = 1 - s$$



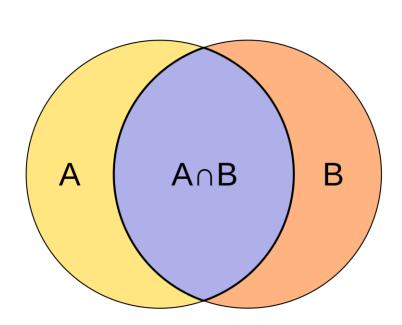


- 자카드 유사도: *Jaccard* Similarity
  - 두 집합 사이의 유사도를 측정하는 대표적인 방법
    - 페이스북에서 친구 추천할 때, 넷플릭스에서 영화 추천할 때.
  - 전체 집합의 크기와 교집합의 크기로 유사도 측정

$$- s(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

• 자카드 거리: Jaccard Distance

$$-d = 1 - s$$



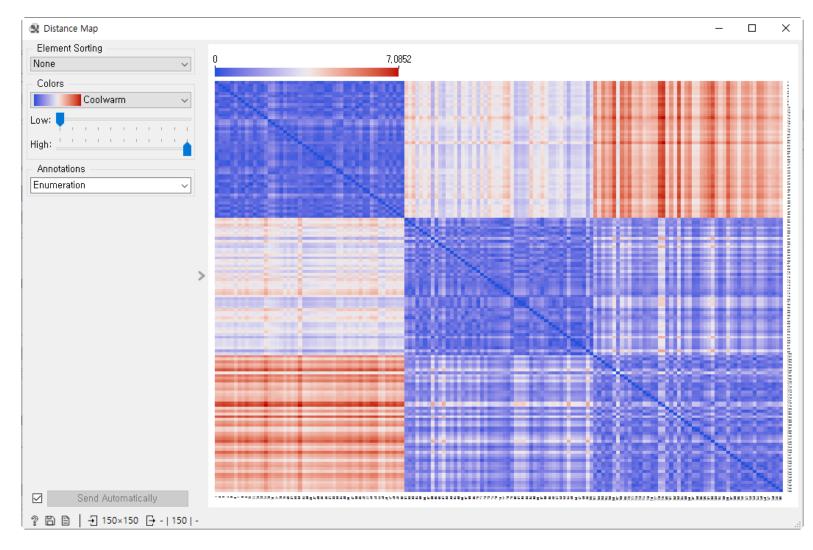


#### ■ 기타 유사도/거리 측정을 위한 지표: Distance Metrics

- 피어슨 거리: Pearson Distance
- 스피어먼 거리: Spearman Distance
- 해밍 거리: Hamming Distance
- 마할라노비스 거리: Mahalanobis Distance
- 바타차야 거리: Bhattacharyya Distance
- 체비셰프 거리: Chebyshev Distance

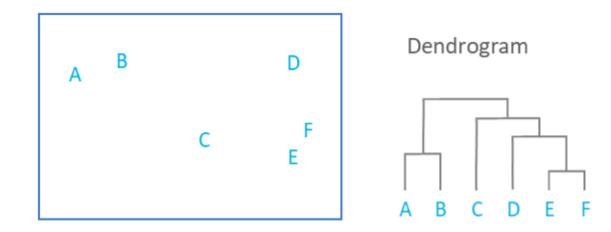


#### ■ IRIS 데이터셋에서의 유사도와 거리:





- 계층적 군집화: *Hierarchical* Clustering
  - 군집간의 거리를 이용하여 계층적으로 군집을 분석하는 방법
    - 병합적 방법: agglomerative, bottom-up approach
    - 분할적 방법: divisive(partitioning), top-down approach
  - 덴드로그램: Dendrogram
    - 군집의 계층적 구조를 그림으로 보여주는 방법



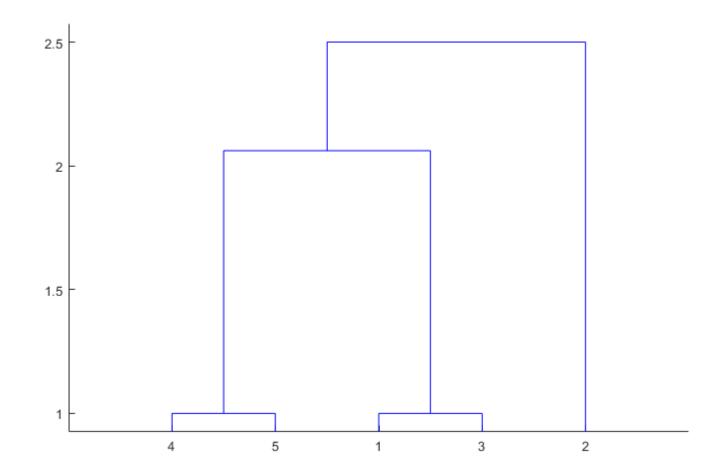


#### ○ O5. 군집 분석

- 군집 간의 거리를 측정하는 방법: *Linkage* Method
  - 최단 연결법: Single Linkage
    - 두 군집에 속하는 데이터 중에서 가장 가까운 데이터 간의 거리로 연결
  - 최장 연결법: Complete Linkage
    - 두 군집에 속하는 데이터 중에서 가장 먼 데이터 간의 거리로 연결
  - 평균 연결법: Average Linkage
    - 두 군집에 속하는 모든 데이터 간의 거리의 평균 거리로 연결
  - 중심 연결법: *Centroid* Linkage
    - 두 군집에서의 중심점(centroid)을 찾아서 두 중심점의 거리로 연결
  - Ward의 연결법: Ward Linkage
    - 두 군집을 합쳤을 때의 분산이 최소화 되는 군집을 합치는 방법

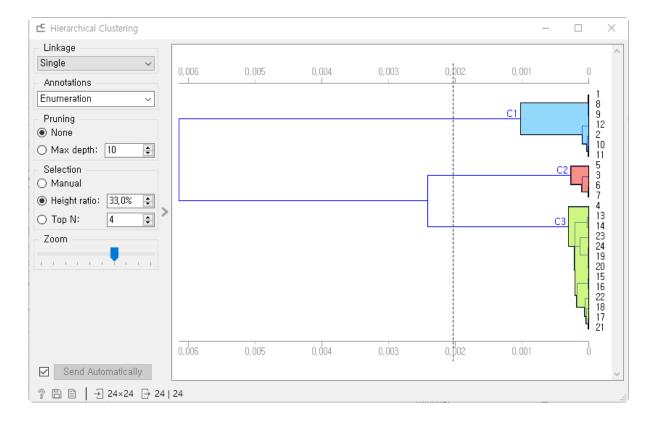


■ 덴드로그램으로 군집을 나누는 방법:

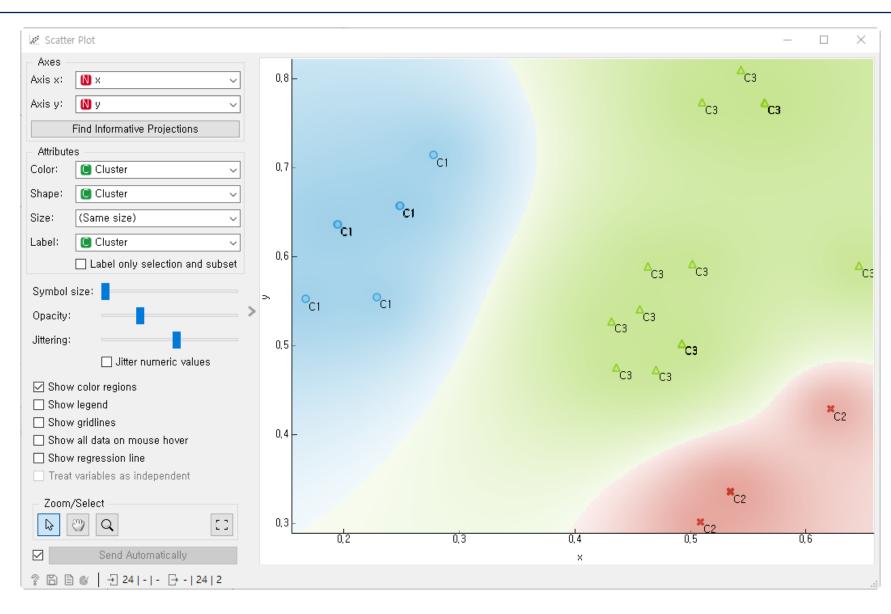


#### Orange: Hierarchical Clustering





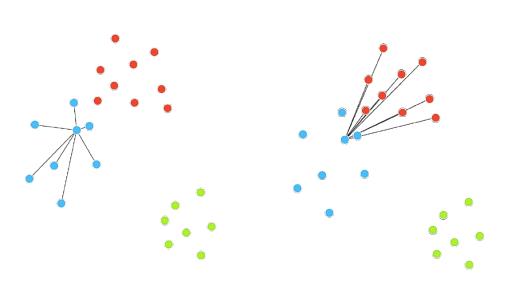
#### ⟨**>** 05. ·





- 군집 모델의 평가: Cluster Evaluation
  - 비지도 학습: 군집이 얼마나 서로 잘 구분이 되었는 지로 평가
  - 실루엣 점수: Silhouette Score
    - 같은 군집의 데이터와는 가깝고, 다른 군집 데이터와의 거리는 멀수록 좋다.
    - a(i): 같은 군집의 데이터와의 거리의 평균값
    - -b(i): 가장 가까운 다른 군집의 데이터와의 거리의 평균값

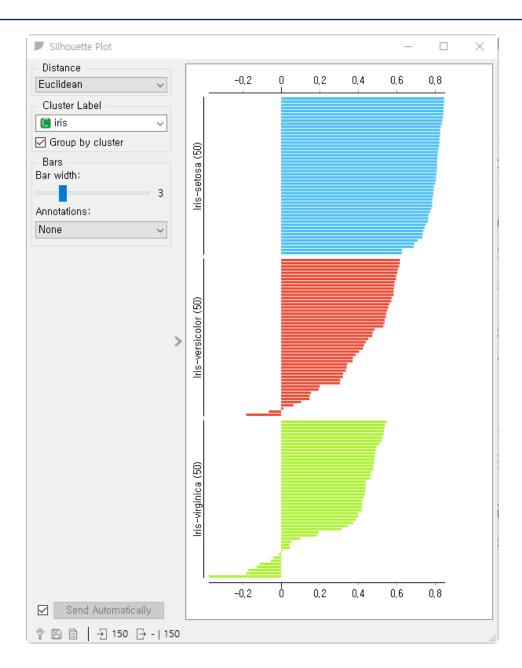
$$- s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$





Orange: Silhouette Plot



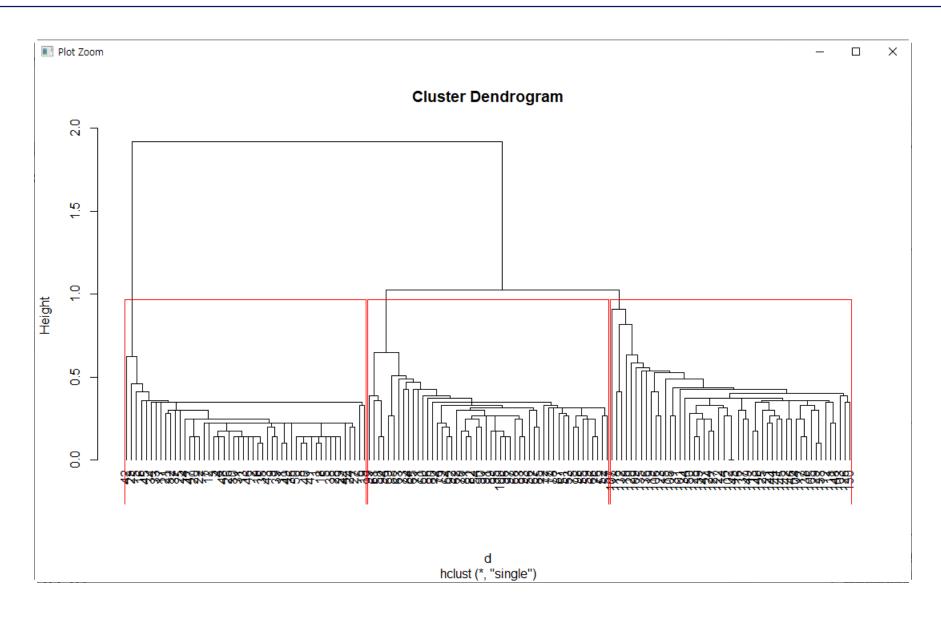




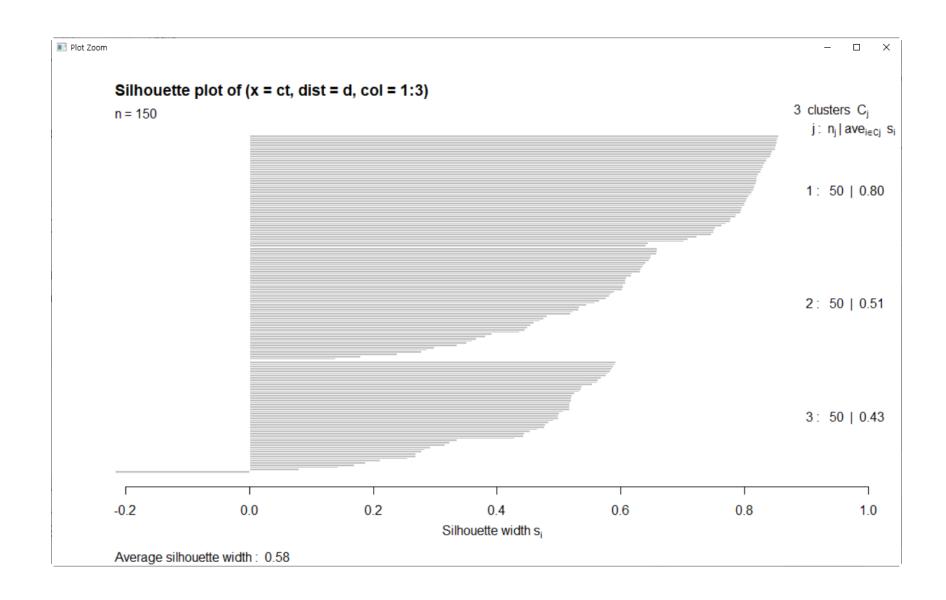
#### R: hclust()

```
#install.packages("cluster")
library(cluster)
df <- iris
df$Species <- as.numeric(df$Species)</pre>
head(df)
d <- dist(df, method="euclidean")</pre>
model <- hclust(d, method="single")</pre>
model
plot(model, hang=-1)
rect.hclust(model, k=3)
ct <- cutree(model, k=3)
plot(silhouette(ct, dist=d, col=1:3))
```



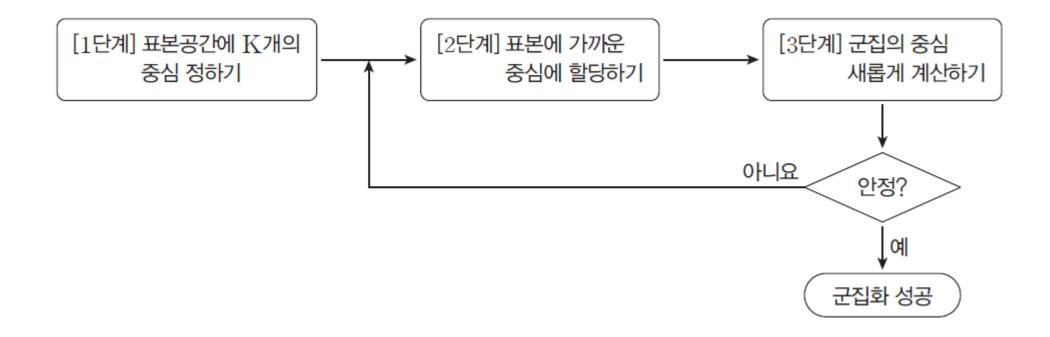


# ○ O5. 군집 분석





- k-평균 군집화: *k-means* clustering
  - 가장 일반적으로 널리 사용되는 군집화 알고리즘





#### ○ O5. 군집 분석

- k-평균 군집화: k-means clustering
  - 먼저 찾고자 하는 클러스터의 개수 k를 지정한다.
  - k개의 중심점(centroid) 위치를 임의로 지정한다.
  - 각 데이터들을 가장 가까운 중심점에 속하는 군집으로 결정한다.
  - 군집이 결정된 데이터들로 다시 중심점을 찾는다.
  - 다시, 각 데이터들을 가장 가까운 중심점에 속하는 군집으로 재편한다.
  - 중심점 이동이 없을 때까지 (혹은, 임계값 이하가 될 때까지) 반복한다.
  - 알고리즘의 동작이 멈추었을 때의 군집을 각 데이터의 군집으로 분류한다



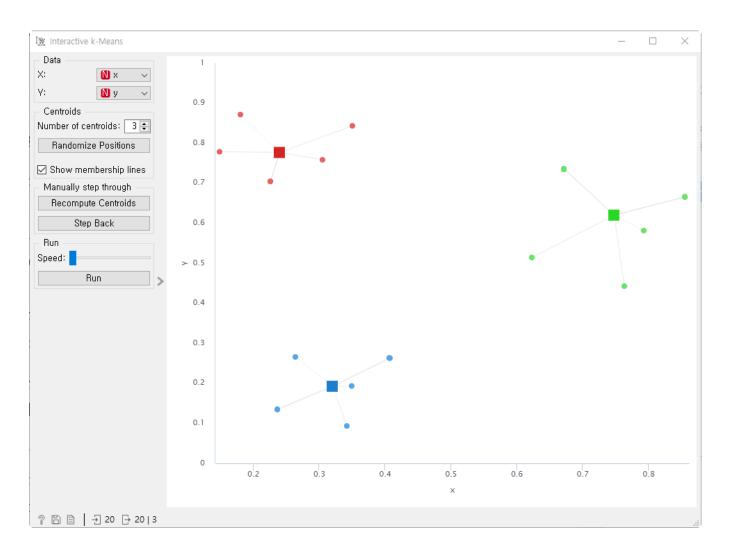
#### ■ k=2 (군집이 두 개)인 경우:





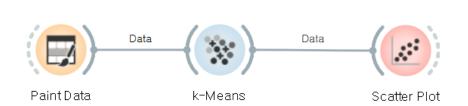
#### Orange: Interactive k-Means

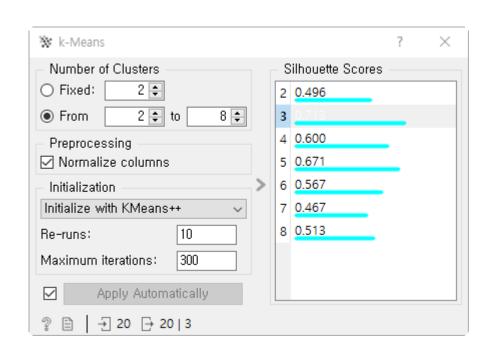






- 적절한 클러스터의 개수 k를 어떻게 선택할까?
  - 실루엣 점수가 가장 높아지는 k를 선택하기





# ○ O5. 군집 분석





- ZOO dataset: zoo.tab
  - UCI M/L Repository: <a href="https://archive.ics.uci.edu/ml/datasets/zoo">https://archive.ics.uci.edu/ml/datasets/zoo</a>

#### **Zoo Data Set**

Download: Data Folder, Data Set Description

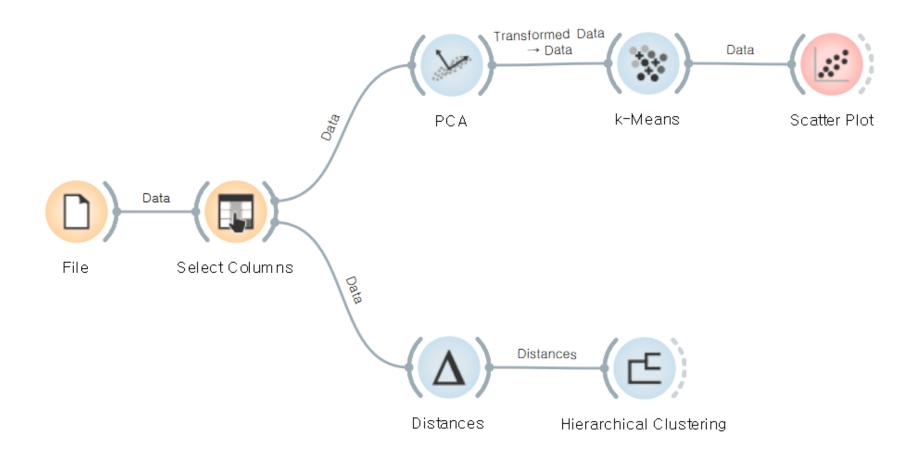
Abstract: Artificial, 7 classes of animals



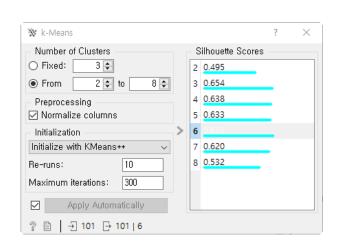
Data Set Characteristics:	Multivariate	Number of Instances:	101	Area:	Life
Attribute Characteristics:	Categorical, Integer	Number of Attributes:	17	Date Donated	1990-05-15
Associated Tasks:	Classification	Missing Values?	No	Number of Web Hits:	365821

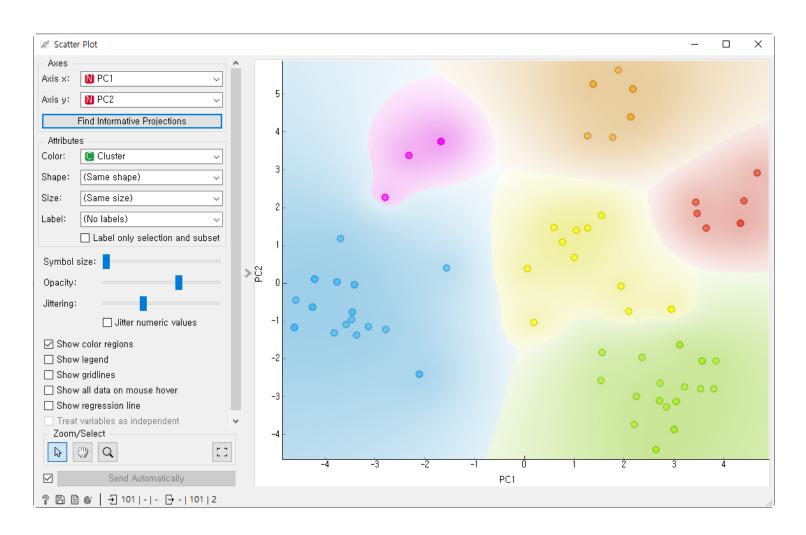


#### ■ ZOO dataset 클러스터링:

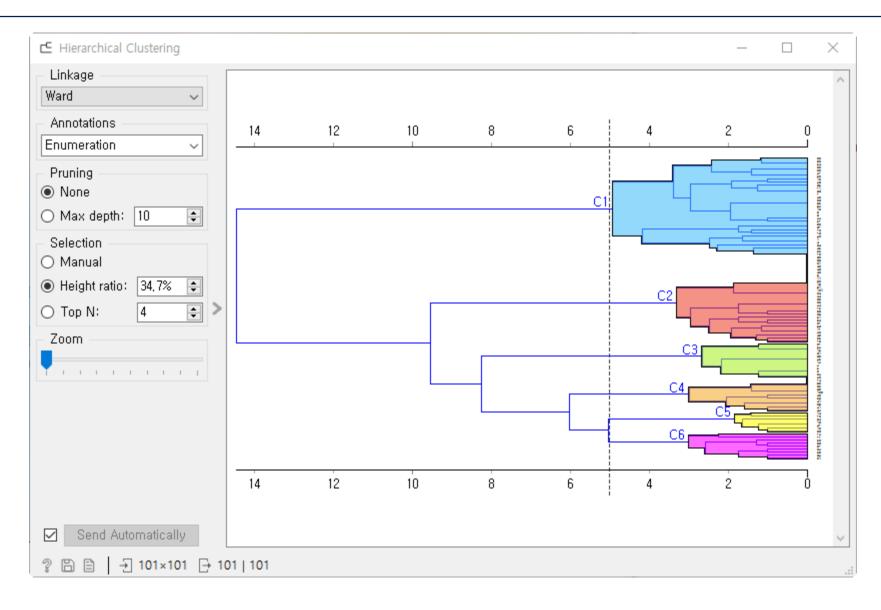














# ○5. 군집 분석

#### R: kmeans()

```
df <- iris[, 1:4]</pre>
model <- kmeans(df, center = 3)</pre>
model
model$centers
model$cluster
clusplot(df, model$cluster, lines=0, color=T, shade=T)
subset(df, model$cluster==1)
```



#### ▶ 05. 군집 분석

```
> model
K-means clustering with 3 clusters of sizes 62, 38, 50
Cluster means:
 Sepal.Length Sepal.Width Petal.Length Petal.Width
   5.901613 2.748387 4.393548 1.433871
   6.850000 3.073684 5.742105 2.071053
   5.006000 3.428000 1.462000 0.246000
Clustering vector:
 Within cluster sum of squares by cluster:
[1] 39.82097 23.87947 15.15100
(between_SS / total_SS = 88.4 %)
```



#### ○ O5. 군집 분석

- 데이터 정규화: Normalization
  - 데이터의 특성들이 비슷한 영향력을 행사할 수 있도록 변환해 줌
  - 단위의 차이, 값의 범위 차이 등으로 인한 차이 발생을 완화함
  - Min-Min 정규화: 값의 범위를 [0, 1] 구간에 있도록 변환

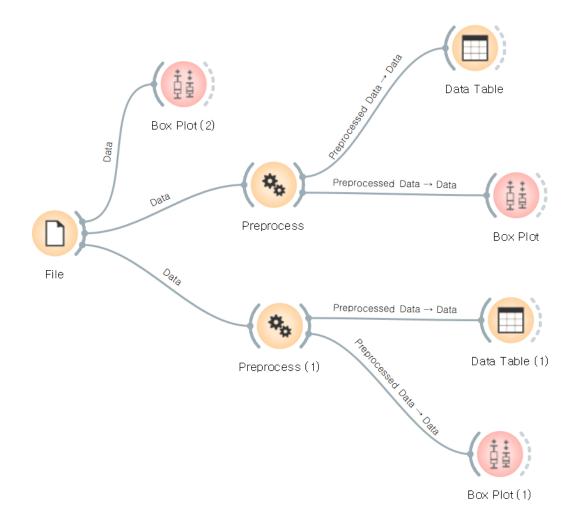
$$- X = \frac{X - MIN}{MAX - MIN}$$

• Z-Score 정규화: 값의 분포가 정규 분포를 따르도록 변환

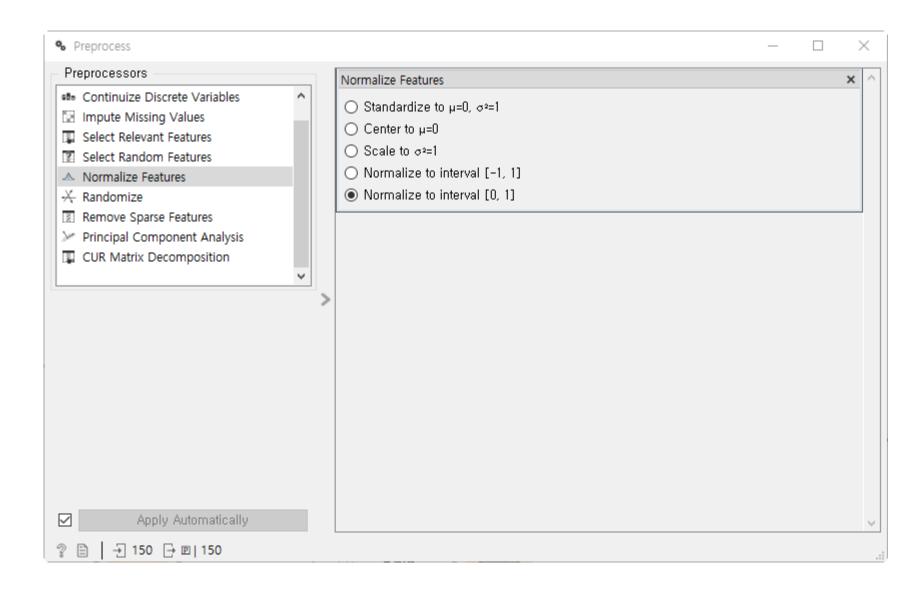
$$X = \frac{X-\mu}{\sigma}$$



#### Orange: Preprocess







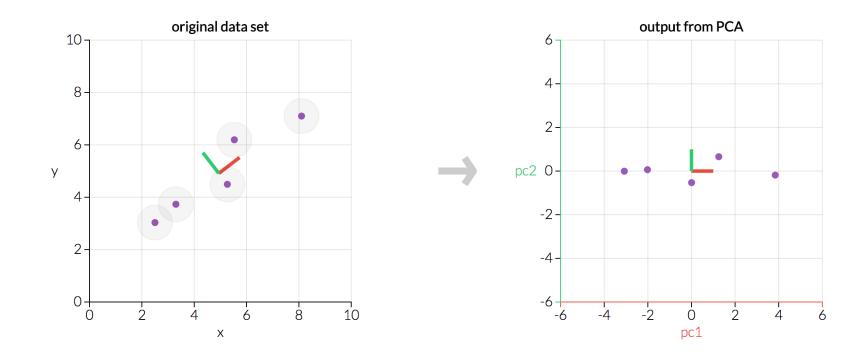


#### ○ O5. 군집 분석

- 차원 축소: Dimension Reduction
  - 고차원의 데이터를 2차원(또는 3차원) 데이터로 축소하는 기법
  - 차원 축소가 필요한 이유:
    - 시각화: 2, 3차원의 데이터는 산점도 등으로 시각화가 가능
    - 노이즈 제거: feature의 수가 줄어들기 때문에 노이즈도 같이 제거됨
    - 성능의 향상: 차원의 수가 작으면 처리 속도도 늘어남
  - 차원 축소의 단점:
    - 정보의 손실: 차원의 축소 과정에서 정보의 왜곡 손실 현상이 발생함

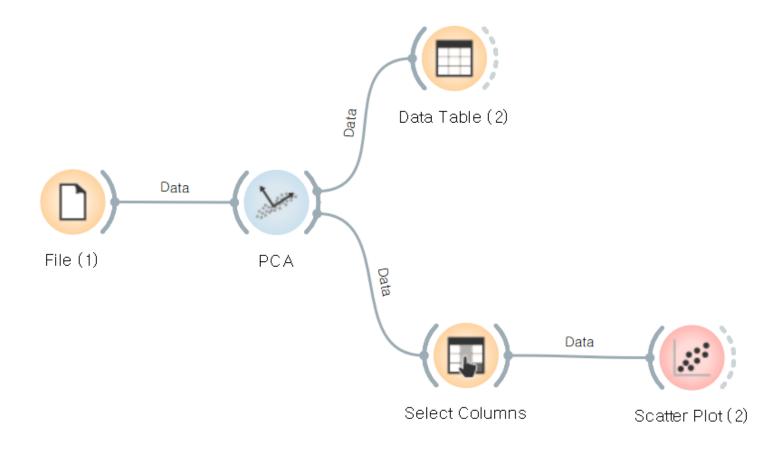


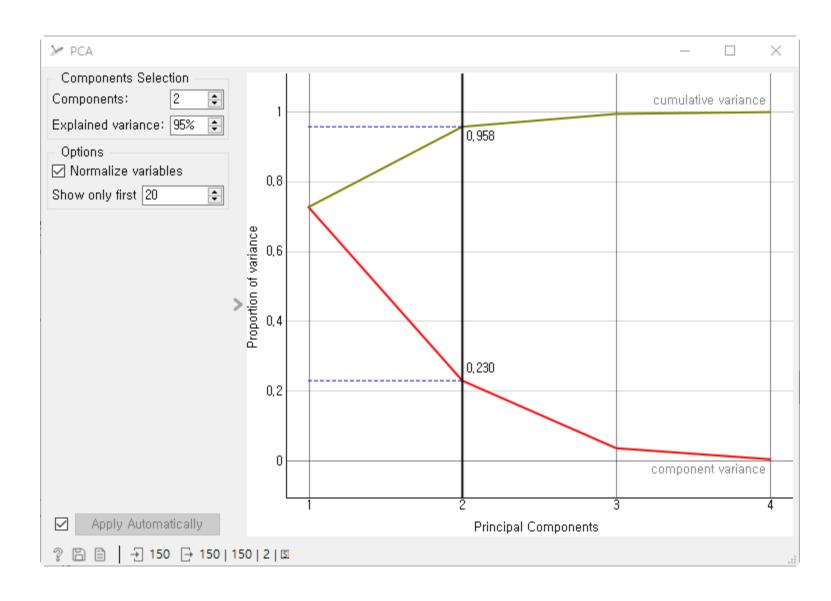
- 주성분 분석: PCA, Principal Component Analysis
  - 어떤 데이터 분포를 분산이 가장 큰 방향으로 정사영하여 주성분을 구함



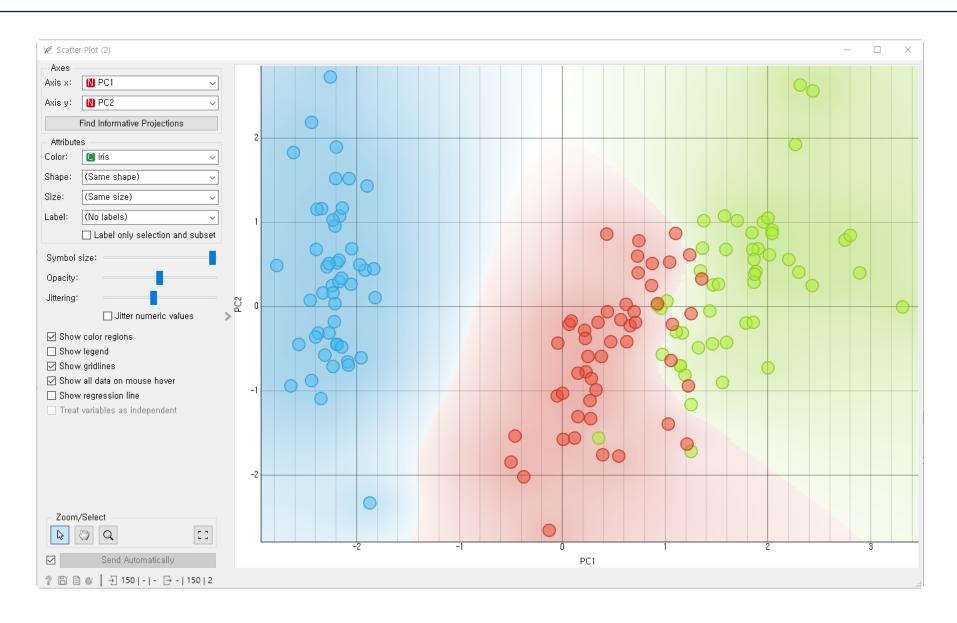


Orange: PCA









# Any Questions?

