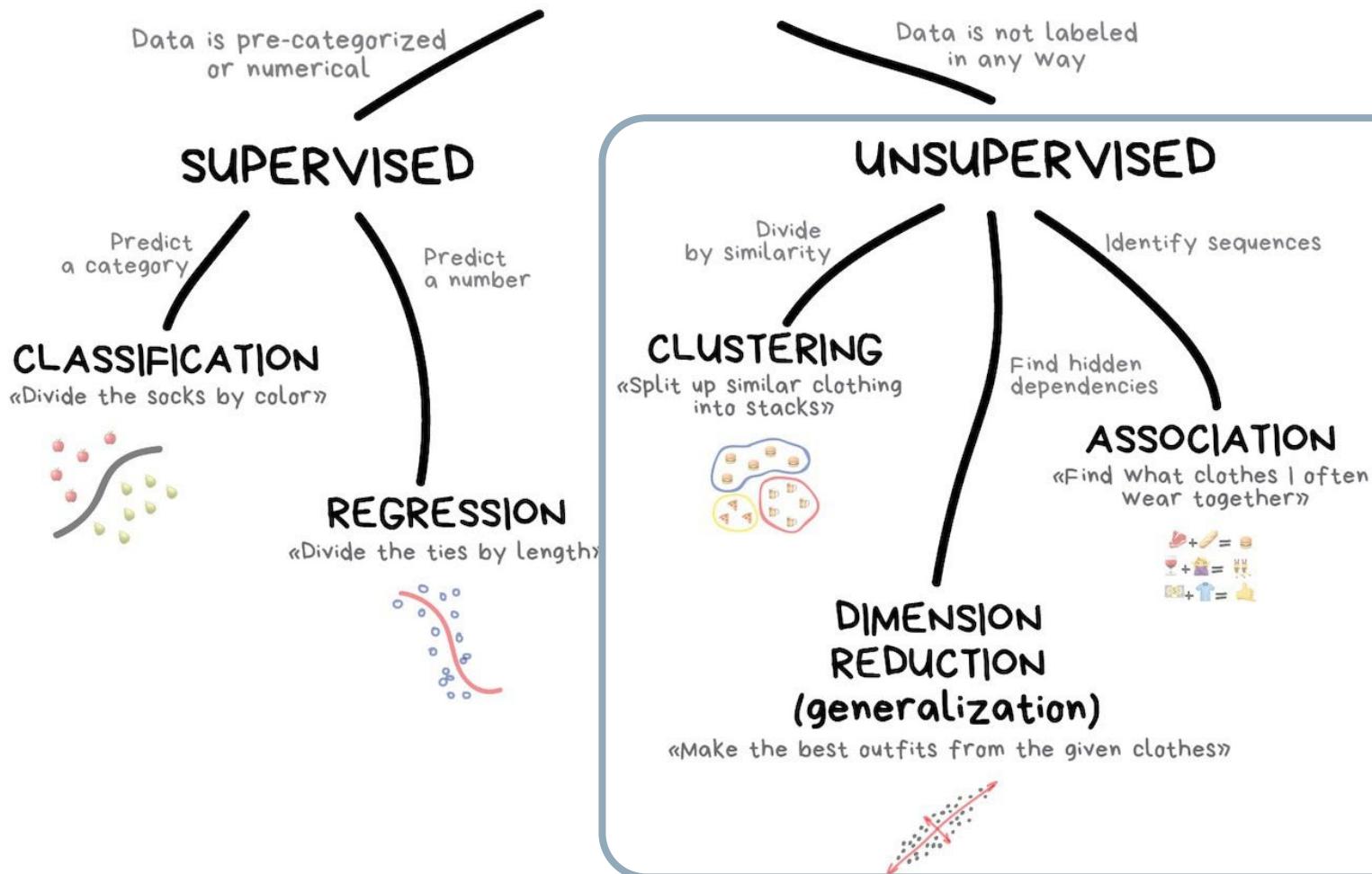


# CLUSTERING

Week 11

# Type of Learning

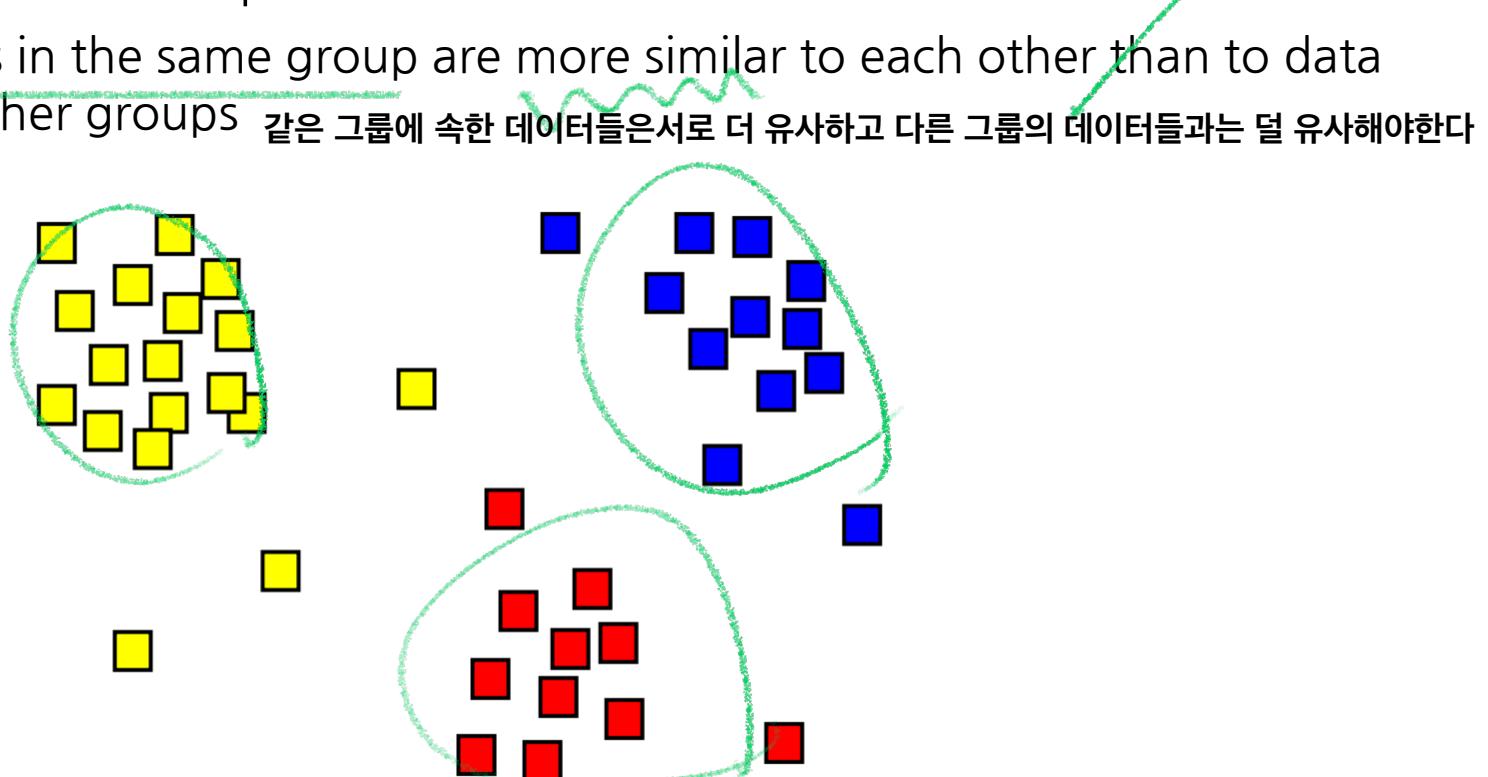
## CLASSICAL MACHINE LEARNING



# Clustering

# Unsupervised Learning: Clustering

- [Remind] Unsupervised learning is learning with unlabeled data
  - No certain output to be estimated
- **Clustering** is to group a set of data points to satisfy following conditions as much as possible
  - Data points in the same group are more similar to each other than to data points in other groups 같은 그룹에 속한 데이터들은 서로 더 유사하고 다른 그룹의 데이터들과는 덜 유사해야 한다



# Clustering

## □ Purpose

- ▣ Discover patterns in data
- ▣ Data compression and summarization 
- ▣ Anomaly detection
- ▣ Preprocessing for other machine learning tasks

1. 데이터 안에 숨겨진 패턴이나 구조 발견
2. 데이터를 압축하거나 요약하는데 사용
3. 이상치 탐지
4. 다른 머신러닝 작업을 위한 전처리 과정

# Clustering

- Types of clustering
  - ▣ Hard Clustering: Each data point belongs to only one cluster 각 데이터 포인트는 오직 하나의 클러스터에만 속함
  - ▣ Soft Clustering (Fuzzy Clustering): Each data point has a probability of belonging to multiple clusters 각 데이터 포인트는 여러 클러스터에 속할 확률을 가짐
  
- Types of clustering methods
  - ▣ Partition-based Clustering (e.g., K-Means): Partition-based clustering assigns data points into  $k$  clusters by minimizing intra-cluster variance. 데이터를  $k$ 개의 군집으로 나누되 군집 내부의 분산을 최소화하는 방식
  - ▣ Hierarchical Clustering: Builds a hierarchy of clusters using a tree-like structure (dendrogram) 트리구조(dendrogram)을 사용해 데이터 간 유사성에 따라 군집을 계층적으로 생성
  - ▣ Density-based Clustering (e.g., DBSCAN): Groups points that are closely packed together based on density 밀도가 높은 지점들을 군집으로 묶고 밀도가 낮은 부분은 노이즈로 간주
  - ▣ Model-based Clustering (e.g., Gaussian Mixture Models): Assumes that data is generated from a mixture of probability distributions 데이터가 확률분포의 혼합에서 생성되었다고 가정

# Clustering

- Data points in the same group are more similar to each other than to data points in other groups

**1. How to know  
some points are more similar than others?**



**Using distance measure**

**2. How to group?**



**Determine certain rule to group**

# $k$ -means Clustering

# *k*-means Clustering

- Objective function of clustering

$$\sum_i \min_j \|x_i - \mu_j\|^2$$

한 데이터가 가장 가까운 클러스터에 속한다고 가정  
데이터 점도  $\rightarrow$  유클리드 거리<sup>2</sup>  
 $\mu_j$ 는 j-th 클러스터의 중심점

## Combinatorial Optimization Problem

조합 최적화 문제:

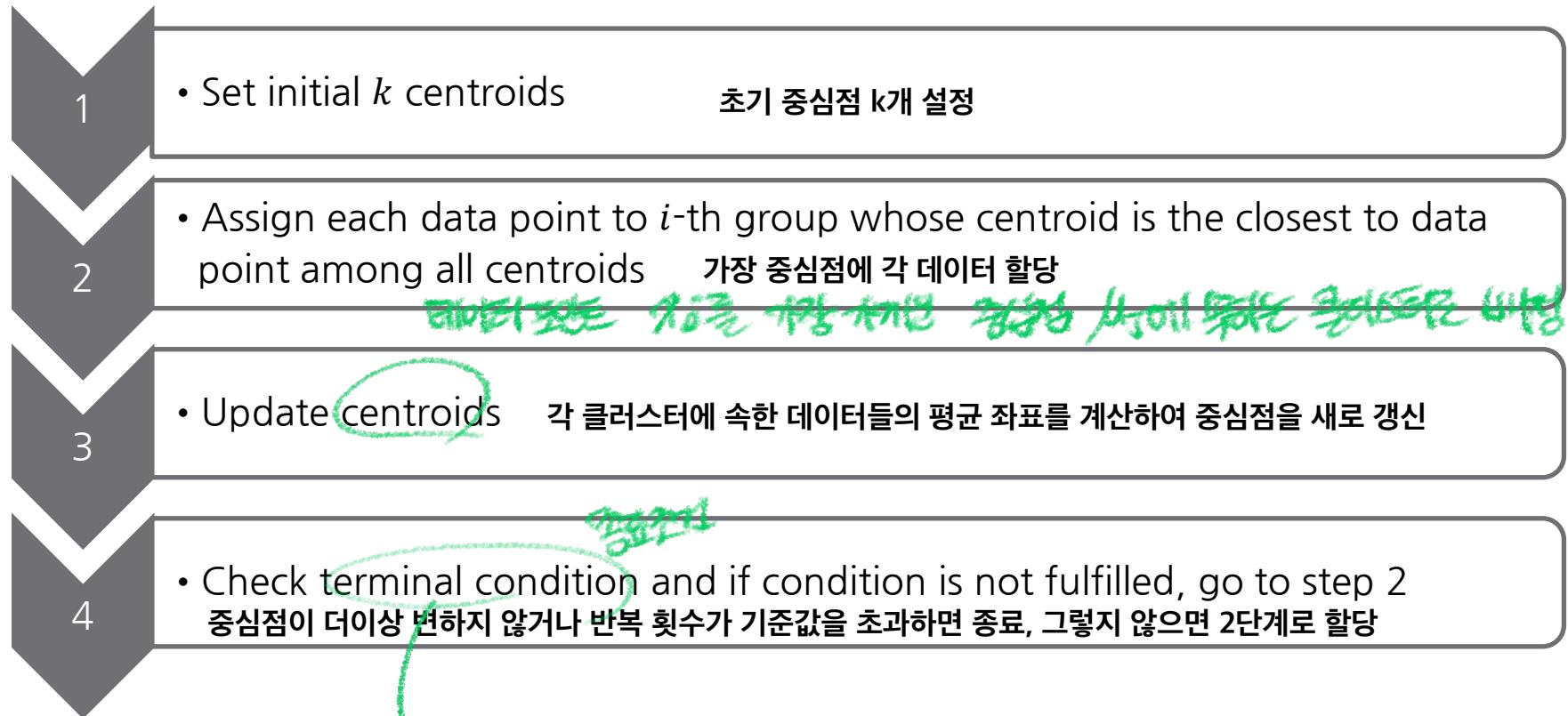
클러스터링은 단순히 수학적으로 미분 가능한 문제가 아니라 데이터를  $k$ 개의 그룹으로 나눈 조합을 선택해야하는 문제

이는 조합적인 성질을 가짐 정확한 해를 찾는 것은 NP-hard 문제일 수 있다.

따라서 K-means는 정확한 최적해를 구하기보다는 반복적으로 갱신하면서 근사해를 찾음

# *k*-means Clustering

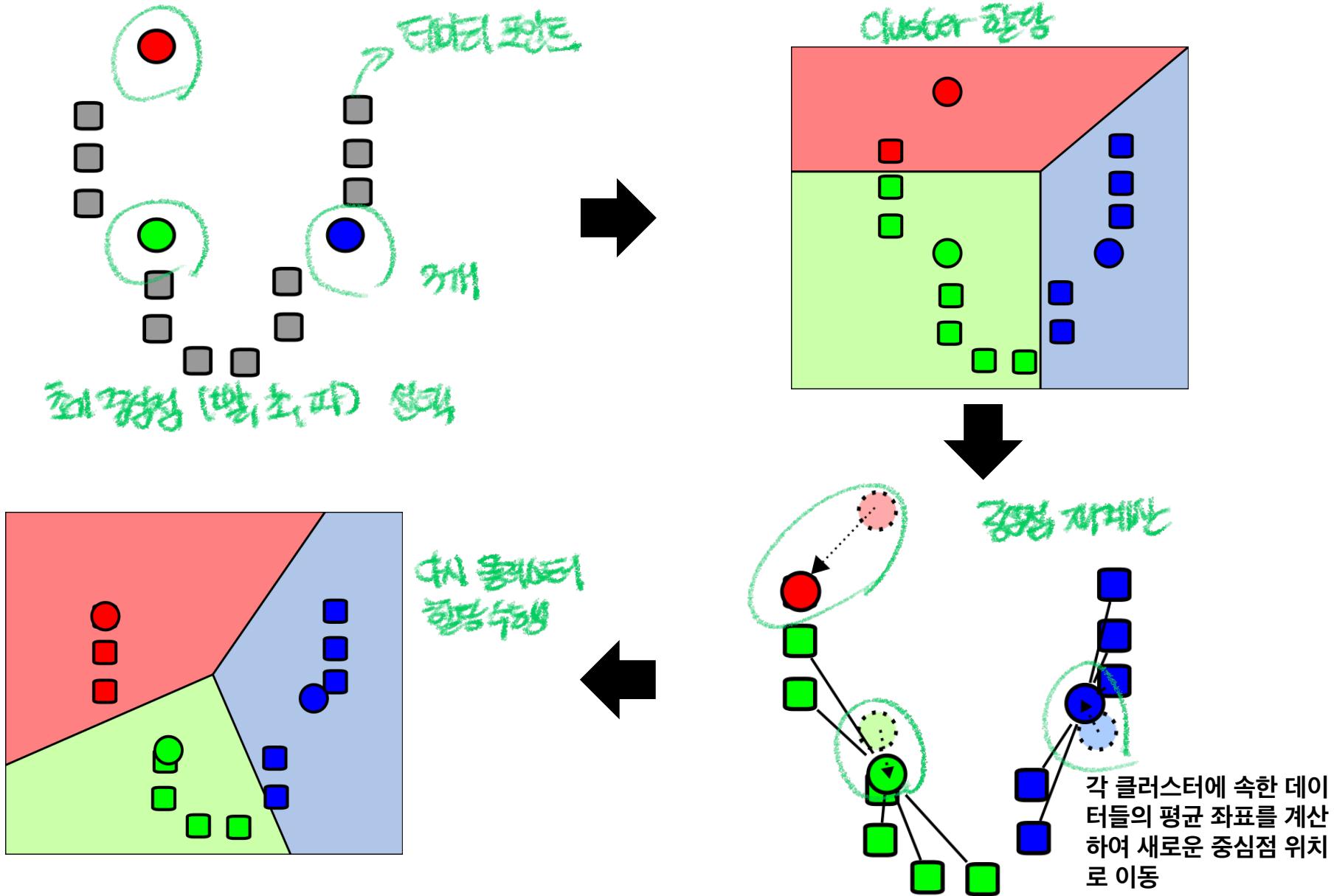
## □ Procedure of *k*-means clustering



### Terminal conditions

- ① [No change in centroids or the number of iteration is over the pre-specified threshold]

# *k*-means Clustering



# How to Update Centroids

중점 업데이트법

## □ Arithmetic mean

중점 업데이트  
방법

$$m_i^{(t+1)} = \frac{1}{|S_i|^{(t)}} \sum_{x_j \in S_i}^{(t)} x_j$$

중복되는 항목은 한번에 더해지지  
않아야 함

- $t$  is iteration
- $m_i$  is  $i$ -th group centroid
- $S_i$  is a set of  $i$ -th group and  $|S_i|$  is size of  $S_i$

## □ Example

- If  $(3,1)$ ,  $(2,2)$ ,  $(4,6)$  belong to group, updated centroid is

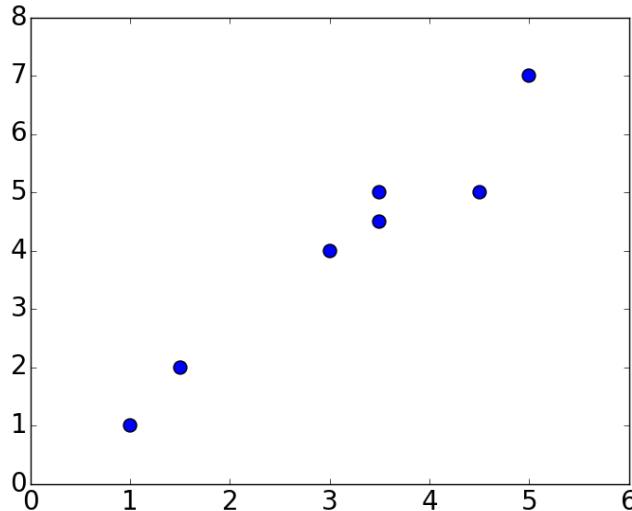
가장 좋은 업데이트 방식

$$\left( \frac{3 + 2 + 4}{3}, \frac{1 + 2 + 6}{3} \right) = (3, 2)$$

중점

# Question

- Clustering for 2D data set



A table with two rows and 8 columns. The first row contains labels 1 through 7. The second row contains values for x and y. The value 1.0 in the x row is circled in green. The value 5.0 in the x row and 7.0 in the y row are also circled in green.

	1	2	3	4	5	6	7
x	1.0	1.5	3.0	5.0	3.5	4.5	3.5
y	1.0	2.0	4.0	7.0	5.0	5.0	4.5

- 1) When  $k = 2$  and initial centroids are  $(1.0, 1.0)$  and  $(5.0, 7.0)$ , determine group of each data point
- 2) What are new centroids of two groups?

$\mathcal{B}_0 = (10, 10) \quad (15, 20) \quad (30, 40)$

$$\left( \frac{1+15+30}{3}, \frac{10+20+40}{3} \right) = (18.3, 23.3)$$

$\mathcal{B}_1 = (50, 7.0) \quad (35, 5.0) \quad (45, 5.0) \quad (3.5, 4.5)$

$$\left( \frac{50+7.0+3.5}{4}, \frac{7.0+5.0+5.0+4.5}{4} \right) = (14.3, 5.38)$$

# *k*-means Clustering: Pros and Cons

- Pros
  - Simple and efficient
  - Works well with large datasets
- Cons
  - Requires specifying  $k$
  - Sensitive to initial centroid selection

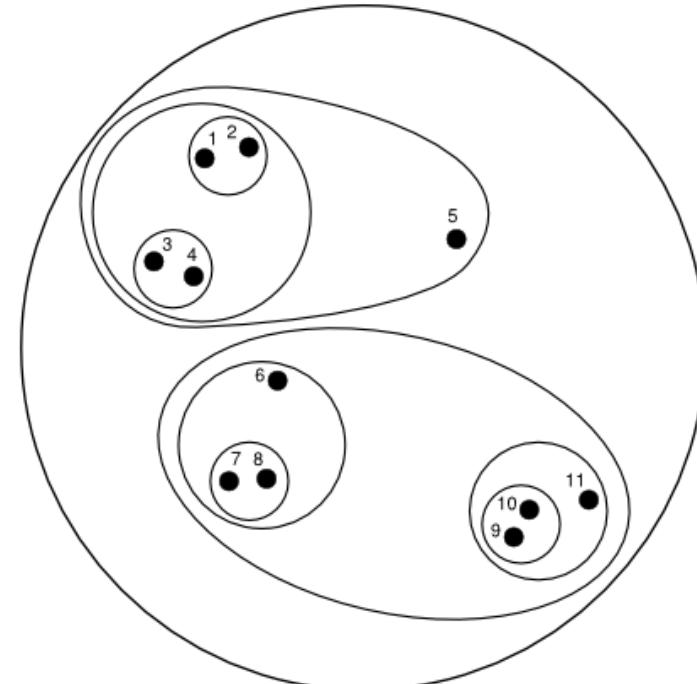
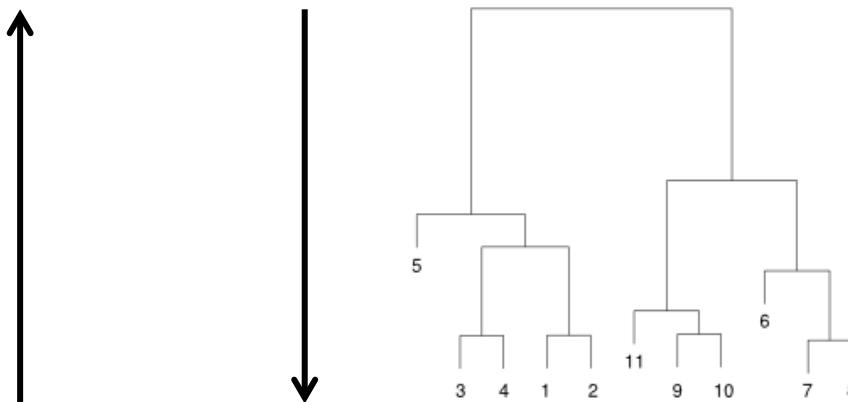
# Hierarchical Clustering

# Hierarchical Clustering

- Hierarchical clustering builds a hierarchy of clusters
  - ▣ Agglomerative: Bottom up approach, each data point starts in its own cluster and pairs of clusters are merged as one moves up the hierarchy
  - ▣ Divisive: Top down approach, all data points start in one cluster and splits are performed recursively as one moves down the hierarchy

모든 데이터가 하나의 클러스터로 시작 -> 조건에 따라 클러스터를 반복적으로 분할 -> 원하는 수 만큼 클러스터가 나올 때까지 재귀적으로 수행

Divisive



Agglomerative

각 데이터 포인트가 처음에는 자시난의 클러스터로 시작 ->  
가장 가까운 클러스터를 반복적으로 병합 -> 마지막에는 모  
든 포인트가 하나의 클러스터가 됨

# Linkage Criteria for Agglomerative Clustering

방법



bottom up

여기 주 응용을 어떻게 정의하나

- Way to calculate similarity between two clusters  $A, B$

Type	Formula
Complete-linkage	$\max\{d(a, b) : a \in A, b \in B\}$ 가장 먼 두점
Single-linkage	$\min\{d(a, b) : a \in A, b \in B\}$ 가장 가까운 두점
Mean linkage	$\frac{1}{ A  B } \sum_{a \in A} \sum_{b \in B} d(a, b)$ 두 클러스터에 속한 모든 짧은 거리의 평균
Centroid linkage	$d(c_A, c_B)$ 각 클러스터 중심간거리
Ward linkage	$\text{Var}(A \cup B) - \text{Var}(A) - \text{Var}(B)$ 전체에 따른 분산 증가량

- $a$  belongs to  $A$ ,  $b$  belongs to  $B$
- $\text{Var}(X)$  is within-cluster variance (variance of cluster  $X$ )

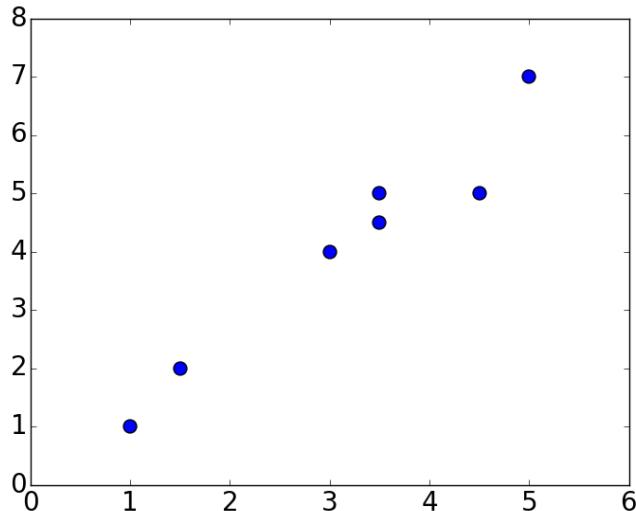
$$\text{Var}(X) = \frac{1}{n_A} \sum_{i \in A} \|\mathbf{x}_i - \mu_A\|^2$$

클러스터 내부분산

- $d(a, b)$  is distance between two data points  $a$  and  $b$

# Question

- Clustering for 2D data set



	1	2	3	4	5	6	7
x	1.0	1.5	3.0	5.0	3.5	4.5	3.5
y	1.0	2.0	4.0	7.0	5.0	5.0	4.5
c	1	1	2	2	2	2	2

1.4

- 1) Using complete-linkage, calculate linkage criterion of cluster 1 and 2

1.250328

$$d(a, b) = \sqrt{(1.0 - 3.0)^2 + (1.0 - 4.0)^2}$$

- 2) Using centroid-linkage, calculate linkage criterion of cluster 1 and 2

2.50328

$$M_1 = \left( \frac{1.0 + 1.5}{2}, \frac{1.0 + 2.0}{2} \right) = (1.25, 1.5)$$

$$M_2 = (3.9, 5.1)$$

8.47619

# Example: Single Linkage Clustering

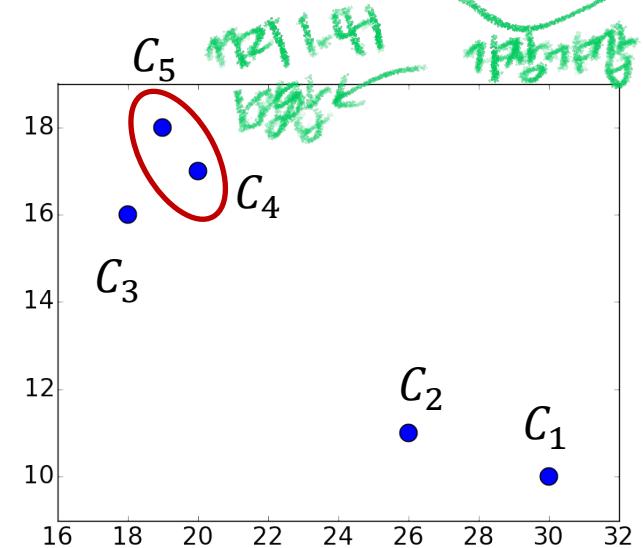
Agglomerative

- Find clusters through single linkage hierarchy clustering
  - Start each data as own cluster
  - Distance measure between two points: Euclidean distance

$$\text{Distance } d(A, B) = \min\{d(a, b) : a \in A, b \in B\}$$

	1	2	3	4	5
1	0				
2	4.12	0			
3	15.23	11.18	0		
4	12.21	8.48	4.12	0	
5	13.60	9.90	3.61	1.41	0

	1	2	3	4	5
x	30	26	16	20	19
y	10	11	16	17	18



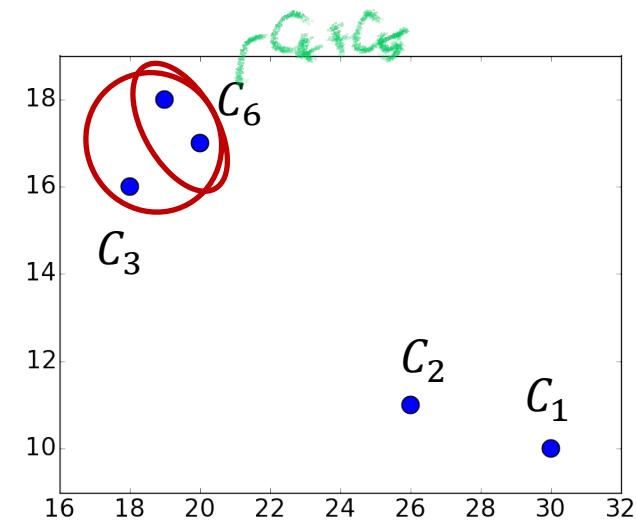
# Example: Single Linkage Clustering

- Find clusters through single linkage hierarchy clustering
  - ▣ Merge cluster 4 and 5 to create new cluster

Distance  $d(A, B) = \min\{d(a, b) : a \in A, b \in B\}$

	1	2	3	6
1	0			
2	4.12	0		
3	15.23	11.18	0	
6	12.21	8.48	3.61	0

	1	2	3	4	5
x	30	26	16	20	19
y	10	11	16	17	18



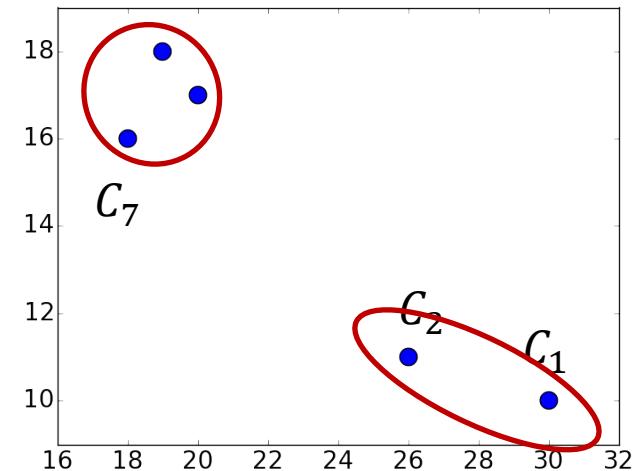
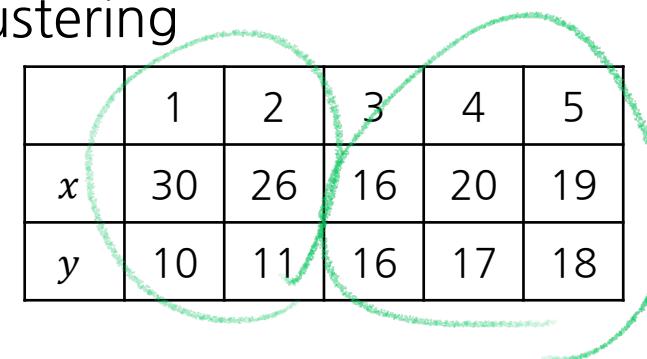
# Example: Single Linkage Clustering

- Find clusters through single linkage hierarchy clustering
  - ▣ Merge cluster 3 and 6 to create new cluster

Distance  $d(A, B) = \min\{d(a, b) : a \in A, b \in B\}$

	1	2	7
1	0		
2	4.12	0	
7	12.21	8.48	0

1.2 w/o 6



# Example: Single Linkage Clustering

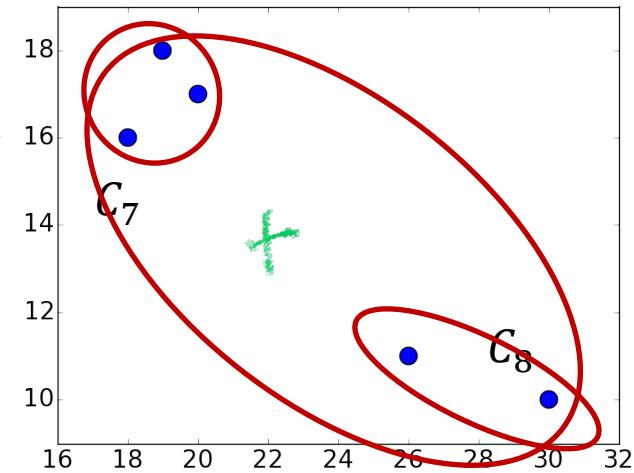
- Find clusters through single linkage hierarchy clustering
  - ▣ Merge cluster 1 and 2 to create new cluster

	1	2	3	4	5
$x$	30	26	16	20	19
$y$	10	11	16	17	18

Distance  $d(A, B) = \min\{d(a, b) : a \in A, b \in B\}$

	7	8
7	0	
8	8.48	0

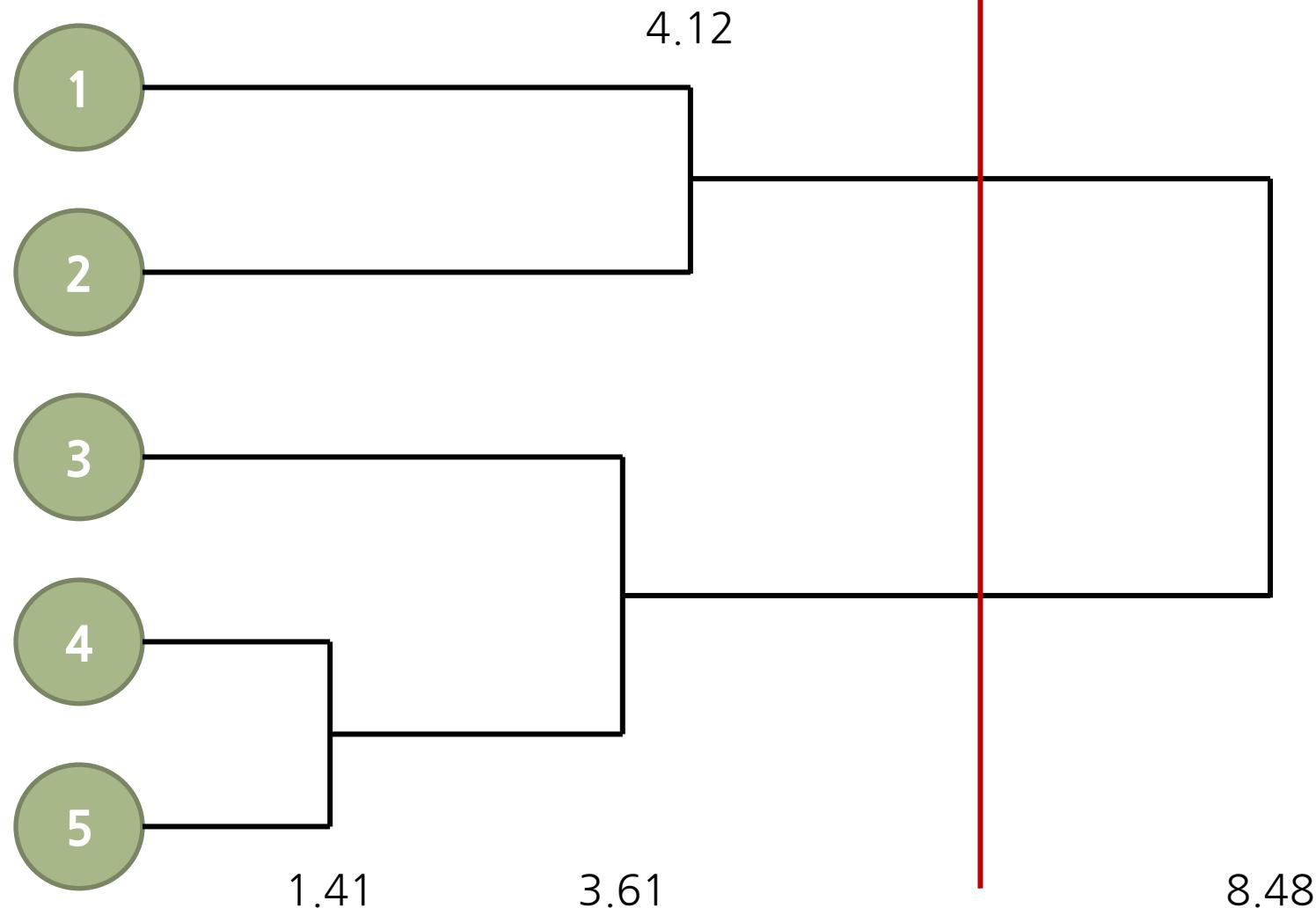
বেসিন ১৫  
১৮



৭

# Example: Single Linkage Clustering

- Dendrogram



# Example: Complete Linkage Clustering

Agglomerative

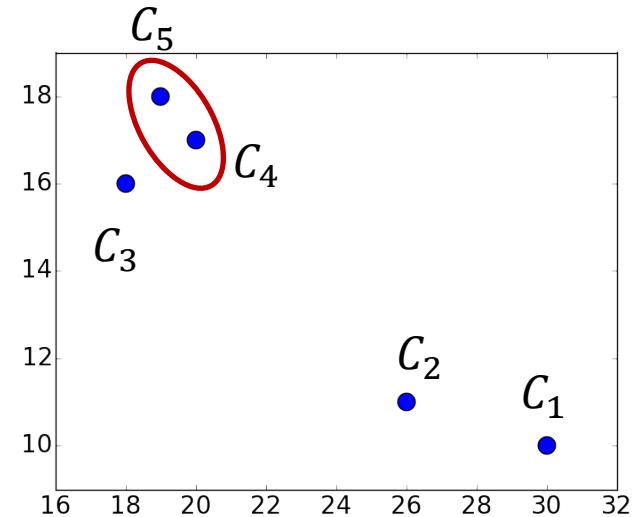
방법은 두 데이터가 다른 클러스터와 가장 먼저 떨어져 앉는 게 아름다

- Find clusters through complete linkage hierarchy clustering
  - Start each data as own cluster

	1	2	3	4	5
x	30	26	16	20	19
y	10	11	16	17	18

$$\text{Distance } d(A, B) = \max\{d(a, b) : a \in A, b \in B\}$$

	1	2	3	4	5
1	0				
2	4.12	0			
3	15.23	11.18	0		
4	12.21	8.48	4.12	0	
5	13.60	9.90	3.61	1.41	0



최초 1개지만 이제는 여러 singleton

# Example: Complete Linkage Clustering

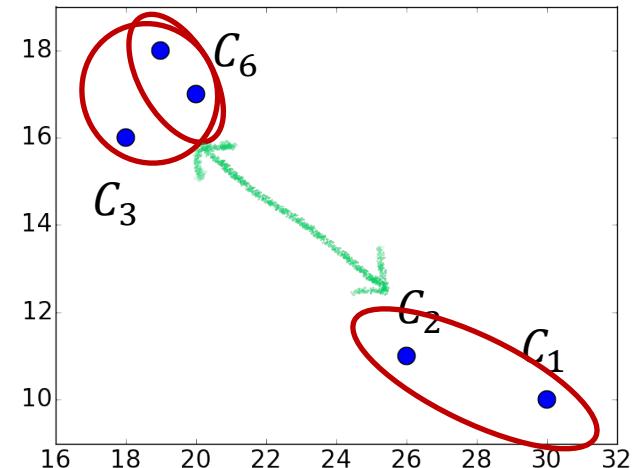
- Find clusters through complete linkage hierarchy clustering
  - Merge cluster 4 and 5 to create new cluster

	1	2	3	4	5
x	30	26	16	20	19
y	10	11	16	17	18

Distance  $d(A, B) = \max\{d(a, b) : a \in A, b \in B\}$

	1	2	3	6
1	0			
2	4.12	0		
3	15.23	11.18	0	
6	13.60	9.90	4.12	0

연쇄적 축소화가 가능하니 Max가 되도록 선택



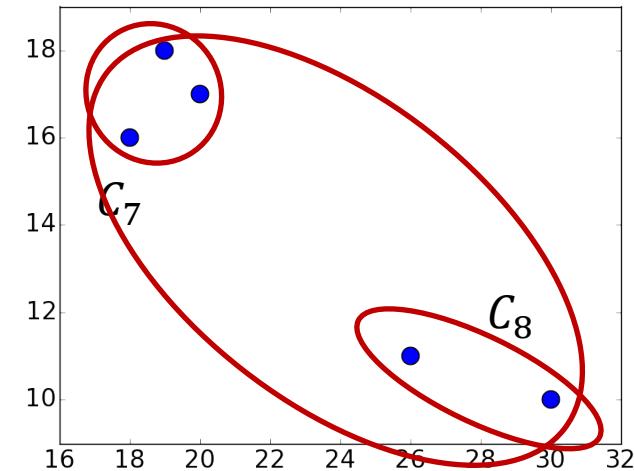
# Example: Complete Linkage Clustering

- Find clusters through complete linkage hierarchy clustering
  - ▣ Merge cluster 1 and 2 to create new cluster
  - ▣ Merge cluster 3 and 6 to create new cluster

	1	2	3	4	5
$x$	30	26	16	20	19
$y$	10	11	16	17	18

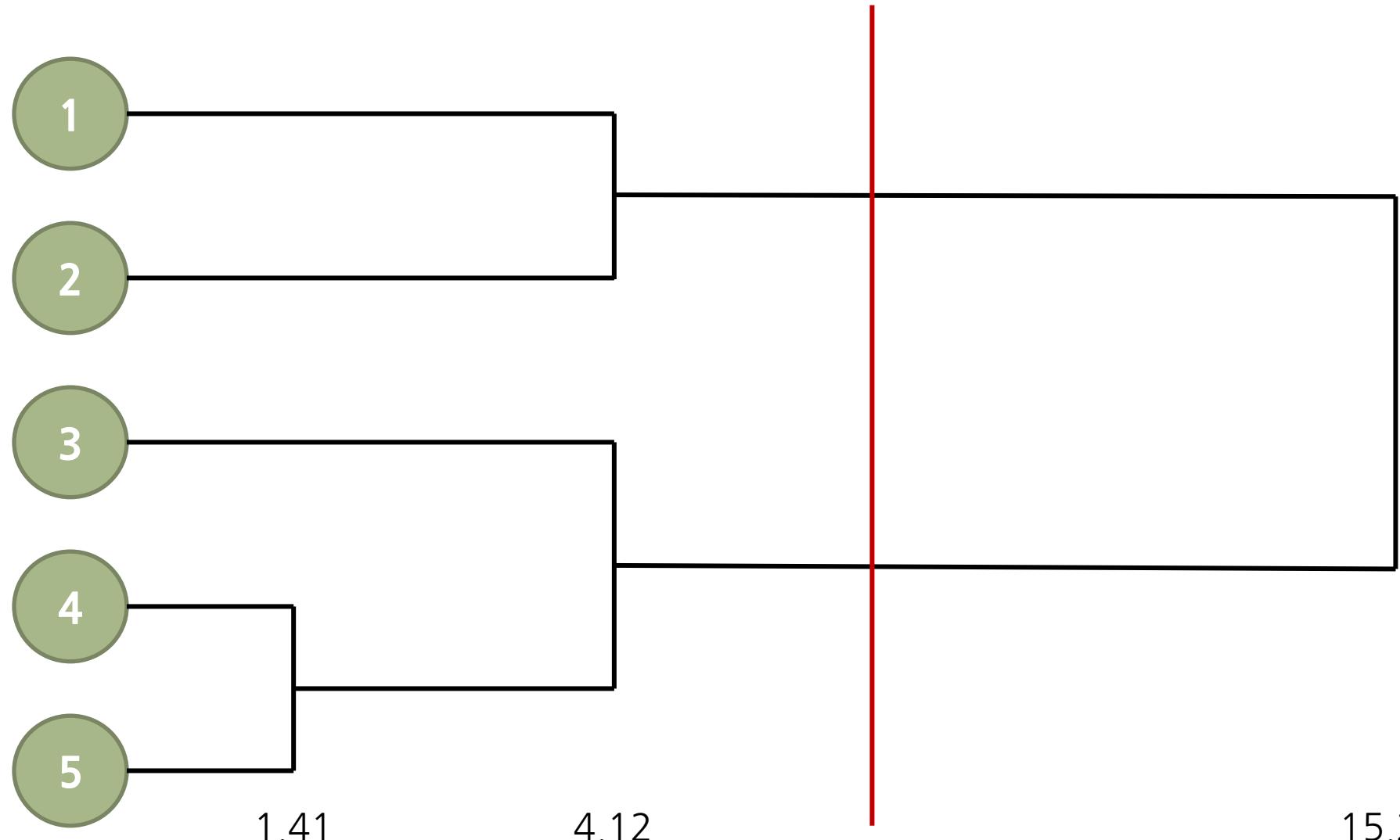
Distance  $d(A, B) = \max\{d(a, b) : a \in A, b \in B\}$

	7	8
7	0	
8	15.23	0



# Example: Complete Linkage Clustering

- Dendrogram



# Divisive Clustering - DIANA

분할정복과 방법

- Divisive method starts with one cluster including all samples
  - At each step, divide cluster into two sub clusters until every cluster consists of one data point 각 단계에서 클러스터를 두개의 하위 클러스터로 나누며, 최종 적으로는 각 클러스터가 하나만의 데이터만 가지도록 한다.
  - This algorithm is based on the average distance between one object and the others 이 알고리즘은 한 데이터 포인트와 다른 포인트들 간의 평균거리를 기준으로 함.

자기 자신과의 평균

$$\bar{d}(i, C) = \begin{cases} \frac{1}{|C|-1} \sum_{j \in C, j \neq i} d(i, j), & \text{if } i \in C \\ \frac{1}{|C|} \sum_{j \in C} d(i, j), & \text{if } i \notin C \end{cases}$$

- i represent i-th object

# DIANA Algorithm

1

- Consider all samples as one cluster

모든 샘플을 하나의 클러스터로 간주

2

- Select the cluster  $C$  containing two objects with the longest distance  
가장 멀리 떨어진 두개의 객체를 포함하는 클러스터  $C$  생성

가장 멀리 떨어진 두개의 객체를 포함하는 클러스터  $C$  생성

3

- Divide cluster  $C$  into two as follows (At first,  $C'$  is empty set( $\emptyset$ ))
  - Find object  $i$  with maximum  $\bar{d}(i, C)$   $C$ 내에서 다른 점들과의 평균 거리가 가장 먼 객체  $i$ 를 찾는다
  - $C \leftarrow C - \{i\}$ ,  $C' \leftarrow C' \cup \{i\}$   $C$ 를  $C'$ 와  $C$ 로 나눔
  - If there exist the objects  $j$  in  $C$  whose  $e(j) = \bar{d}(j, C) - \bar{d}(j, C') > 0$ , select one of them with maximum  $e(j)$ , remove  $j$  from  $C$  and add  $j$  into  $C'$   $C$ 에 남아있는 도는 객체의  $e(j)$ 가 음수이면 분할 종료, 더이상 옮길만한 데이터 없음
  - If  $e(j) < 0$  for all objects in  $C$ , finish this step

4

- Repeat step 2 and 3 until the number of clusters is the same as the number of samples  
클러스터 수가 데이터 수와 같아질 때까지 2,3단계 반복

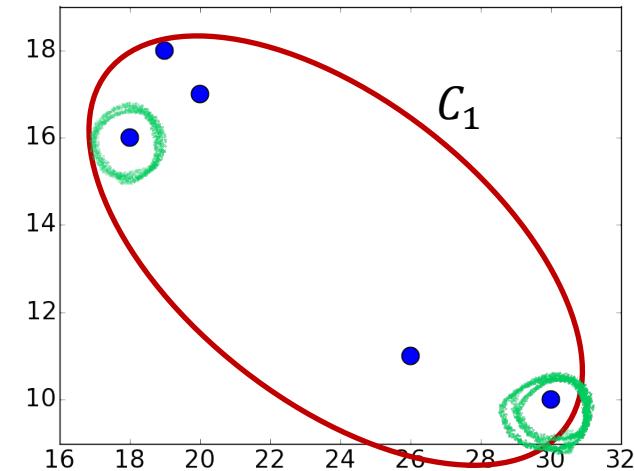
# Example: DIANA

- Find clusters through DIANA
  - Start with a cluster consisting of all objects
  - $C_1 = \{1,2,3,4,5\}$
  - $C_2 = \{\}$

Step 2: Find pair of objects wit the longest distance

$d(i,j)$	1	2	3	4	5
1	0				
2	4.12	0			
3	<b>15.23</b>	11.18	0		
4	12.21	8.48	4.12	0	
5	13.60	9.90	3.61	1.41	0

	1	2	3	4	5
x	30	26	16	20	19
y	10	11	16	17	18



# Example: DIANA

- Find clusters through DIANA
  - $C_1$  is the selected cluster

Step 3

$d(i, j)$	1	2	3	4	5
1	0	4.12	15.23	12.21	13.60
2	4.12	0	11.18	8.48	9.90
3	15.23	11.18	0	4.12	3.61
4	12.21	8.48	4.12	0	1.41
5	13.60	9.90	3.61	1.41	0

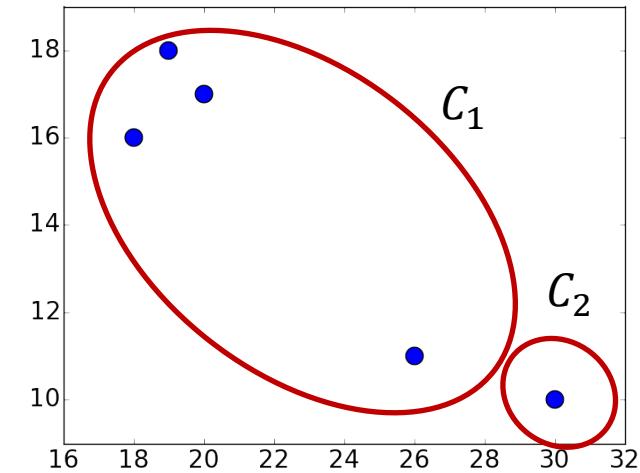
↓ Average except 0

	1	2	3	4	5
$\bar{d}(i, C_1)$	11.29	8.42	8.54	6.56	7.13

제일 큼

비율 시도는 흥미로운 이동

	1	2	3	4	5
$x$	30	26	16	20	19
$y$	10	11	16	17	18



계산방법

$$J(C_i, C_j) = \frac{\sum_{i \in C_i} \sum_{j \in C_j} d(i, j)}{\left(\sum_{i \in C_i} \sum_{j \in C_j} d(i, j)\right)^2}$$

# Example: DIANA

- Find clusters through DIANA

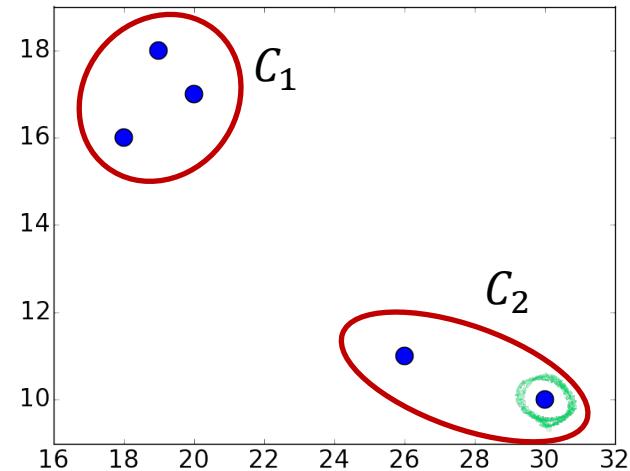
- $C_1 = \{2, 3, 4, 5\}, C_2 = \{1\}$

Step 3

	2	3	4	5
$\bar{d}(i, C_1)$	9.85	6.30	4.67	4.97
$\bar{d}(i, C_2)$	4.12	15.2	12.2	13.6
$e(i)$	5.73	-8.9	-7.53	-8.63

한국어로 쓰여진 설명입니다.  
C2에 속하는 데이터는 1입니다.

	1	2	3	4	5
x	30	26	16	20	19
y	10	11	16	17	18



# Example: DIANA

- Find clusters through DIANA

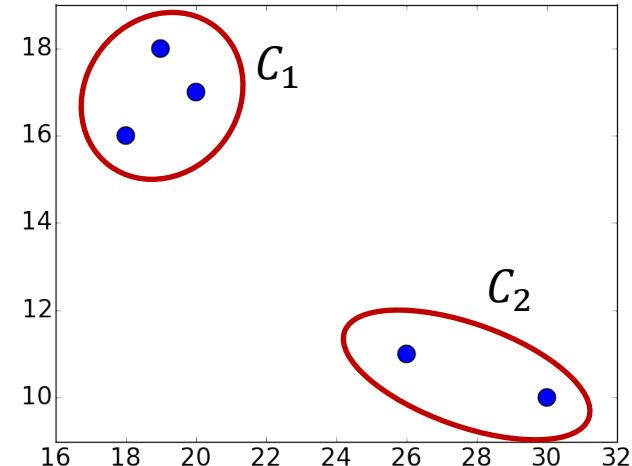
- $C_1 = \{3, 4, 5\}, C_2 = \{1, 2\}$

Step 3

	3	4	5
$\bar{d}(i, C_1)$	3.87	2.77	2.51
$\bar{d}(i, C_2)$	13.21	10.35	11.75
$e(i)$	-9.34	-7.58	-9.24

특정 초기 이동에 유리한  
기준이 없음  
→ 초기화

	1	2	3	4	5
$x$	30	26	16	20	19
$y$	10	11	16	17	18



# Example: DIANA

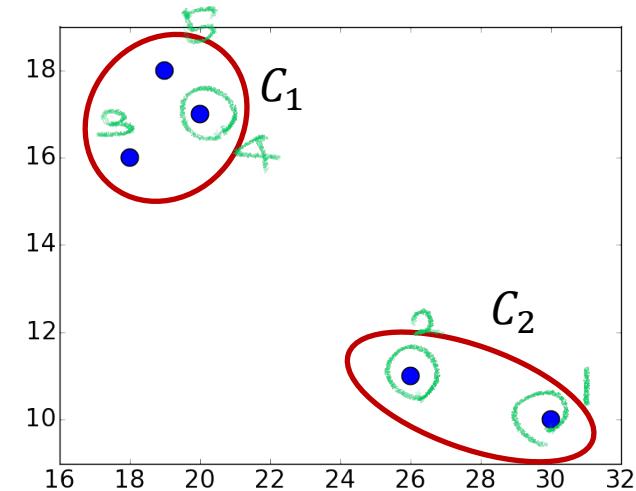
- Find clusters through DIANA
  - ▣  $C_1 = \{3,4,5\}, C_2 = \{1,2\}$
  - ▣ Find pair of objects wit the longest distance

방법: 가장 멀리 떨어져

Step 2: Find pair of objects wit the longest distance

$d(i,j)$	1	2	3	4	5
1	0				
2		0			
3			0		
4				0	
5			3.61	1.41	0

	1	2	3	4	5
$x$	30	26	16	20	19
$y$	10	11	16	17	18



# Example: DIANA

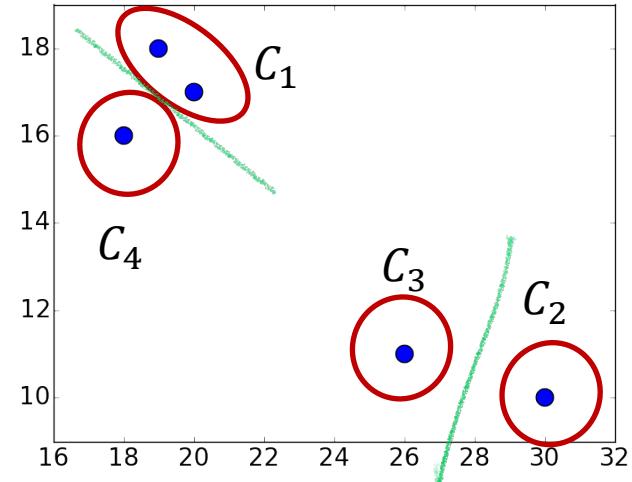
- Find clusters through DIANA
  - ▣  $C_1 = \{4, 5\}, C_4 = \{3\}$

Step 3

	4	5
$\bar{d}(i, C_1)$	1.41	1.41
$\bar{d}(i, C_4)$	4.12	3.61
$e(i)$	-2.71	-2.20

图2，324是

	1	2	3	4	5
$x$	30	26	16	20	19
$y$	10	11	16	17	18



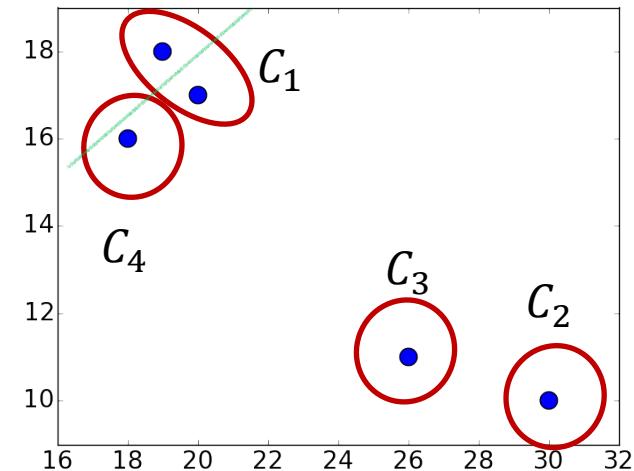
# Example: DIANA

- Find clusters through DIANA
  - ▣ Select  $C_2$
  - ▣  $C_1 = \{1,2\}, C_3 = \{\}$
  - ▣  $C_2$  contains only two object, so divide  $C_2$  into two clusters directly:  $C_2 = \{1\}, C_3 = \{2\}$
  
- ▣ Select  $C_1$
- ▣  $C_1 = \{3,4,5\}, C_4 = \{\}$

Step 3

	3	4	5
$\bar{d}(i, C_1)$	<b>3.87</b>	2.77	2.51

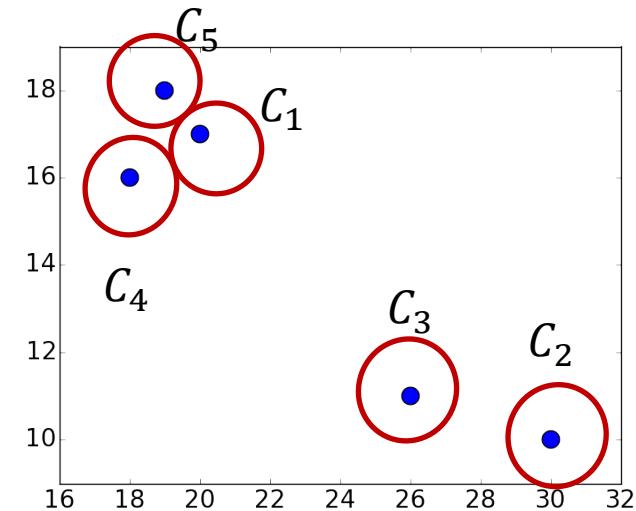
	1	2	3	4	5
$x$	30	26	16	20	19
$y$	10	11	16	17	18



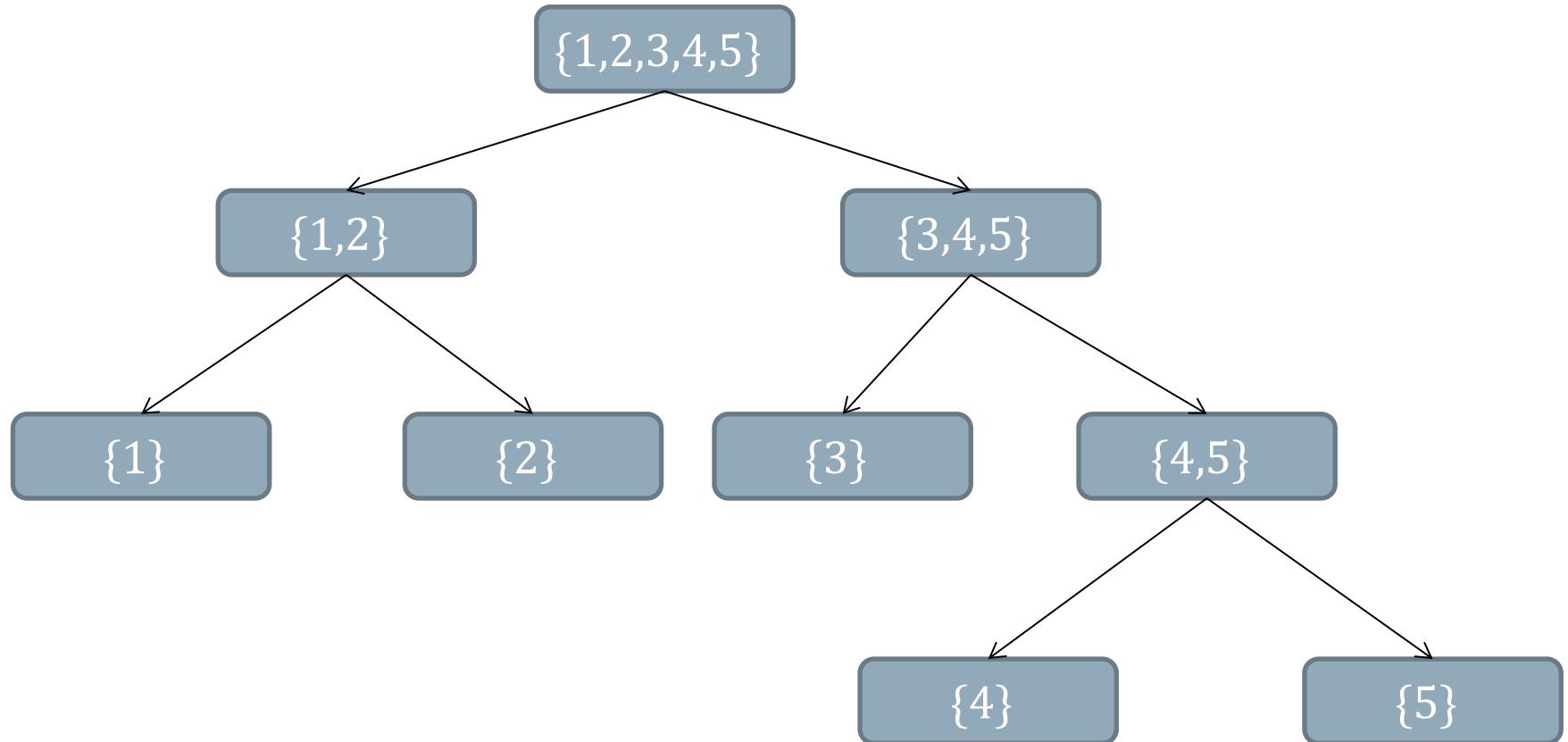
# Example: DIANA

- Find clusters through DIANA
  - ▣ Select  $C_1$
  - ▣  $C_1 = \{4,5\}, C_5 = \{ \}$
  - ▣  $C_1$  contains only two object, so divide  $C_1$  into two clusters directly:  $C_1 = \{4\}, C_5 = \{5\}$

	1	2	3	4	5
$x$	30	26	16	20	19
$y$	10	11	16	17	18

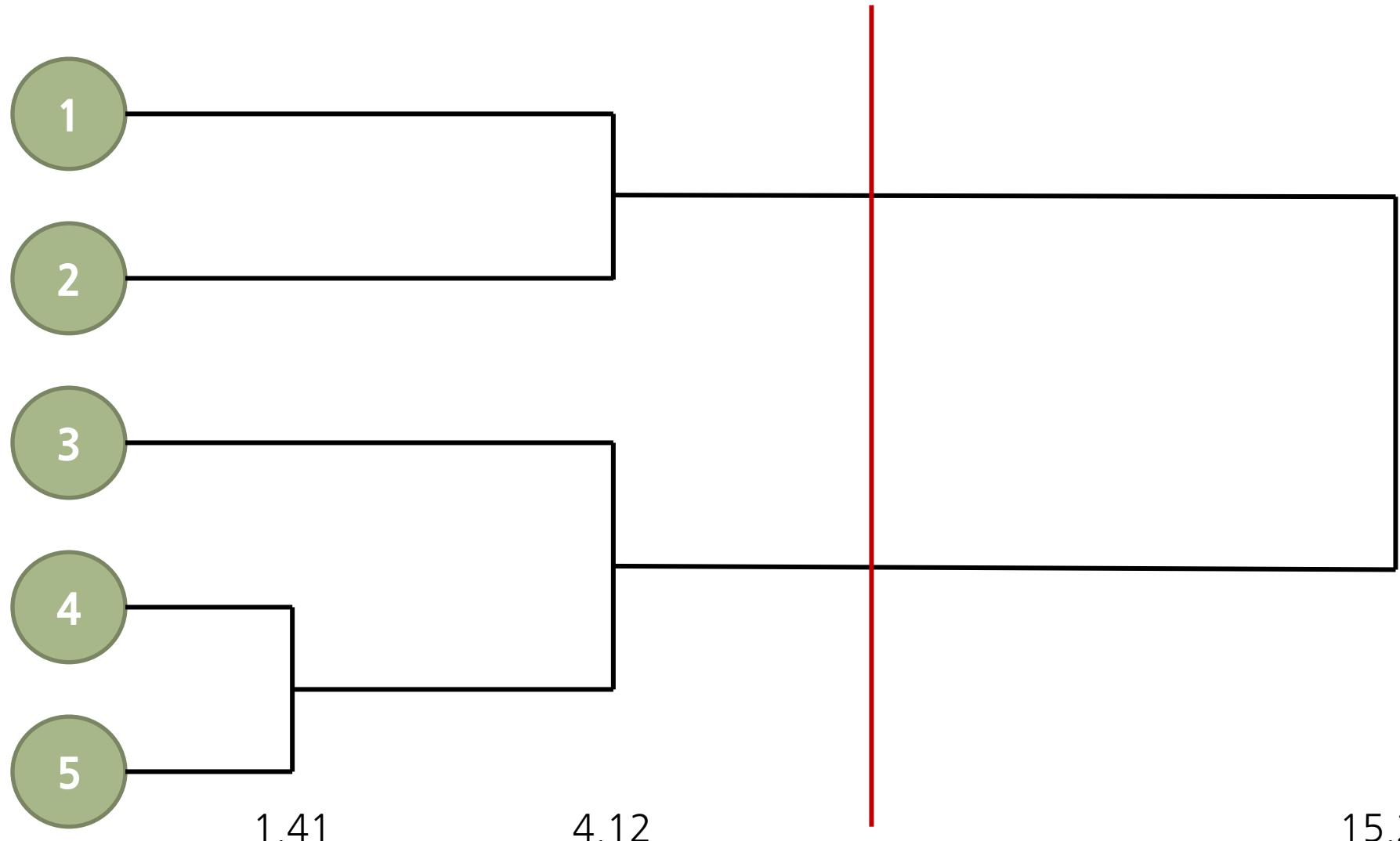


# Example: DIANA



# Example: DIANA

- Dendrogram



# Hierarchical Clustering: Pros and Cons

- Pros
  - ▣ No need to specify  $k$  in advance 사전에 k값을 지정할 필요 X
  - ▣ Produces a hierarchical structure 계층 구조 생성
- Cons
  - ▣ Computationally expensive for large datasets 대량 데이터셋에서 계산비용 증가
  - ▣ Sensitive to noise 소음에 민감

# Evaluation Measures

# How to Measure Clustering Quality

- Clustering problem is unsupervised problem
  - No explicit answer for learning
  - We need to define a method to measure quality of clustering

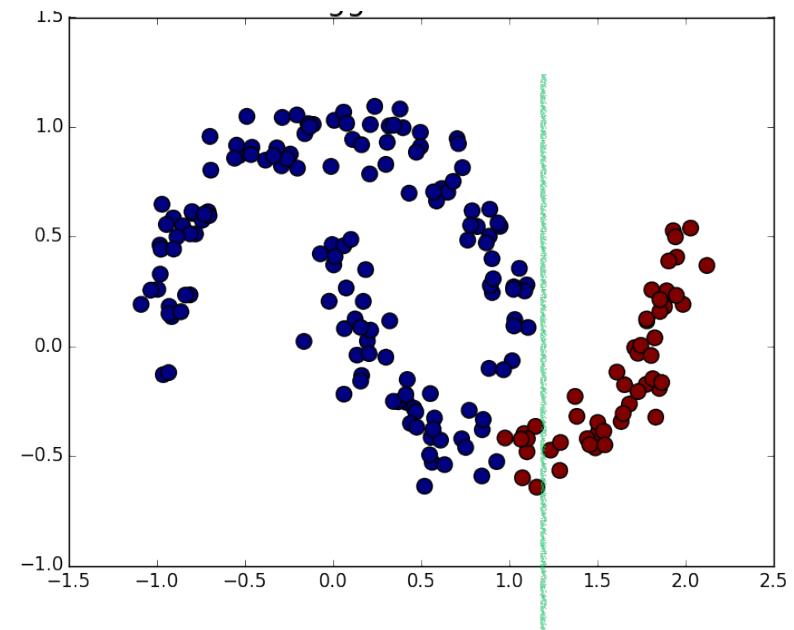
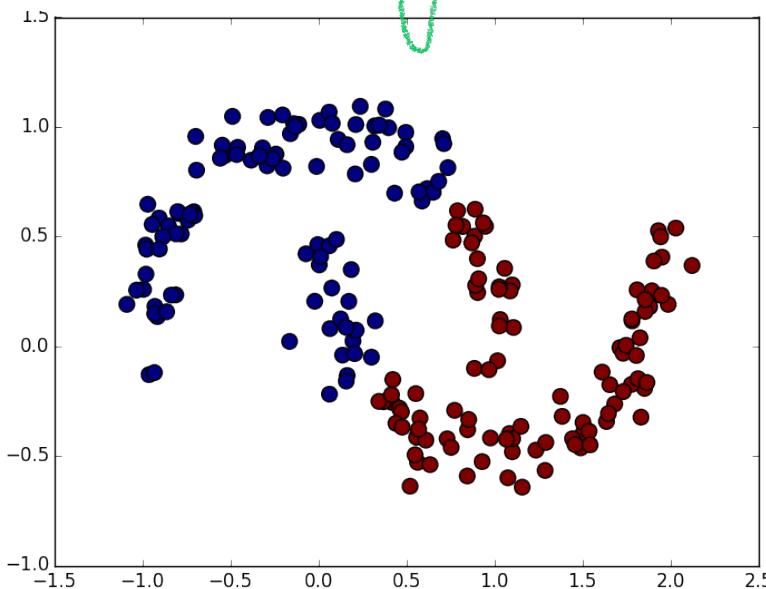
비지도 학습

학습을 위한 명확한 정답 X

클러스터링의 품질을 확정할 수 있는 기준을 정의해야 함

## Which one is better?

V



# How to Measure Clustering Quality

- Measures that do not require ground truth labels

- Inertia

凝聚means에서 적용

- Within-cluster sum-of-squares

값↑ 더 좋은 클러스터링

클러스터 내의 오차 측정

$$\sum_{i=0}^n \min_{\mu_j \in C} \|x_j - \mu_i\|^2$$

각 데이터 포인트가 가장 가까운 클러스터 중심과  
얼마나 가까운지

- Silhouette Coefficient

형상지수

- $s(i)$ : Silhouette coefficient of  $i$ -th sample

- $a(i)$ : The mean distance between a sample and all other points in the same class

- $b(i)$ : The mean distance between a sample and all other points in the next nearest cluster

에가기울수록↑

클러스터간

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

음수면 잘못된 클러스터

$$0 \leq s(i) \leq 1$$

- Overall clustering quality can be obtained by averaging  $s(i)$  for all samples

# How to Measure Clustering Quality

- Clustering performance evaluation measure
    - ▣ Homogeneity: each cluster contains only members of a single class
- 정의
- $$h = 1 - \frac{H(C|K)}{H(C)}$$
- $H(C)$  is the entropy of the classes
  - $H(C|K)$  is the conditional entropy of the classes given the cluster assignments
- 정의한 정의에 여러 클래스가  
섞여있으면  $\rightarrow 1$
- $$H(C) = - \sum_{c=1}^{|C|} \frac{n_c}{n} \cdot \log\left(\frac{n_c}{n}\right)$$
- $$H(C|K) = - \sum_{c=1}^{|C|} \sum_{k=1}^{|K|} \frac{n_{c,k}}{n} \log\left(\frac{n_{c,k}}{n_k}\right)$$
- $n$  is the total number of samples,  $n_c$  and  $n_k$  are the number of samples respectively belonging to class  $c$  and cluster  $k$
  - $n_{c,k}$  is the number of samples from class  $c$  assigned to cluster  $k$
- ▣ Completeness: all members of a given class are assigned to the same cluster

$$c = 1 - \frac{H(K|C)}{H(K)}$$

# How to Measure Clustering Quality

- Clustering performance evaluation measure

- ▣ Adjusted Rand Index(ARI)

- Given the knowledge of the ground truth class assignments and our clustering algorithm assignments of the same samples, the adjusted Rand index is a function that measures the similarity of the two assignments

$$ARI = \frac{RI - \mathbb{E}[RI]}{\max(RI) - \mathbb{E}[RI]}$$

- $C$  is a ground truth class assignment,  $K$  is the clustering
  - $a$  is the number of pairs of elements that are in the same set in  $C$  and in the same set in  $K$
  - $b$  is the number of pairs of elements that are in different sets in  $C$  and in different sets in  $K$
  - Raw Rand index  $RI = \frac{a+b}{C_2^n}$  ( $C_2^n$  is the total number of possible pairs in the dataset)

$$C_2^n = \frac{n!}{2! (n-2)!}$$

# How to Measure Clustering Quality

- Contingency table

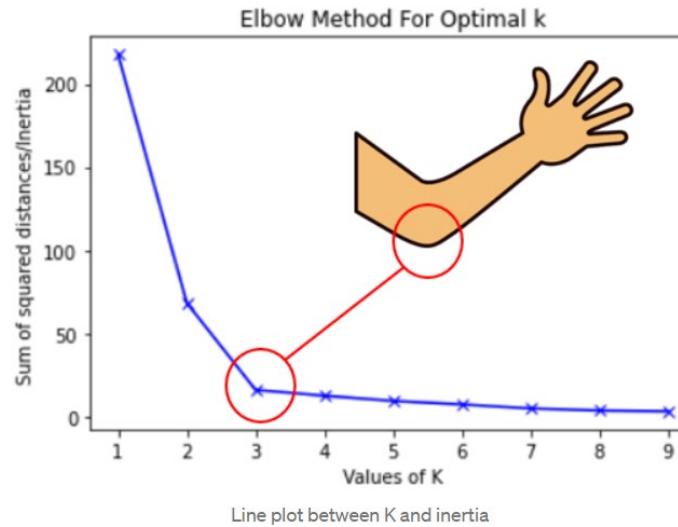
	$K_1$	$K_2$	...	$K_s$	$sums$
$C_1$	$n_{11}$	$n_{12}$	...	$n_{1s}$	$a_1$
$C_2$	$n_{21}$	$n_{22}$	...	$n_{2s}$	$a_2$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$C_r$	$n_{r1}$	$n_{r2}$	...	$n_{rs}$	$a_r$
$sums$	$b_1$	$b_2$	...	$b_s$	$n$

$$ARI = \frac{RI - \mathbb{E}[RI]}{\max(RI) - \mathbb{E}[RI]} = \frac{\sum_{ij} \binom{n_{ij}}{2} - \left[ \sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[ \sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \left[ \sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}$$

$$\binom{n}{k} = {}_n C_k$$

# Determining the Optimal Number of Clusters

- Elbow method
  - ▣ Plot the total within-cluster sum of squares against the number of clusters
  - ▣ Look for an “elbow” point where the rate of decrease sharply slows



- Silhouette Score
  - ▣ Measures how similar an object is to its own cluster compared to other clusters
  - ▣ Value ranges from -1 to 1 and a higher value indicates better clustering