

A background image of a kitchen wall. At the top, a wooden shelf holds several teapots and cups in various colors like blue, red, and white. Below the shelf, a small framed picture hangs on the wall. To the left, a wooden knife block and a large abstract painting are visible. The text 'Machine Learning' is overlaid in white, and 'K-Means Algorithm' is overlaid in yellow, both between two horizontal red lines.

Machine Learning

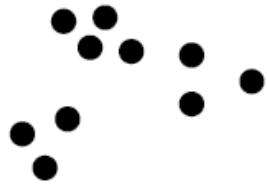
K-Means Algorithm

김선녕(ksycafe@gmail.com)

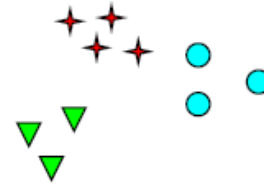
- 지식(정보)가 없는 상태에서 유사(similar)한 데이터를 모아 같은 그룹으로 묶는 일
 - 군집간 분산 최대화, 군집 내 분산 최소화
 - 유사도(similarity) 계산
- 레이블이 없는 대표적인 **비지도학습(unsupervised learning)**
 - 분류(classification) : 지도학습(supervised learning)
- Hierarchical Clustering 과 Partitional Clustering 가 있다.

분할 군집화(Partitional Clustering)

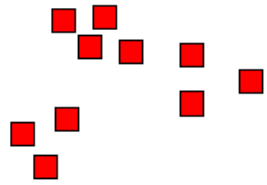
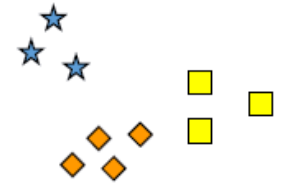
3



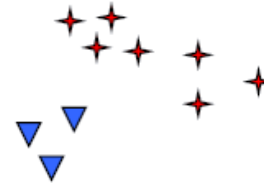
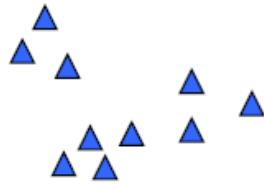
How many clusters?



Six Clusters



Two Clusters



Four Clusters



계층형 군집화(Hierarchical Clustering)

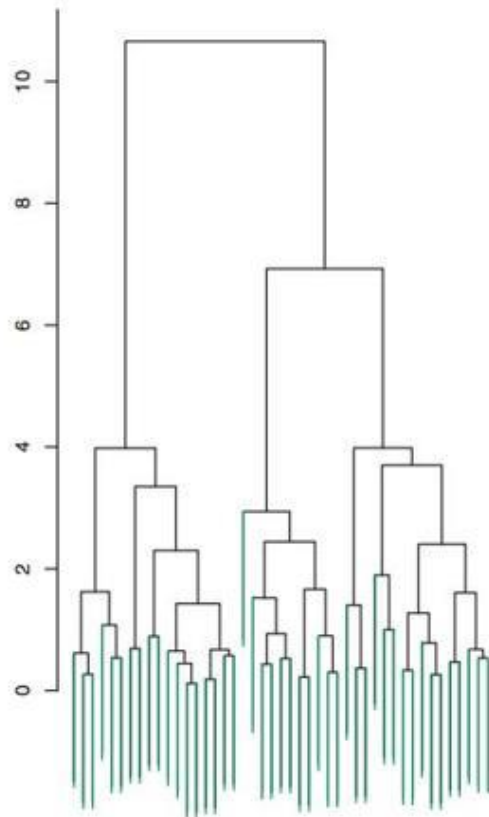
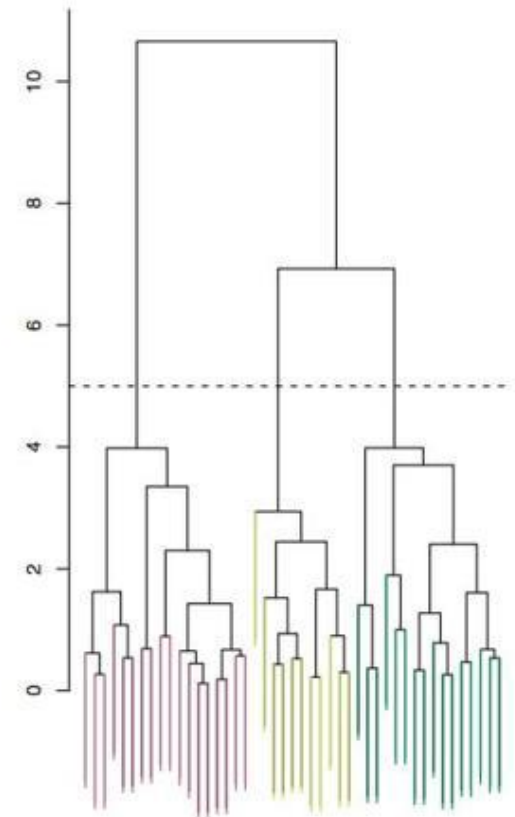
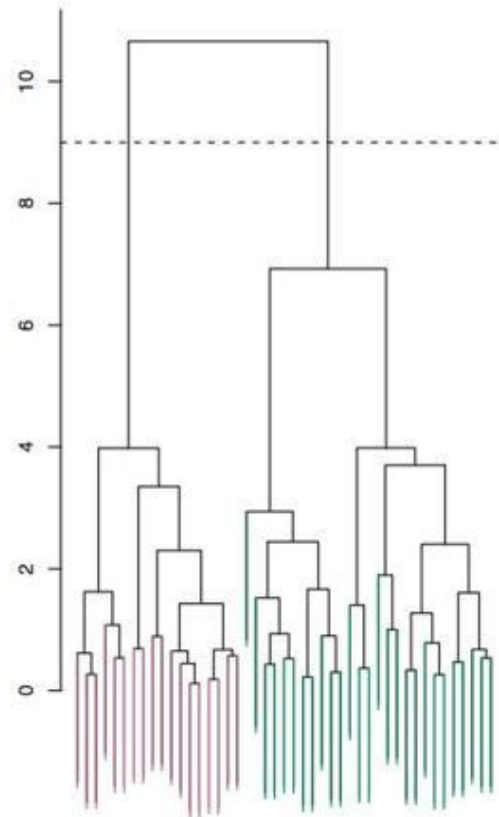
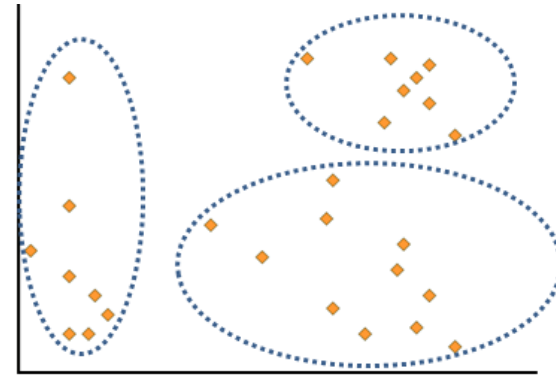
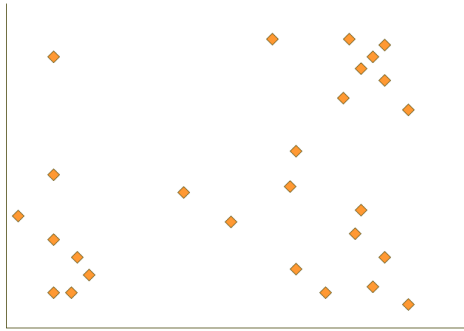


Figure 1







정지영상을 비슷한 것끼리 군집화



영상 분할(segmentation)

k -평균 알고리즘(k -means algorithm)

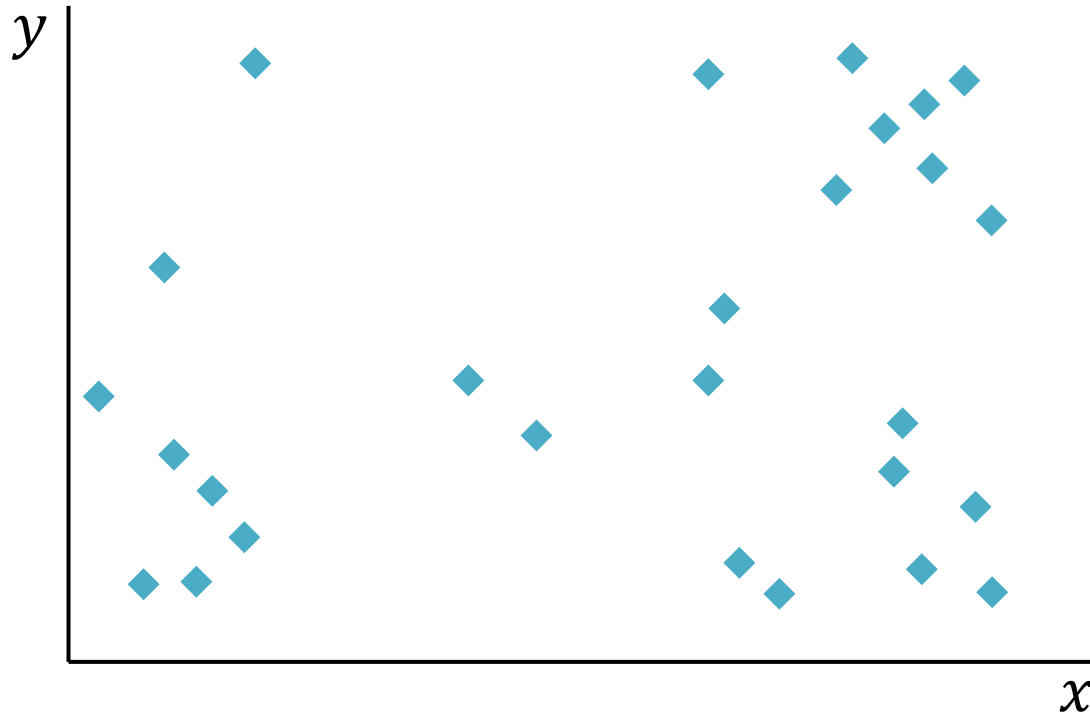
7

입력 : 훈련집합 $X = \{x_1, x_2, \dots, x_n\}$, 군집의 개수 k
출력 : 군집집합 $C = \{c_1, c_2, \dots, c_k\}$

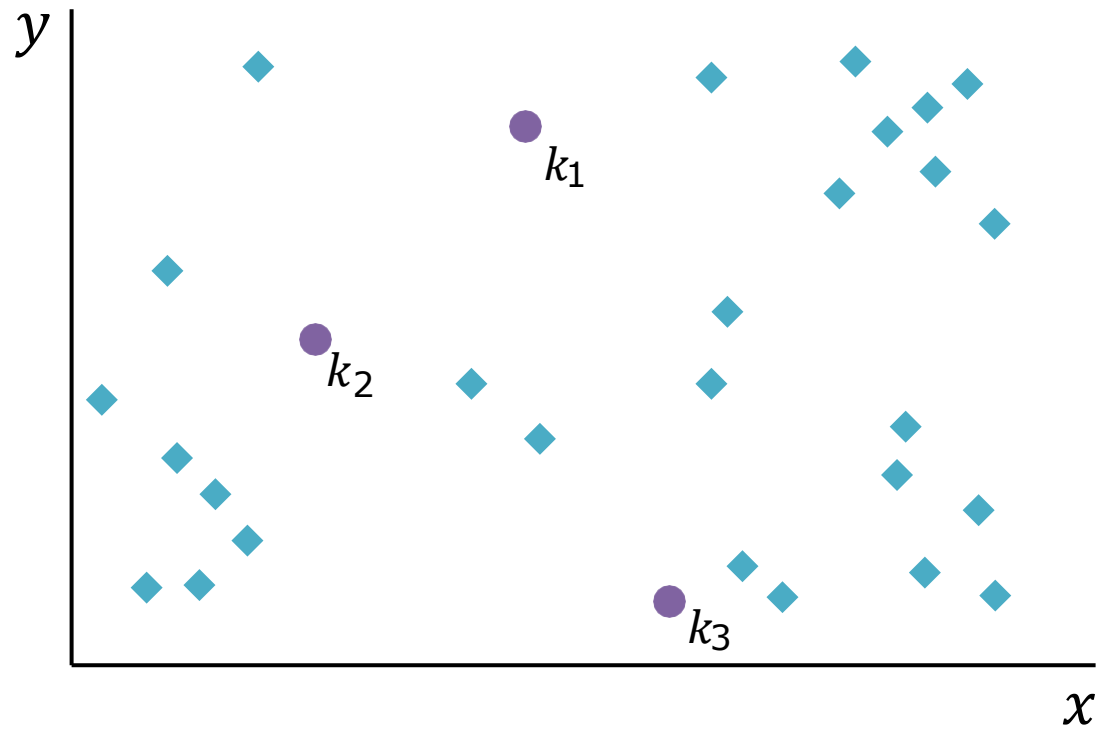
[pseudo code]

- 1: k 개의 점들을 초기 중심점(*centroid*)으로 선택한다.
- 2: *repeat*
- 3: 각 점을 가장 가까운 중심점(*centroid*)에 할당하여 k 개의 군집을 형성한다.
- 4: 각 군집의 중심점을 다시 계산한다.
- 5: *until* 중심점이 바뀌지 않을 때까지

Suppose we wish to cluster these items

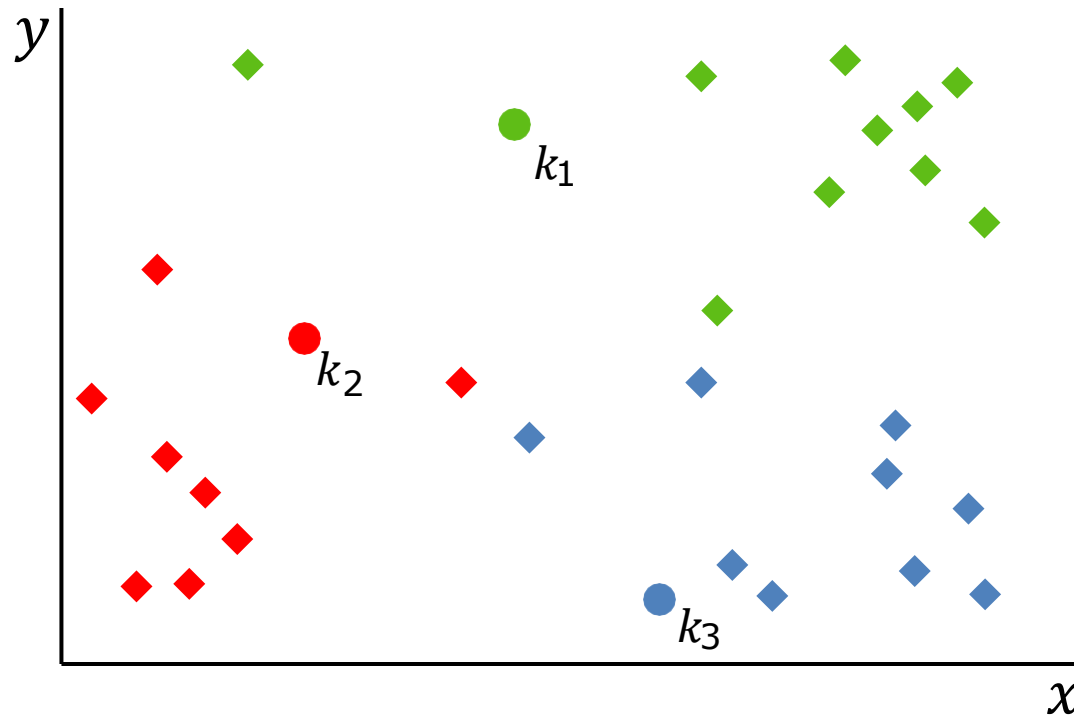


Pick 3 initial cluster centers at random



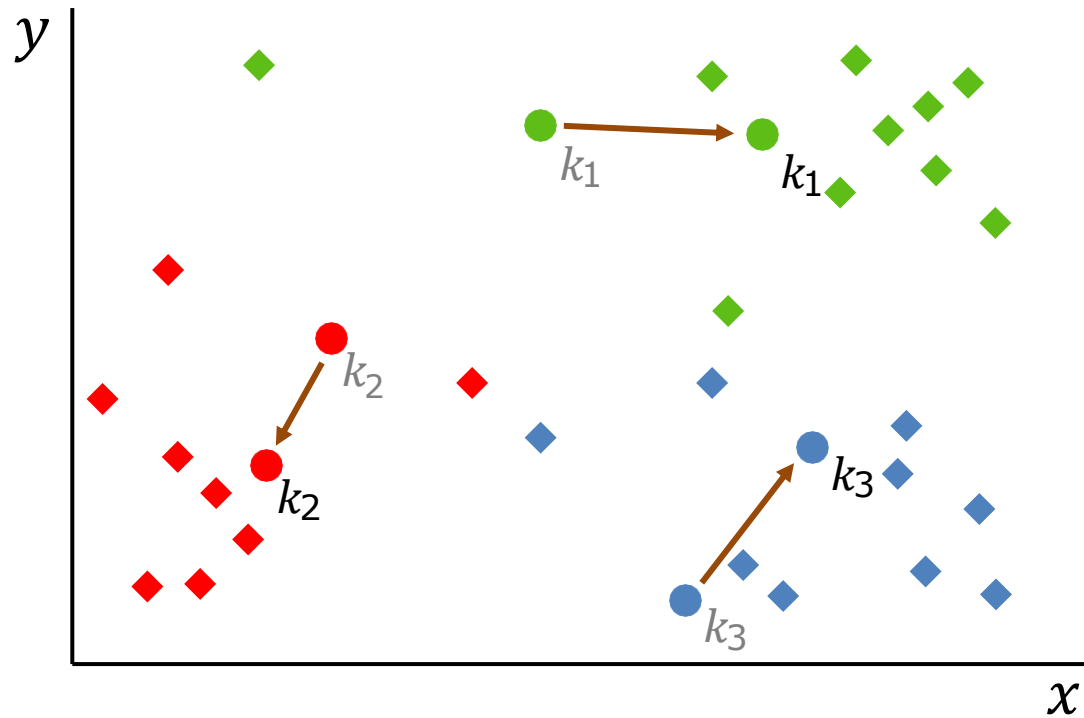
Assign each instance to the closest cluster center

10



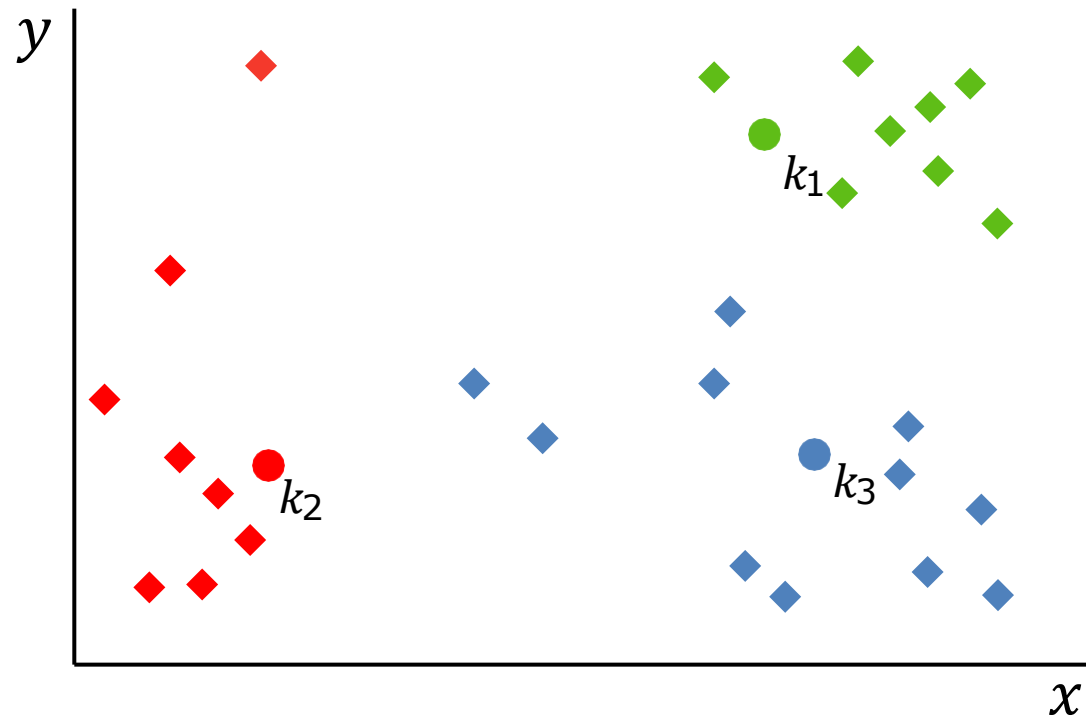
Move each cluster center to the mean of each cluster.
Reassign points closest to a different new cluster center.

11

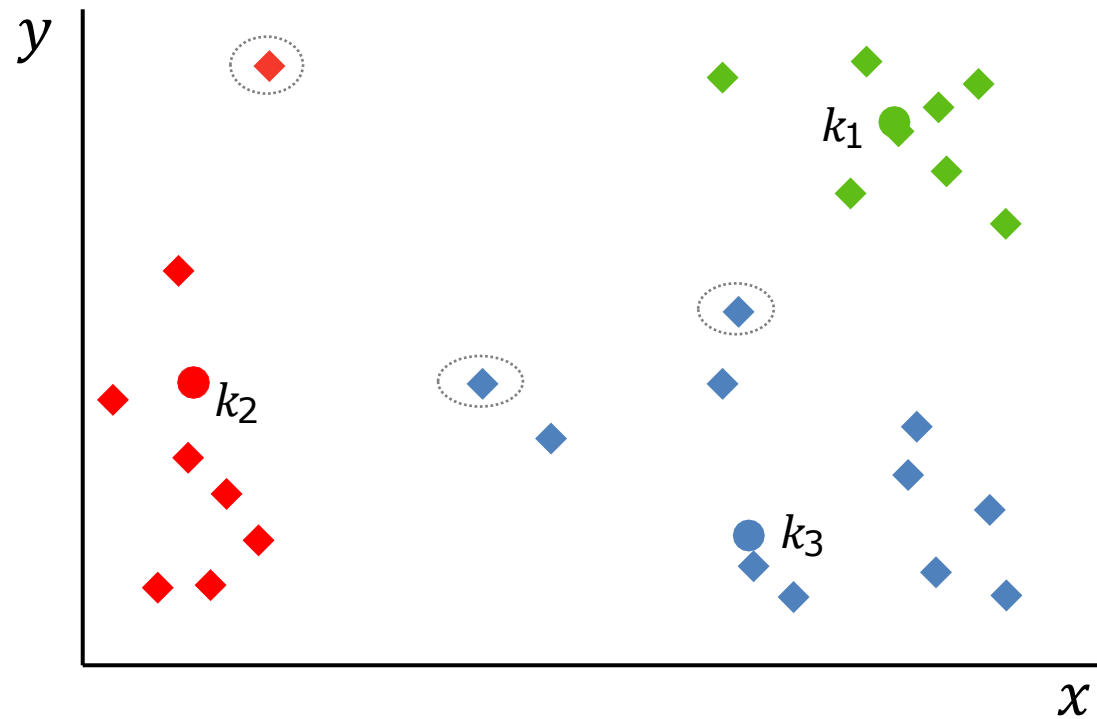


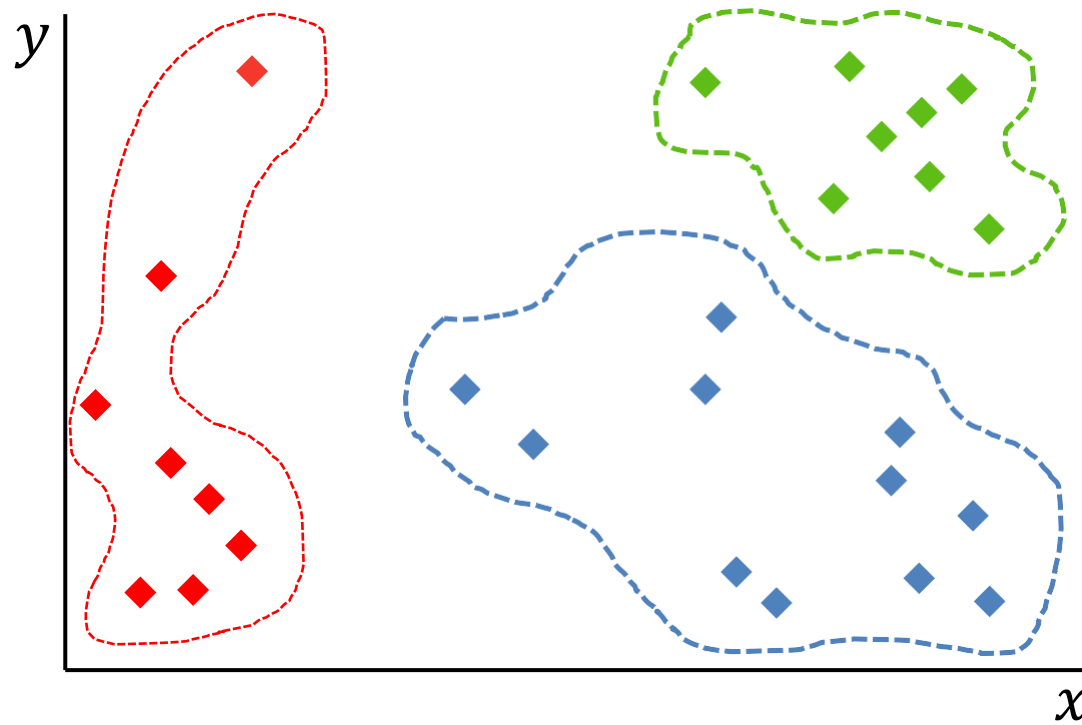
Recompute cluster means

12



No change. Done!





- 목적함수(*objective function*)
 - 점들 사이 또는 점과 군집 중심점 간의 인접성에 의존
 - 클러스터의 품질은 오차 제곱의 합(*ESS : Error Sum of Squares*)으로 평가

$$ESS = \sum_{i=1}^k \sum_{x \in c_i} dist(c_i, x)^2$$

i 번째 군집의 중심점(평균)은

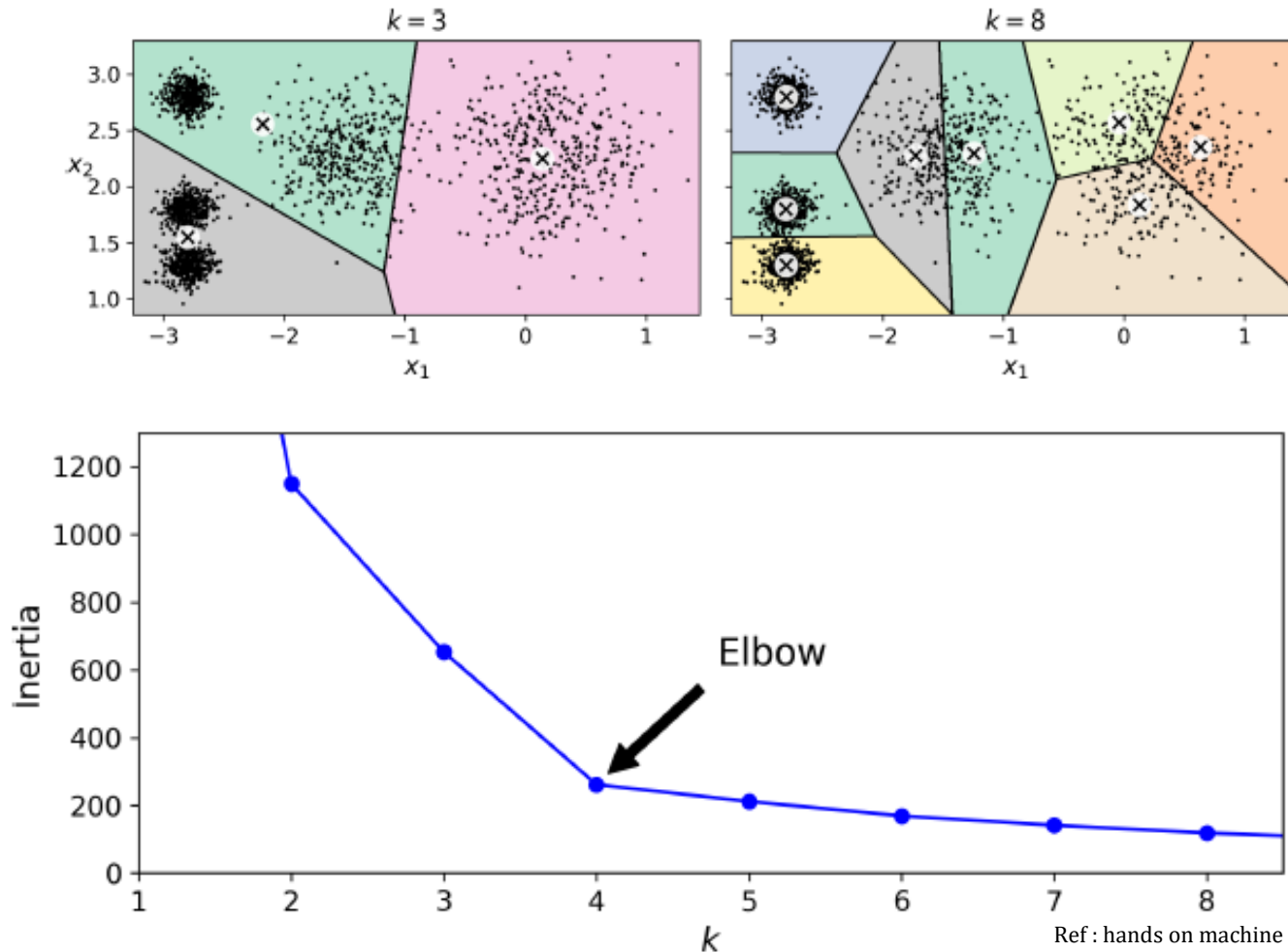
$$c_i = \frac{1}{m_i} \sum_{x \in c_i} x$$

- SSE 의 최소화하는 k 값을 선택하는 것이 중요

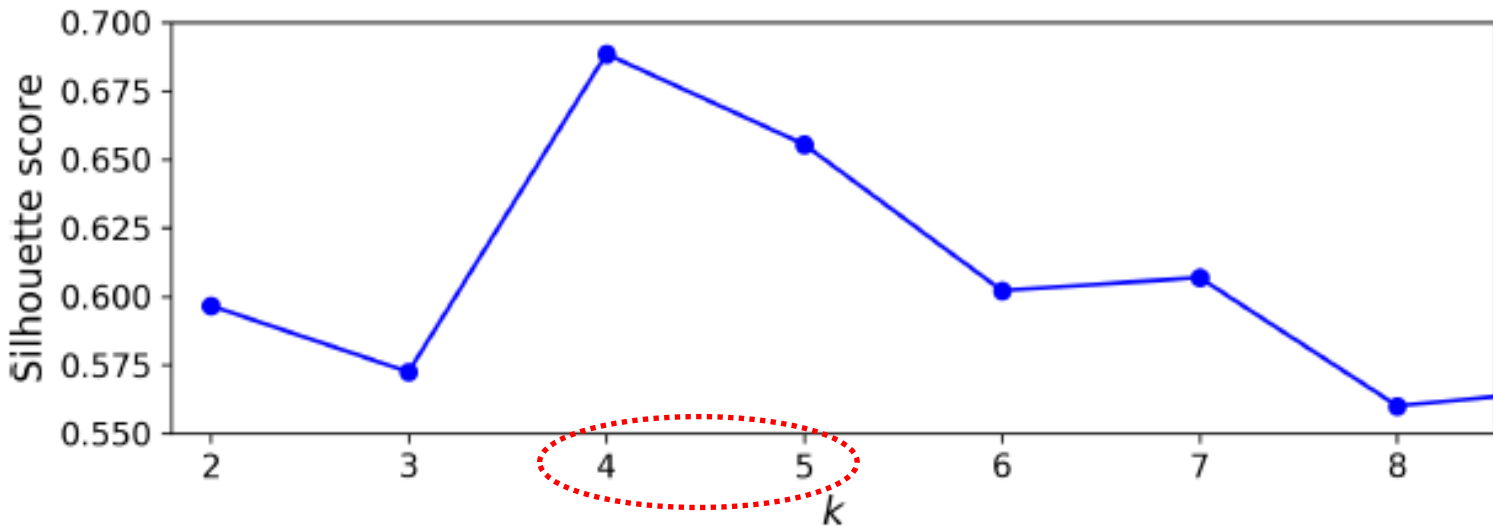
잘못된 클러스터 개수 선택

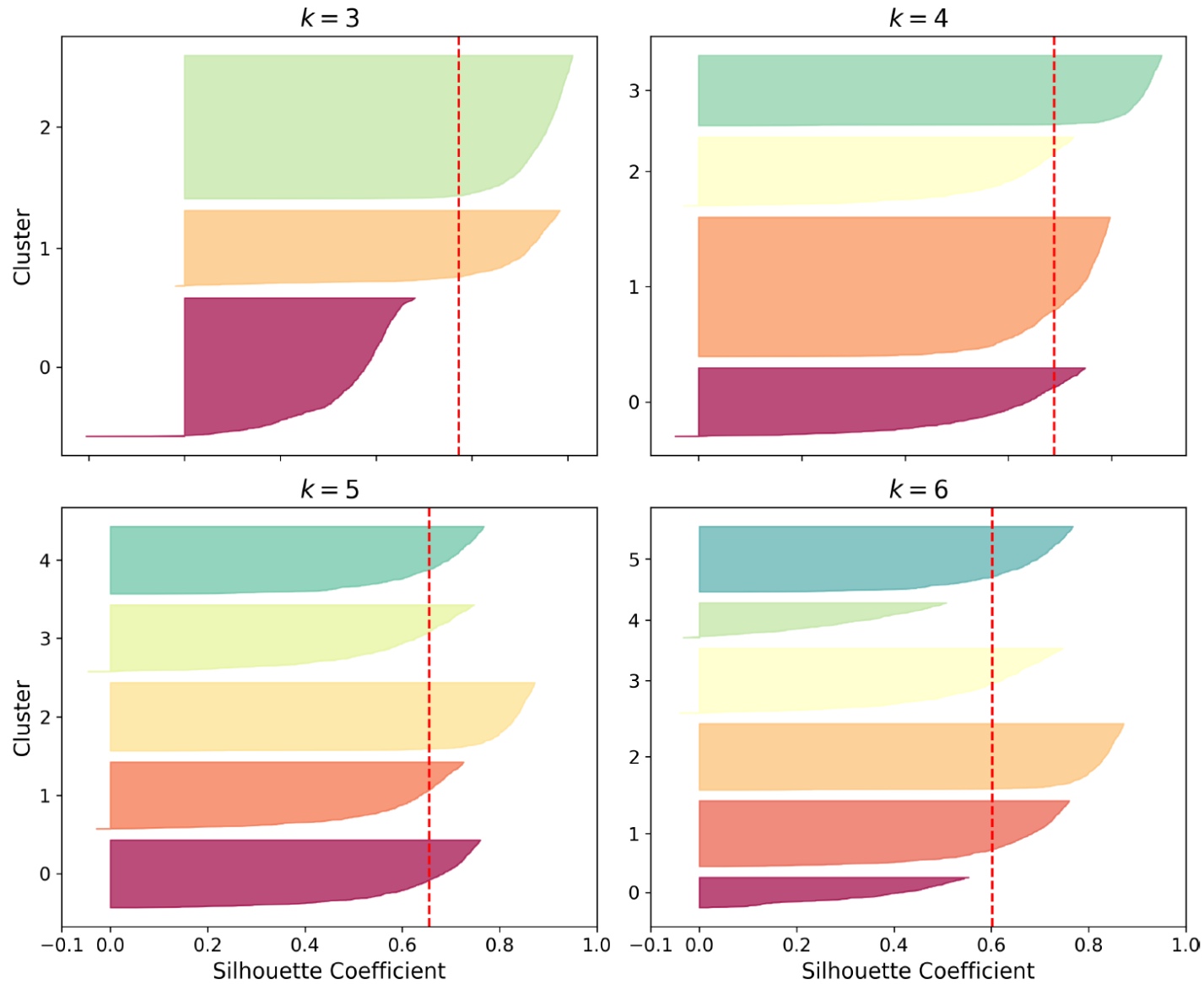
16

- k 가 너무 작으면 별개의 클러스터를 합치고 k 가 너무 크면 하나의 클러스터가 여러 개로 나뉜다
- 이너셔(inertia) : 각 샘플과 가장 가까운 센트로이드사이의 평균제곱거리

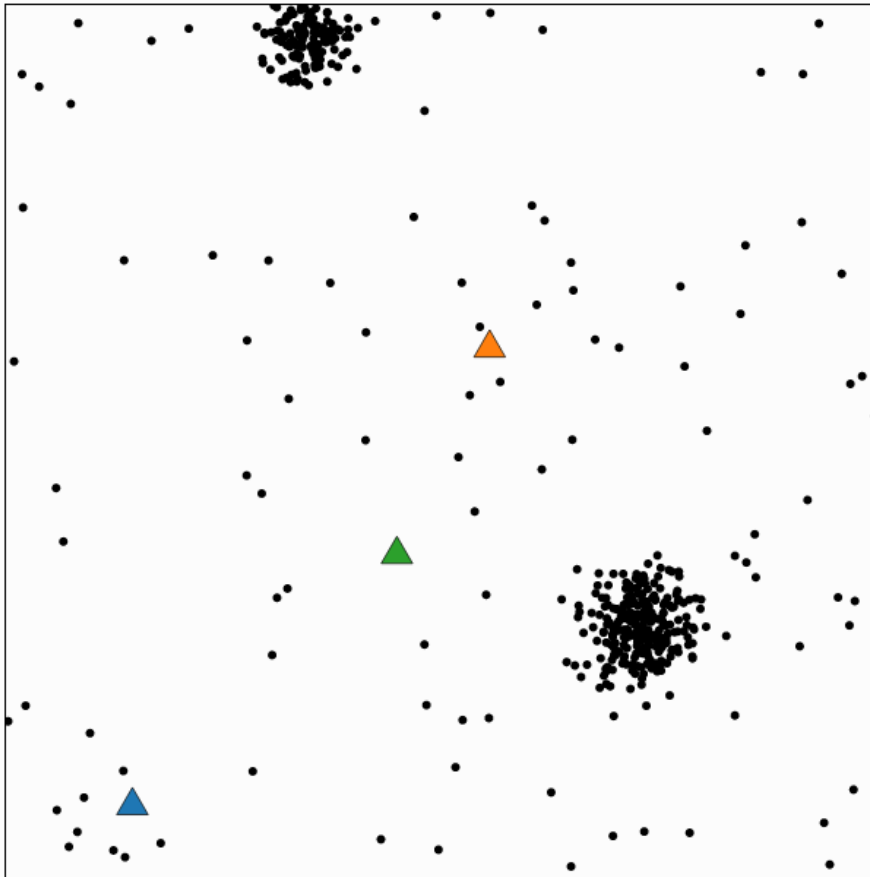


- 실루엣 계수 = $(b - a) / \max(a, b)$
 - a : 동일한 클러스터에 있는 다른 샘플까지 평균거리(클러스터 내부의 평균 거리)
 - b : 가장 가까운 클러스터까지 평균거리(가장 가까운 클러스터의 샘플까지 평균 거리)
 - +1에 가까우면 자신의 클러스터안에 잘 속해 있고 다른 클러스터와는 멀리 떨어져 있다는 뜻.
 - 0에 가까우면 클러스터 경계에 위치한다는 의미
 - -1에 가까우면 잘못된 클러스터라는 의미
- 실루엣 점수 : 실루엣 계수의 평균





<https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html>



Mean square point-centroid distance: not yet calculated

The k -means algorithm is an iterative method for clustering a set of N points (vectors) into k groups or clusters of points.

Algorithm

Repeat until convergence:

Find closest centroid

Find the closest centroid to each point, and group points that share the same closest centroid.

Update centroid

Update each centroid to be the mean of the points in its group.

Find closest centroid

Data

Clustered points ☐ Random

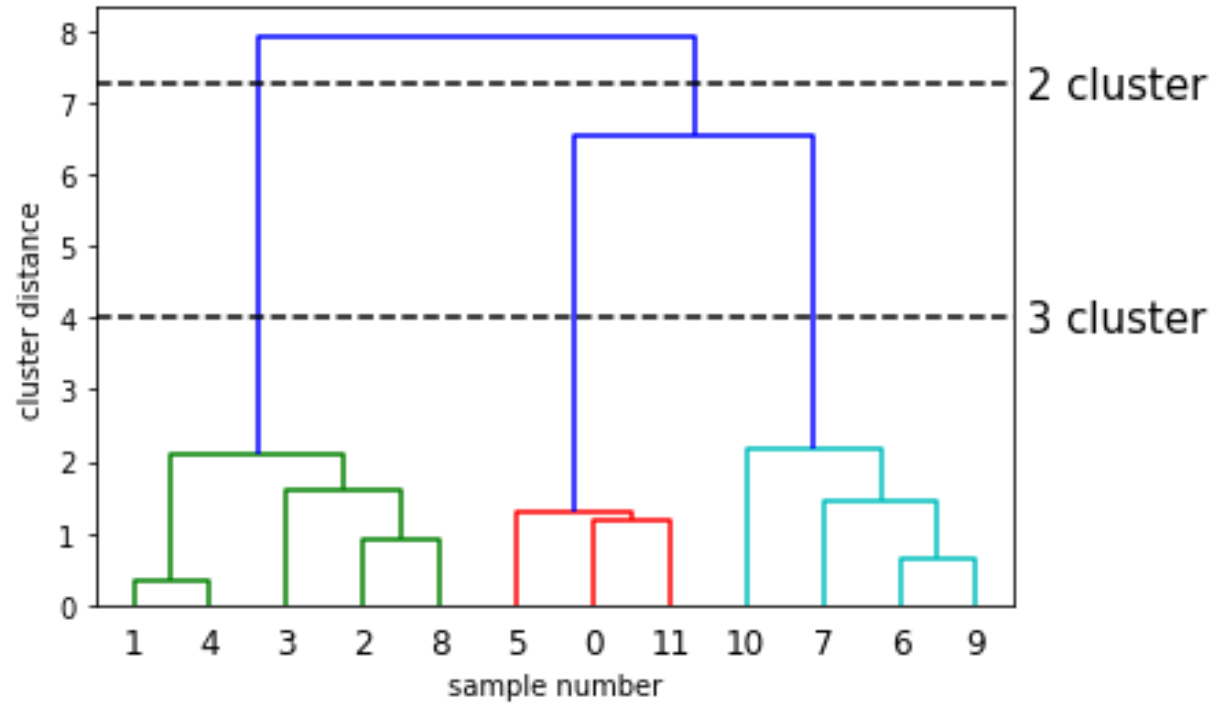
Number of clusters :

Number of centroids :

New points

New centroids

계층형 군집화(Hierarchical Clustering)



DBSCAN(*Density-Based Spatial Clustering of Application with Noise*)

21

