

k-means clustering.

- cluster의 개수 → 우리가 미리 설정.
- cluster의 개수를 어떻게 선택할까?
- 목적이나 더 좋은 결과를 뽑는지.
↳ measure, 척도.

Inclass 16: Model Evaluation and Selection *

[SCS4049] Machine Learning and Data Science

Seongsik Park (s.park@dgu.edu)

Department of Artificial Intelligence, Dongguk University

k-means clustering.

- ① inertia (cost).
- ② silhouette score.

#cluster
개수.

- 중간고사 점수 → 2개월(5), 가령, eclass.
- 이번주 과제 → 구글 설문, 강의 의견, 전반적 의견 제출. eclass.
기명문 (제출여부 확인) 이음. 학번입력 한번은.
- 쉬는시간 → 음악 · MV · live 틀어볼까? 극찬곡 · 영상. 제보.

K-means clustering: inertia

A performance metric, called inertia

- K-means clustering minimize cost.

inertia = the mean squared distance
between each instance and its closest centroid

각각의 샘플들.

그 샘플에서 가장 가까운 평균군.

Elbow rule: any lower value would be dramatic, while any higher value would not help much.

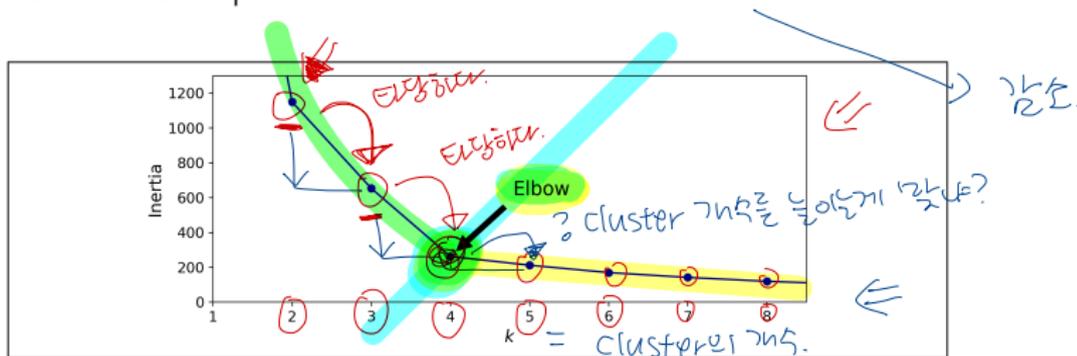
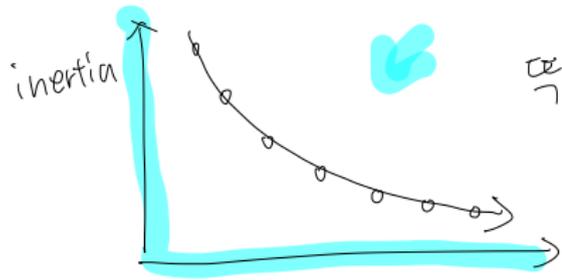


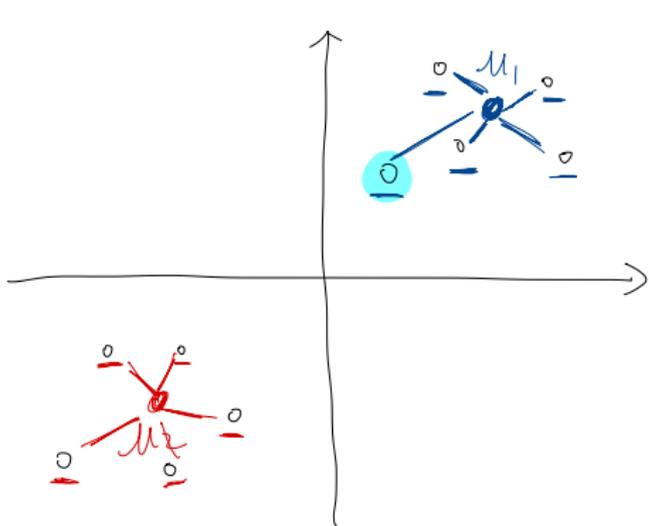
Figure 9-8. Selecting the number of clusters k using the “elbow rule”

k = 4 다, 경계점을 할 수 있음. 절대적이지는 않음.



때로는 $\frac{1}{k}$ 이 $\frac{1}{k+1}$ 이보다 클 수 있다.

clusterings



2D

cluster의 개수 = 2개.

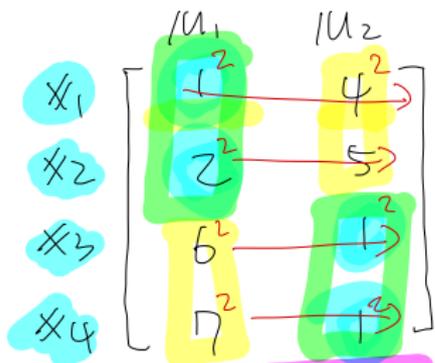
inertia

= 전체
샘플에 대한
평균

$$\left(\frac{\text{거리}^2}{\text{나갈 빈도} = \text{샘플 수}} \right)$$

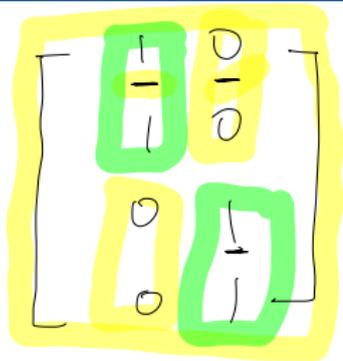
나갈 빈도 = 샘플 수
평균

$$\text{inertia} = \frac{1}{N} \sum_{n=1}^N \underbrace{r_{nk}}_{\begin{cases} 0 \\ 1 \end{cases}} \underbrace{\|x_n - \mu_k\|^2}_{\begin{array}{l} n\text{-번째 sample이} \\ k\text{-번째 cluster에 속했다면.} \end{array}}$$



$$\Rightarrow \begin{bmatrix} 1^2 \\ 2^2 \\ 6^2 \\ 7^2 \end{bmatrix}$$

영벡터



np.min(dist² axis=1)

0 inertia = $\sum_{n,k} r_{nk} \|x_n - \mu_k\|^2$

(클러스터 - 가장 가까운 클러스터 사이의 거리의 제곱.)

$$= \frac{1}{4} (1^2 + 2^2 + 6^2 + 7^2)$$

$$= \frac{1}{N} \sum r_{nk} \|x_n - \mu_k\|^2$$

K-means clustering: silhouette score

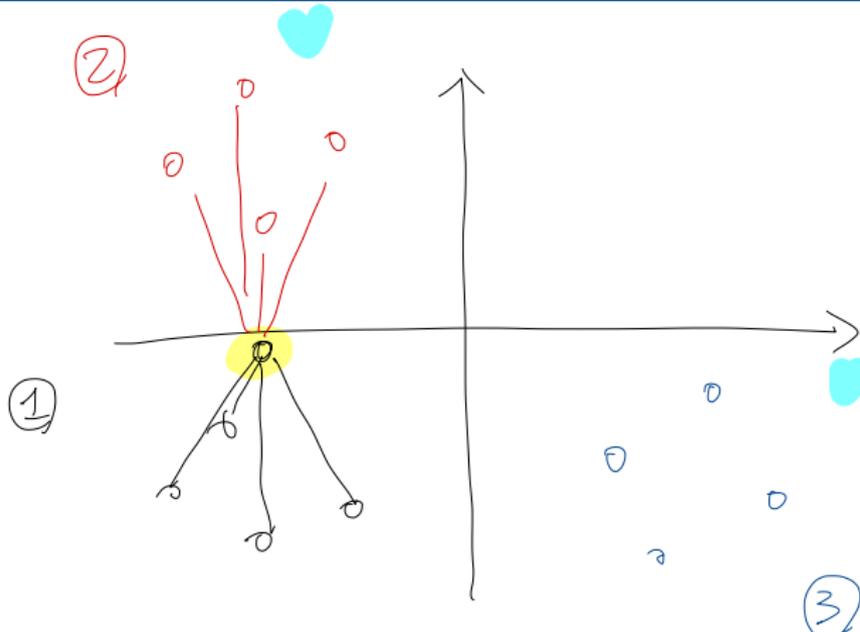
A more precise approach (but also more computationally expensive) is to use the silhouette score, which is the mean silhouette coefficient over all the instances. → 각각의 샘플에 대해 계산.

- An instance's silhouette coefficient is equal to $(b - a) / \max(a, b)$ +1
-1
- where a is the mean distance to the other instances in the same cluster (i.e., the mean intra-cluster distance) 나 = 샘플 VS. 나랑 같은 군에 속한 샘플들
- and b is the mean nearest-cluster distance (i.e., the mean distance to the instances of the next closest cluster, defined as the one that minimizes b , excluding the instance's own cluster). 나 = 샘플 VS. 다른 군에 속한 샘플들

실루엣 스코어 = 평균 (실루엣 coefficient) ↓

-1 ~ +1
사이인 실수.

실루엣 스코어 = 평균 (실루엣 coefficient)



$a = 4$ vs. 내가 속한 군의 샘플 수.

$b = 4$ vs. 두 번째 = 가장 가까운 cluster에 속한 샘플 수.

의 silhouette coeff =
$$\frac{b-a}{\max(a,b)}$$

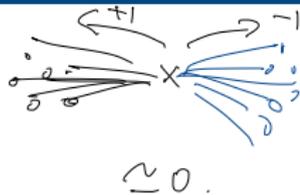
K-means clustering: silhouette score

The silhouette coefficient can vary between -1 and +1.

- A coefficient close to +1 means that the instance is well inside its own cluster and far from other clusters,
- while a coefficient close to 0 means that it is close to a cluster boundary,
- and finally a coefficient close to -1 means that the instance may have been assigned to the wrong cluster.

Silhouette
coeff =

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



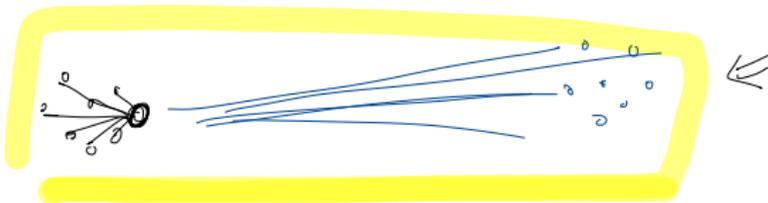
$a, b \geq 0.$

① $a \downarrow, b \uparrow =$

나는 너가속한애들이랑 뭉쳐있고 ($a \downarrow$)

그다음 가까운 cluster까지는 충분히 멀리 있음 ($b \uparrow$)

$$\frac{b - a}{\max(a, b)} \approx \frac{b}{b} = 1$$



② $a \gg b, a \uparrow, b \downarrow =$

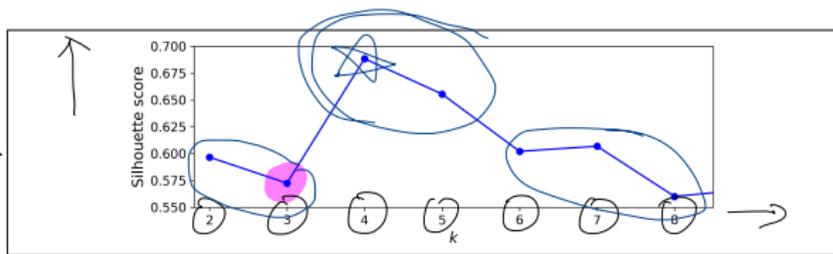
나는 너가속한애랑 멀고 ($a \uparrow$)

가장 이웃한 곳이란 거깝다 ($b \downarrow$)

$$\frac{b - a}{\max(a, b)} \approx \frac{-a}{a} = -1.$$

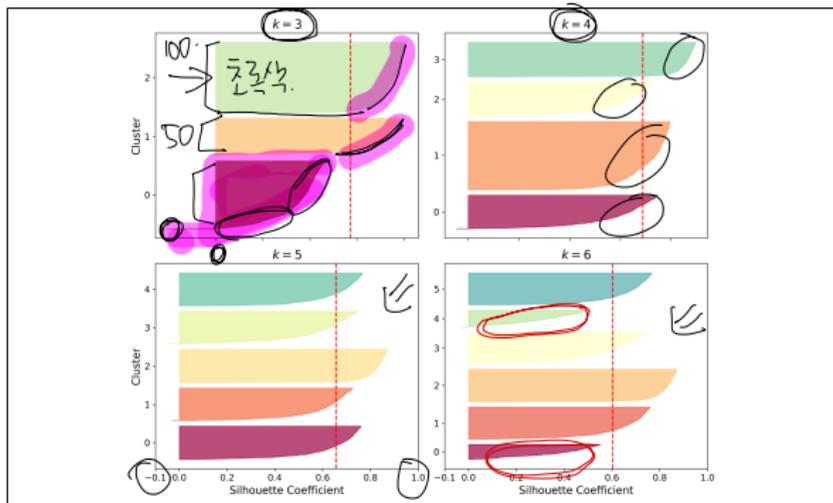
K-means clustering: silhouette score

Handwritten notes: $\frac{1}{2}, \frac{2}{2}, \frac{3}{2}, \frac{4}{2}, \frac{5}{2}, \frac{6}{2}$



Handwritten notes: $\frac{1}{2}, \frac{2}{2}, \frac{3}{2}$ score (정답이 3)

Figure 9-9. Selecting the number of clusters k using the silhouette score coefficient. This is called a *silhouette diagram* (see Figure 9-10):



Handwritten notes: $\frac{1}{2}, \frac{2}{2}, \frac{3}{2}$ coefficient.

Figure 9-10. Silhouette analysis: comparing the silhouette diagrams for various values of k