

~~이 단락~~ ~ clustering, model eval. selection.
kmeans, MoG. → # cluster.

Inclass 13: Model Evaluation and Selection

[SCS4049] Machine Learning and Data Science

Seongsik Park (s.park@dgu.edu)

이론 · python

AI Department, Dongguk University

kmeans : inertia, 실루엣

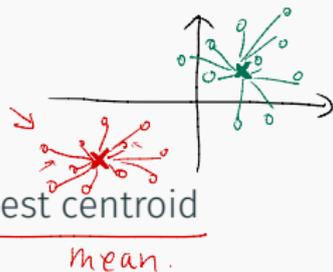
MoG : Information criterion ~ likelihood.

K-means clustering: inertia

A performance metric, called inertia 각 인스턴스의 중심점까지 거리의 제곱의 합

inertia = the mean (squared distance)

between each instance and its closest centroid



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

잘된 clustering \Rightarrow inertia 작아야.

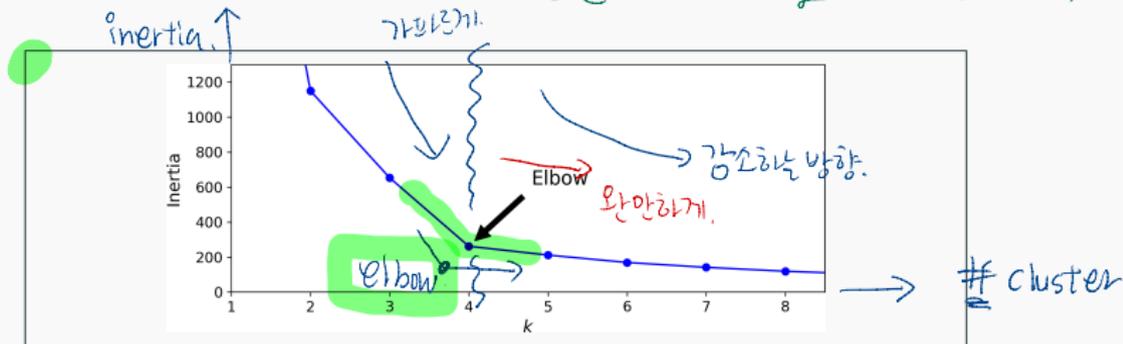


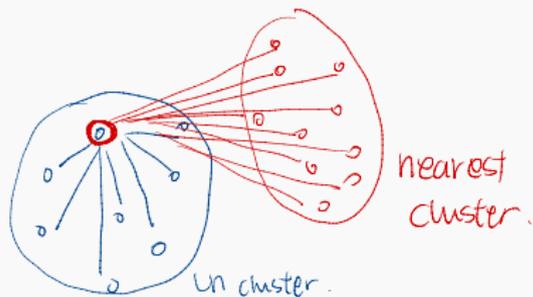
Figure 9-8. Selecting the number of clusters k using the "elbow rule"

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.

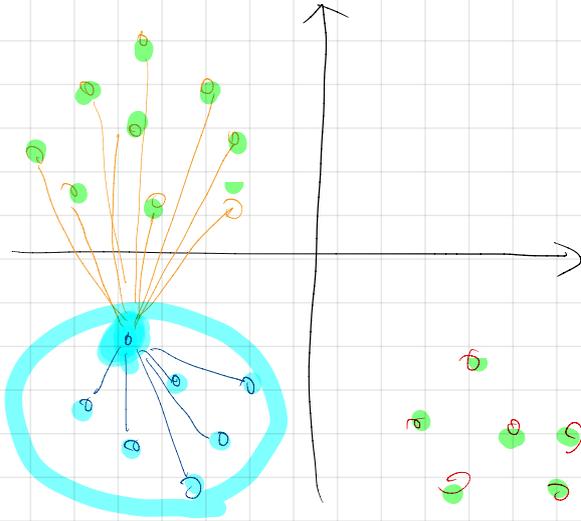
실루엣 계수 (coeff) = 각각의 sample (instance)의
가장 가까운 이웃.

- An instance's silhouette coefficient is equal to $(b - a) / \max(a, b)$
- where a is the mean distance to the other instances in the same cluster (i.e., the mean intra-cluster distance)
- 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).



a = 같은 cluster의 다른 sample까지의 거리.

b = 그 다음 가까운 cluster의 sample까지의 거리.



a = 피라미트 거리의 평균

b = 즉흥상자 거리의 평균

이 때의 실수 coeff

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

하나의 클러스터에 가까울수록, 자음 크기 비율
 $b \gg a$

주요 경우: $a \approx 0$ $b \gg 0$ \Rightarrow $\frac{b-a}{\max(a,b)} \approx \frac{b}{b} = +1$

안주 경우: $a \gg b$ $b \approx 0$ \Rightarrow $\frac{b-a}{\max(a,b)} \approx -\frac{a}{a} = -1$

하나의 클러스터와 다른 클러스터에 더 가깝다.

K-means clustering: silhouette score

The silhouette coefficient can vary between -1 and +1. $[-1, +1]$ 실수.

- 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, $a \approx b$.
- and finally a coefficient close to -1 means that the instance may have been assigned to the wrong cluster.

$$\text{Silhouette score} = \frac{1}{N} \sum_{n=1}^N \text{Silhouette coeff.}$$

K-means clustering: silhouette score

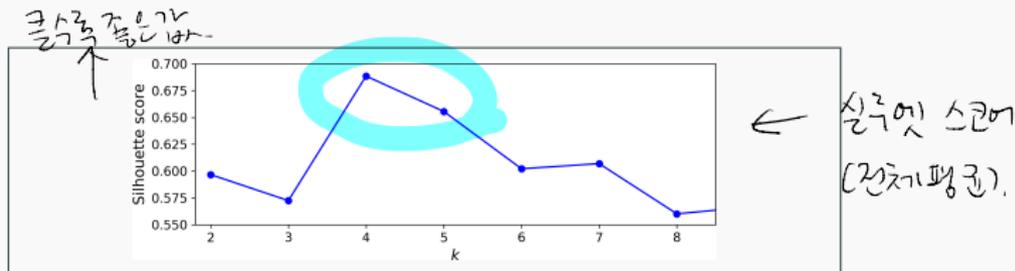


Figure 9-9. Selecting the number of clusters k using the silhouette score coefficient. This is called a silhouette diagram (see Figure 9-10):

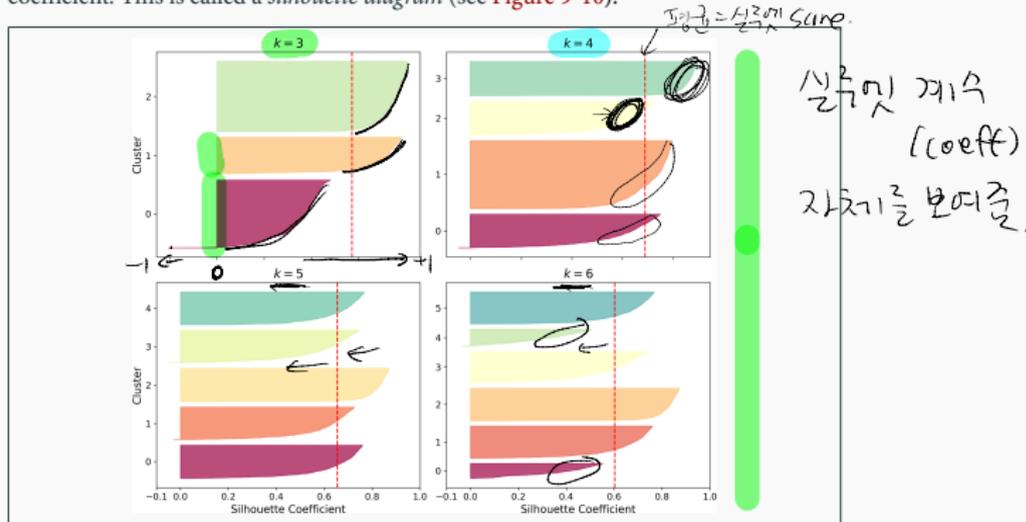


Figure 9-10. Silhouette analysis: comparing the silhouette diagrams for various values of k

MoG clustering: information criterion

- Akaike information criterion (AIC)

model complexity.

작은게 좋은거 정의.

$$\text{AIC} = -2 \cdot \frac{1}{N} \cdot \log \Pr(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta) + 2 \cdot \frac{d}{N} \quad (1)$$

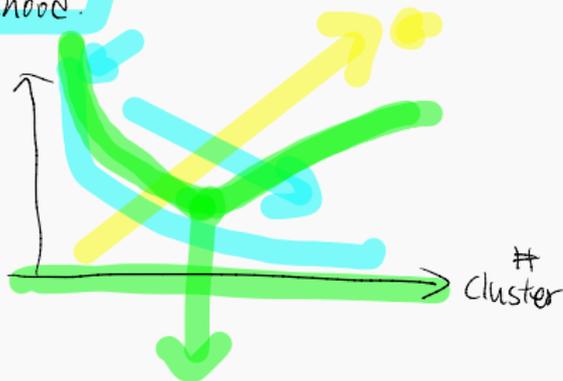
- Bayesian information criterion (BIC)

MoG 항함 $\theta = \{\{\pi_n, \mu_n, \sigma_n\}\}$

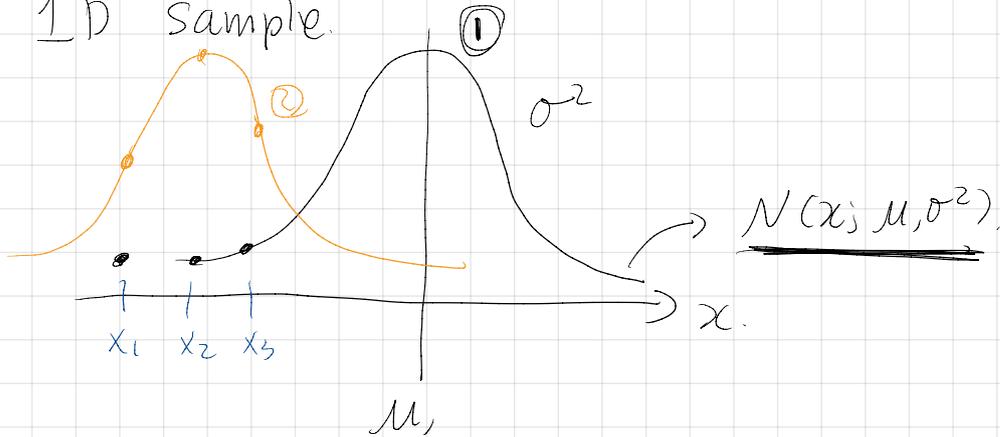
$$\text{BIC} = -2 \cdot \frac{1}{N} \cdot \log \Pr(\mathbf{x}_1, \dots, \mathbf{x}_N | \theta) + d \cdot \log N \quad (2)$$

log-likelihood.

- log-likelihood = -2가 클수록 좋은거
= 작을수록 좋은거



1D Sample.

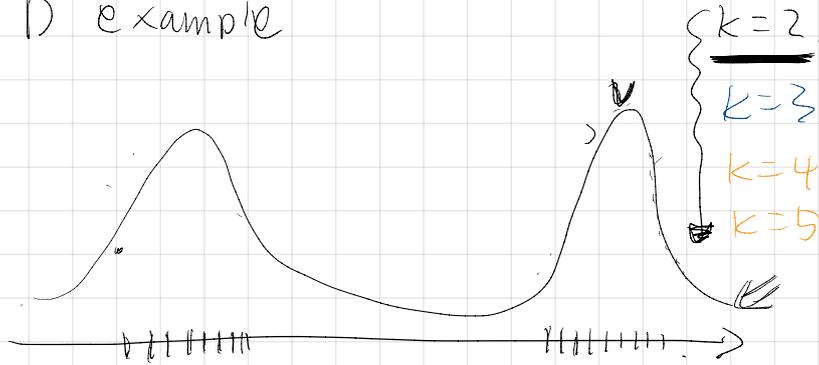


$$\textcircled{1} \mathcal{L}(\mu, \sigma^2; x_1, x_2, x_3) = \prod N(x_n; \mu, \sigma^2),$$

$$\textcircled{2} \mathcal{L}(\mu, \sigma^2; x_1, x_2, x_3) = \prod N(x_n; \mu, \sigma^2) \leftarrow \exists \mu, \sigma^2$$

1D example

model eval, selection



① $k=2$, $f(x) = f(x; \pi_1, \pi_2, \mu_1, \mu_2, \sigma_1, \sigma_2)$

$$\mathcal{L}(\pi_1, \pi_2, \mu_1, \mu_2, \sigma_1, \sigma_2; x_1, x_2, \dots, x_N) = \prod_{n=1}^N f(x_n)$$

② $k=3$ $f(x) = f(x; \pi_1, \pi_2, \pi_3, \mu_1, \mu_2, \mu_3, \sigma_1, \sigma_2, \sigma_3)$

$$\mathcal{L}(\dots; x_1, \dots, x_N) = \prod f(x_n)$$