



K-means and K-Nearest Neighbor

Kwangwoon University MI:RU Machine Learning Study
오민성

Contents

K-means

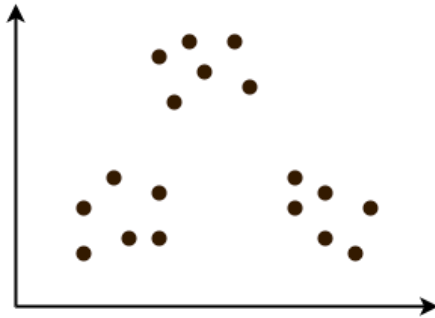
- Algorithms
- How to find the optimal K
- Problems

K-nearest neighbor

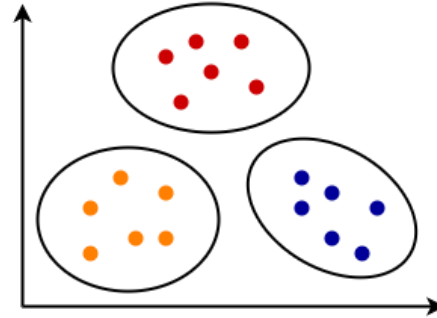
- Algorithms
- How to find the optimal K
- Problems

What algorithms is **K-means**?

- ▷ Clustering
- ▷ Unsupervised Learning



Before K-Means

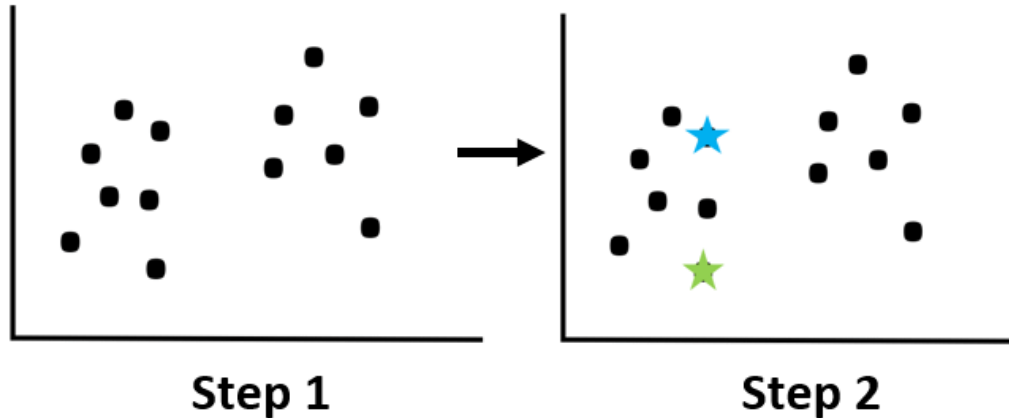


After K-Means

K-Means Clustering Algorithms

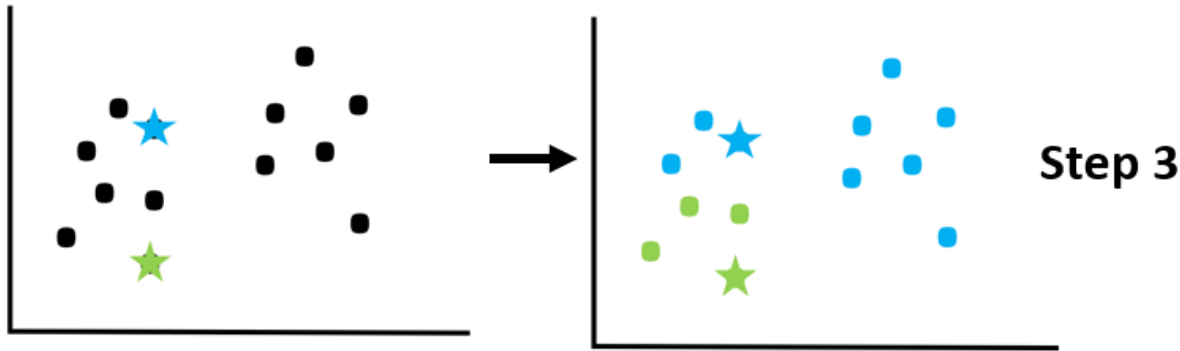
Step 1. Choose the number of clusters ($K = 2$)

Step 2. mark 2 data points randomly for 2 clusters.



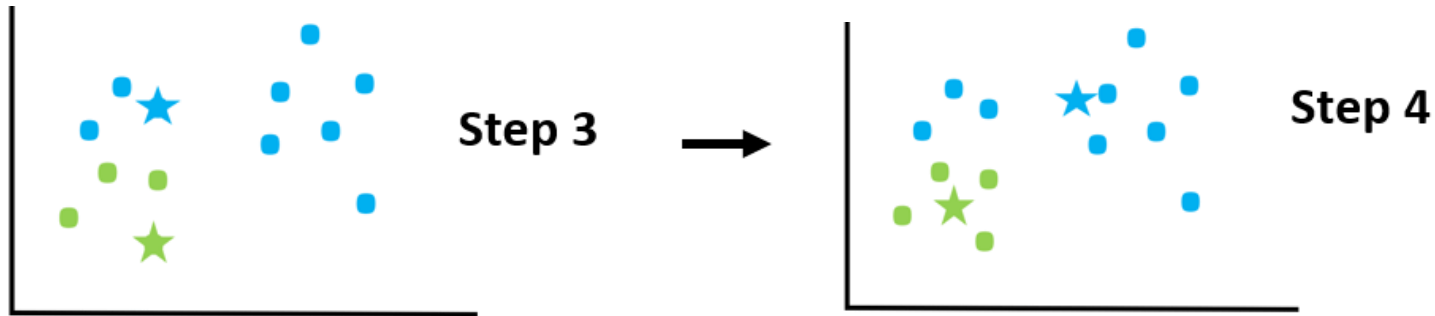
K-Means Clustering Algorithms

Step 3. Assign each data points to the nearest centroid.



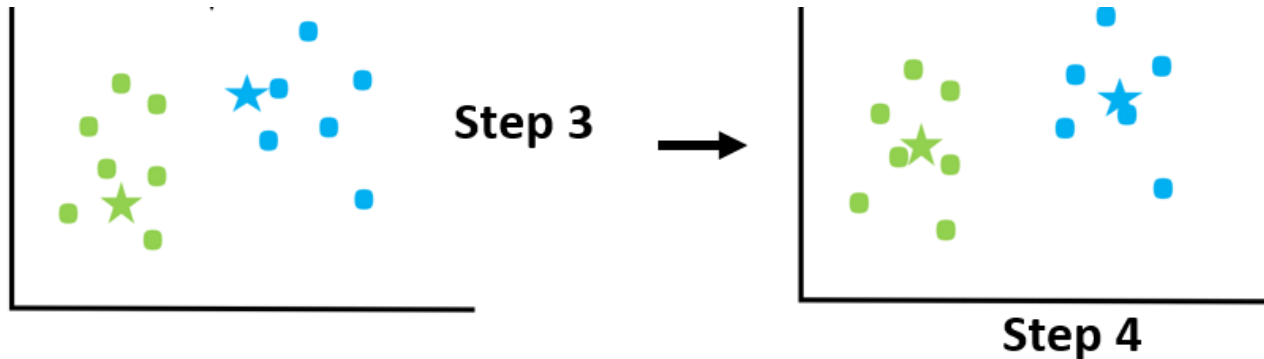
K-Means Clustering Algorithms

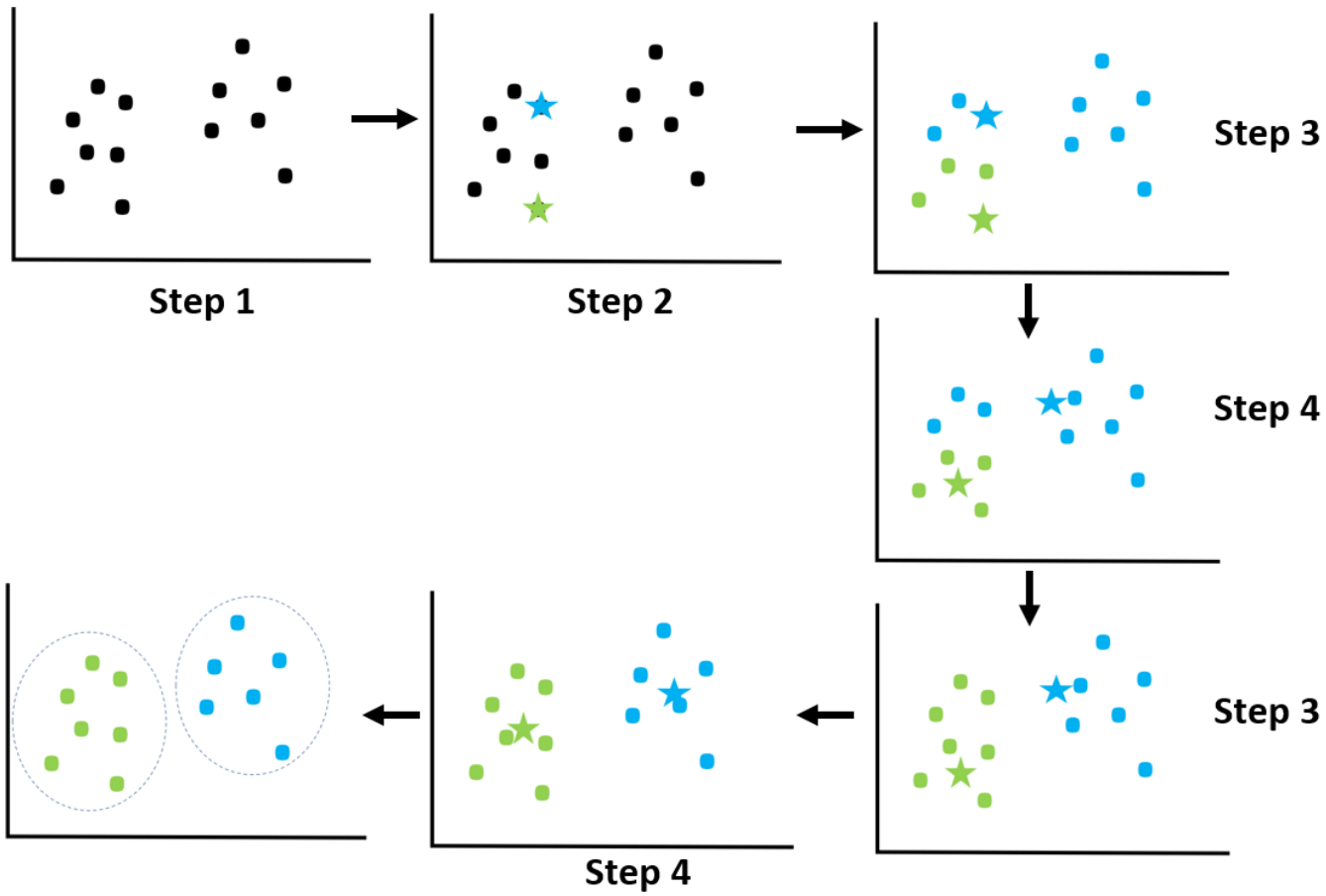
Step 4. Calculate the average of the data points belonging to each cluster and move its cluster centroid to the average location.

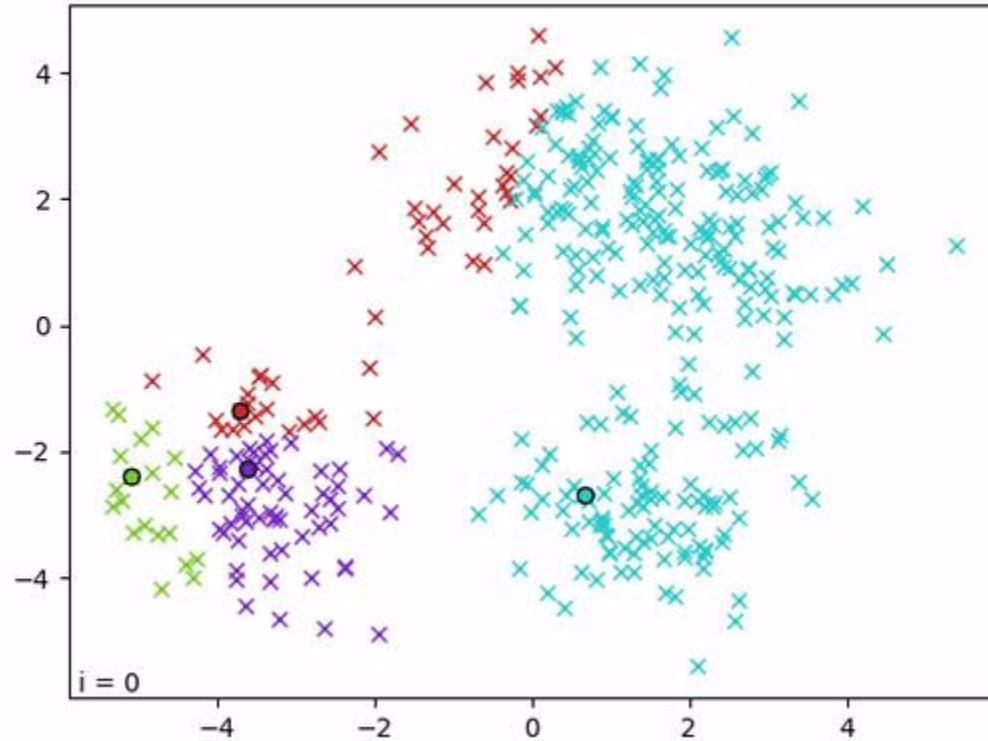


K-Means Clustering Algorithms

Step 5. Repeat Step 3 and Step 4 until cluster centroids do not change.

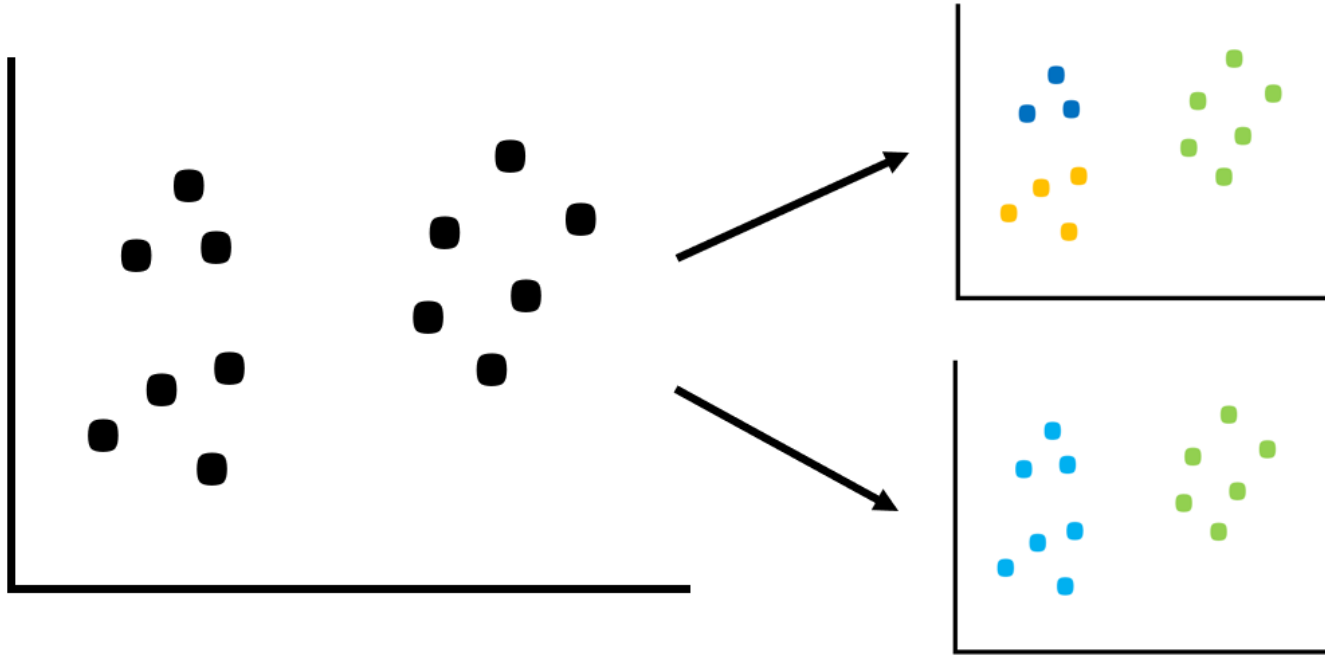






Problem

How to choose the right number of Clusters

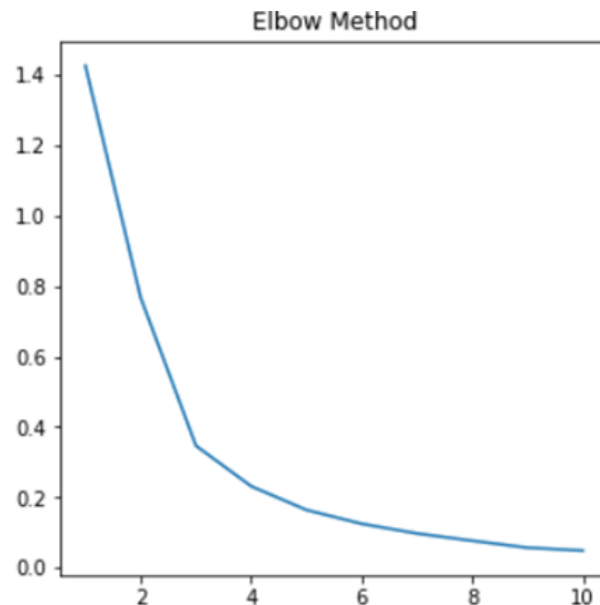


Elbow method

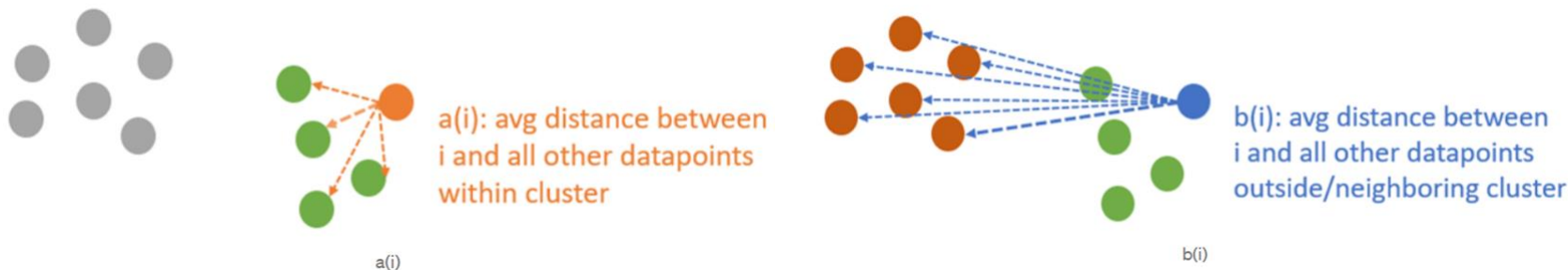
The sum of the squares of the distances of the points in the centroid of the cluster.

The sum is the minimum -> optimal K

But... ambiguous



Silhouette method

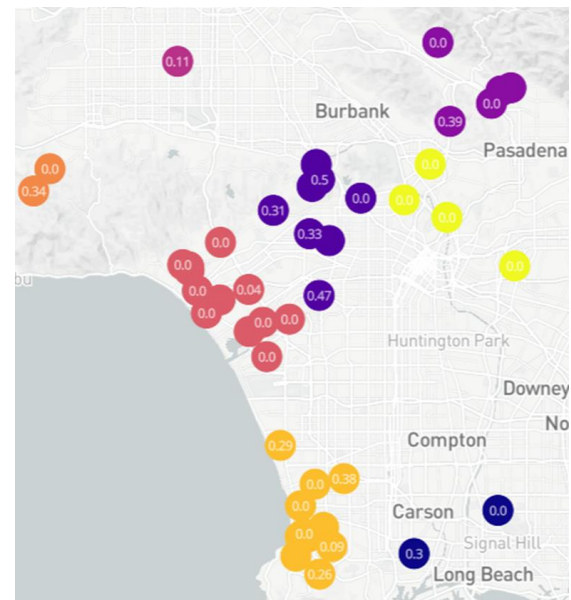
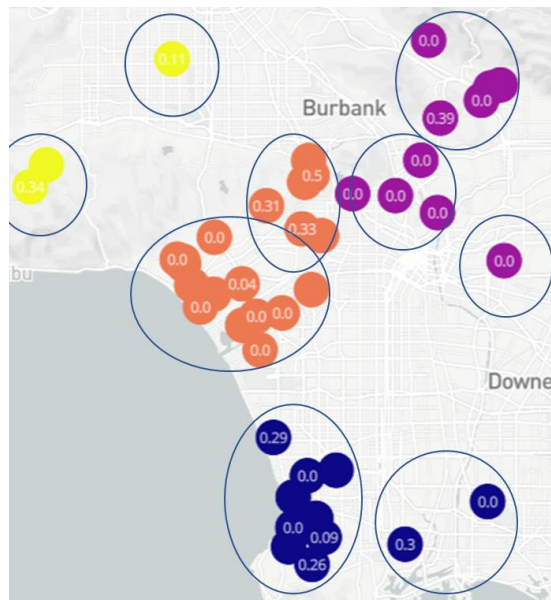
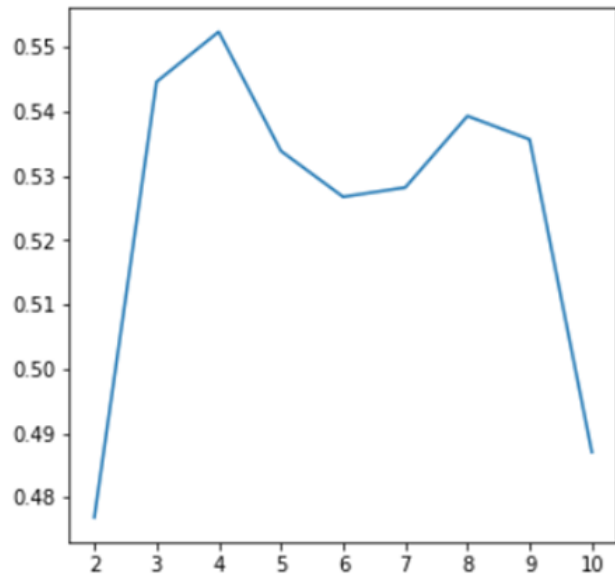


$$a(i) = \frac{1}{|C_I| - 1} \sum_{j \in C_I, j \neq i} d(i, j) \quad b(i) = \min_{J \neq I} \frac{1}{|C_J|} \sum_{j \in C_J} d(i, j) \quad S(i) = \frac{b(i) - a(i)}{\max \{a(i), b(i)\}}$$

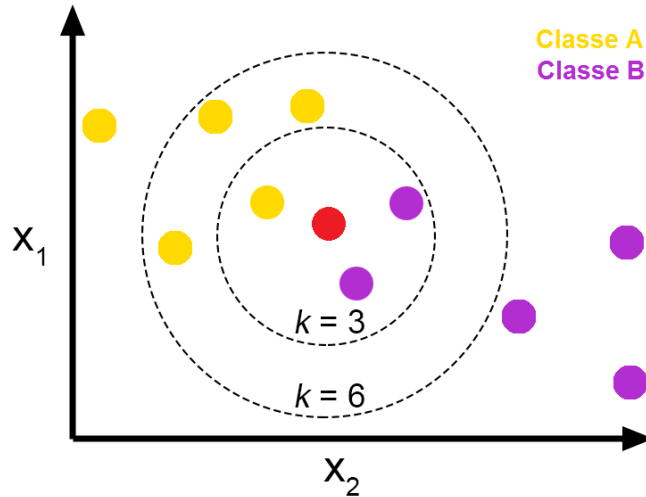
$-1 \leq S(i) \leq 1$

*CI = number of data in cluster
 $d(i, j)$ = distance from i to j

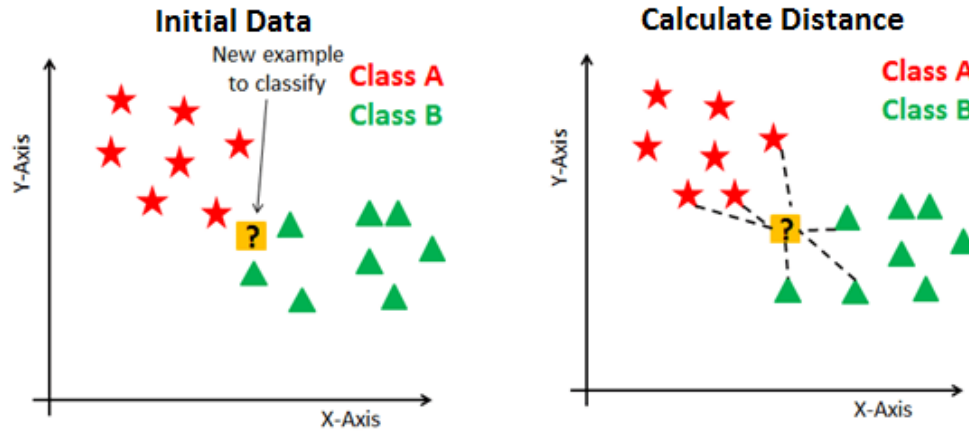
Silhouette Method



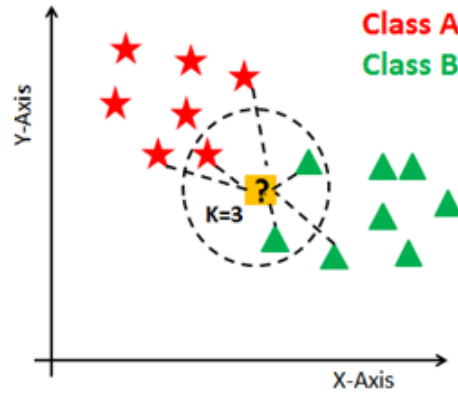
What algorithms is **K-NearestNeighbor**?



- ▷ Classification: Predict by majority vote
- ▷ Regression: Predict by mean value

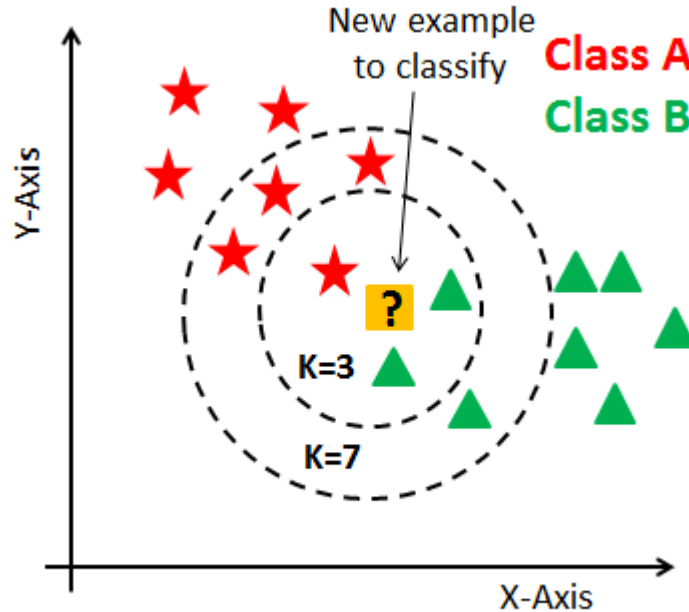


Finding Neighbors & Voting for Labels

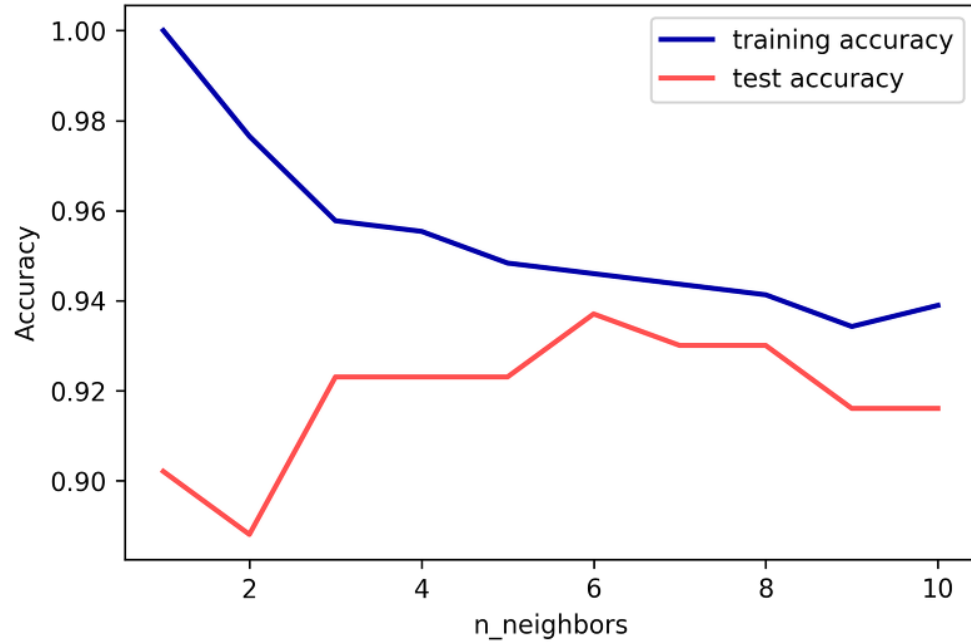


Problem

How to choose the optimal K

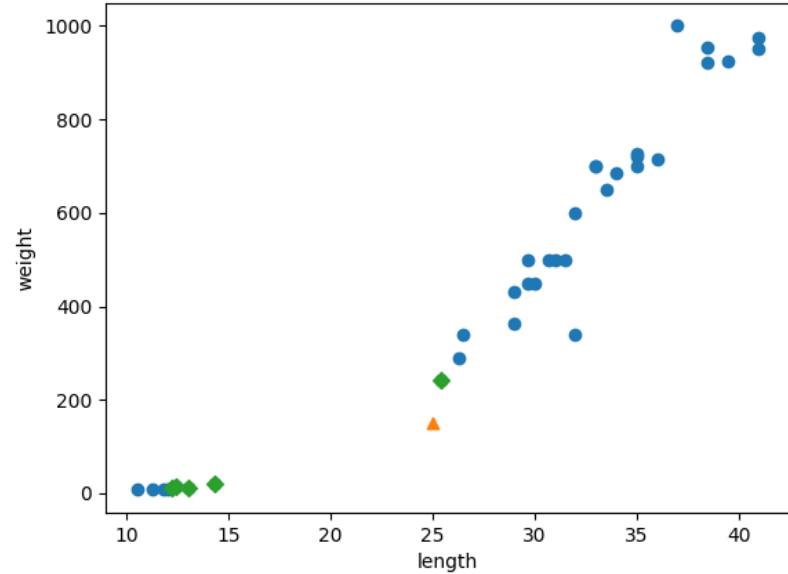
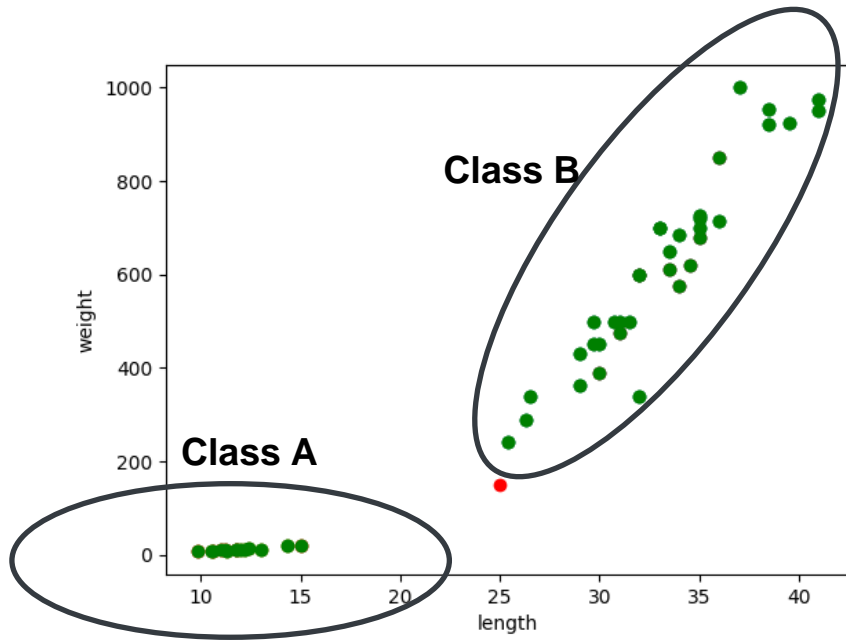


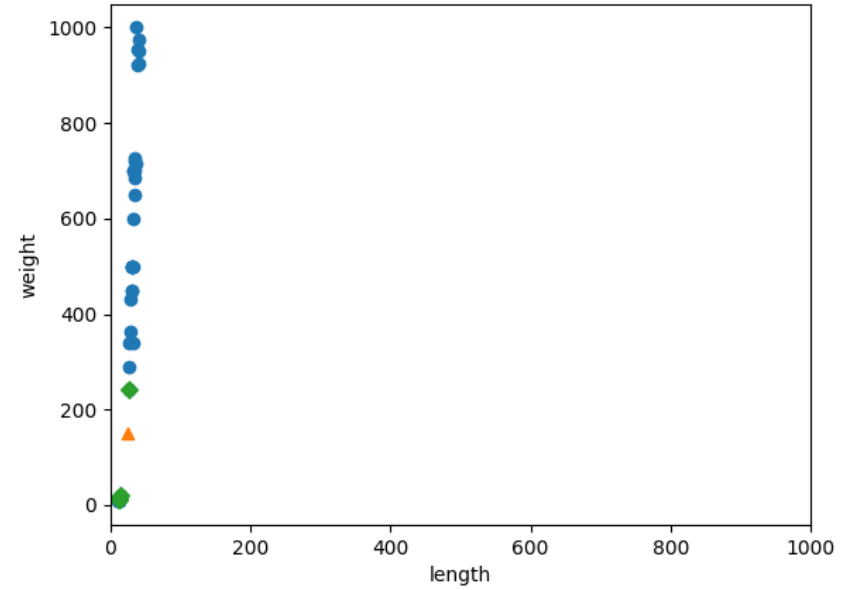
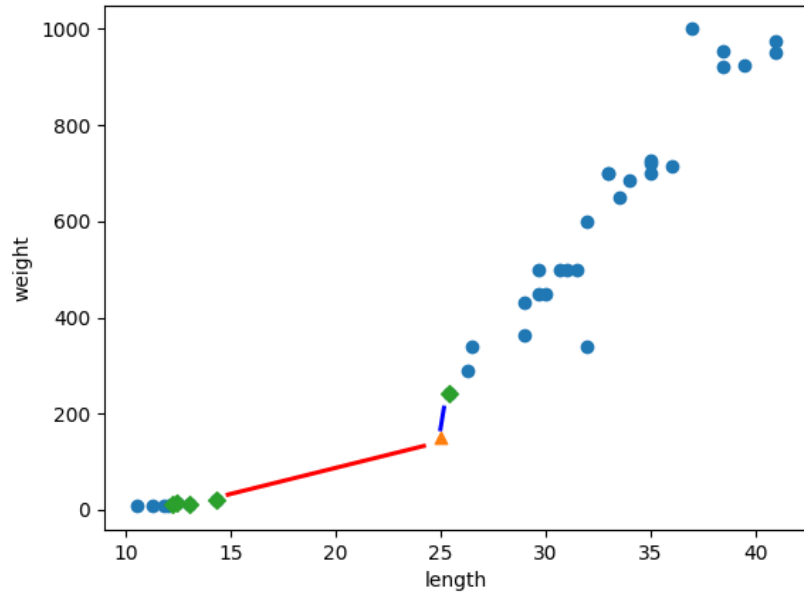
Find optimal K : Parameter tuning

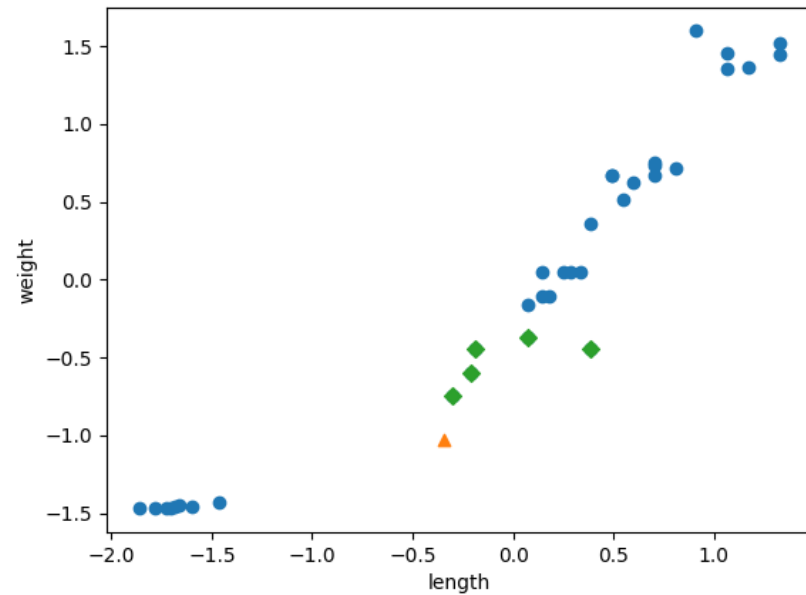
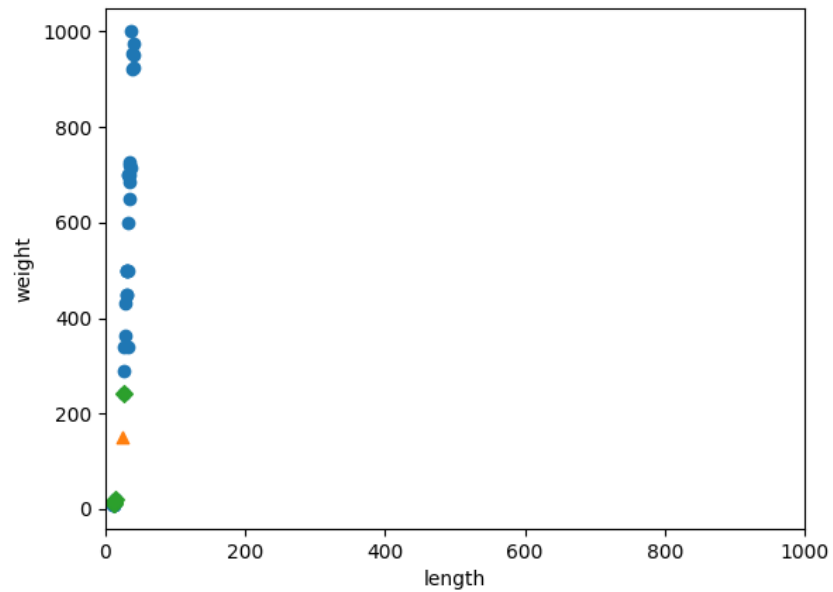


Problem2

(K=5)







Thanks!

Any questions?