

# K-Means

02 March 2022 20:51

## Unsupervised Learning (No Target Variable)

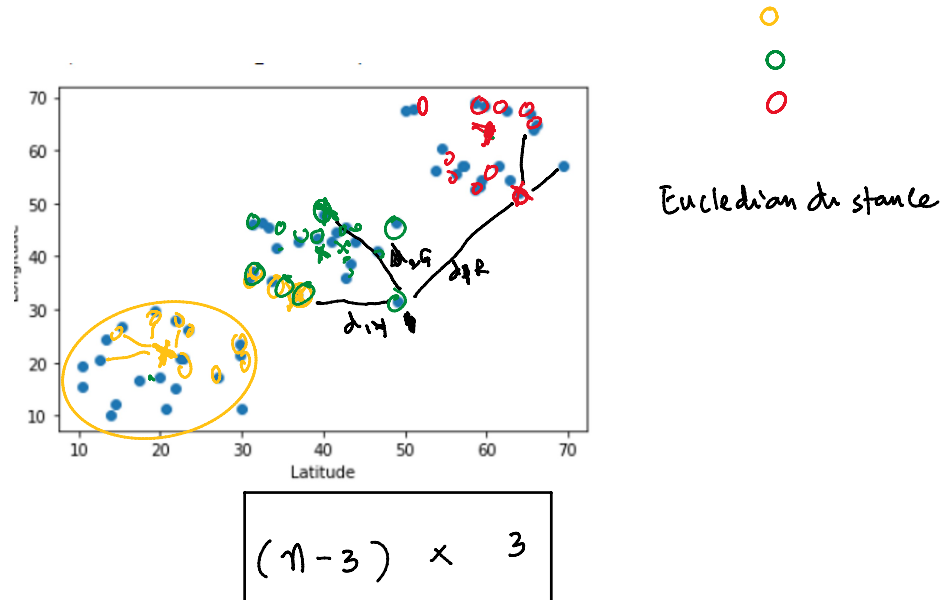
### → Clustering

group observations such that within same group observations are similar (Homogenous)  
Across groups different (Heterogenous)

### → K-Means ✓

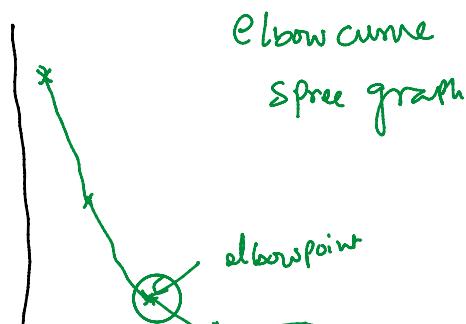
K - no. of clusters "K=5"

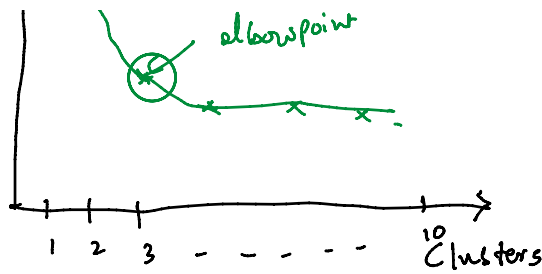
K



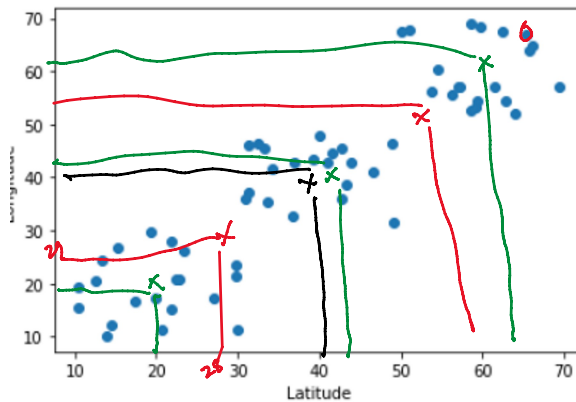
### Elbow Method

3 metrics  
WSS





$$S_1 = 1000$$



$$x_1 = (10-40)^2 + (11-40)^2 + \dots$$

$$x_2 = (10-40)^2 + (11-40)^2 + \dots$$

$$S_2 \quad \left. \begin{array}{l} x_1 \quad (10-29)^2 + \dots \\ x_2 \quad (11-28)^2 + \dots \end{array} \right| \begin{array}{l} x_1 \quad (70-55)^2 + \dots \\ x_2 \quad (70-60)^2 + \dots \end{array} = 700$$

600                      400

$$S_3$$

$$150$$

$$200$$

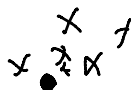
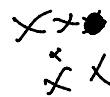
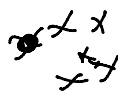
$$300$$

$$S_3 = 650$$

$$S_4 = 620$$

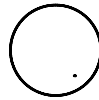
Initialization = (Kmeans++)

It chooses cluster  
Center one by one





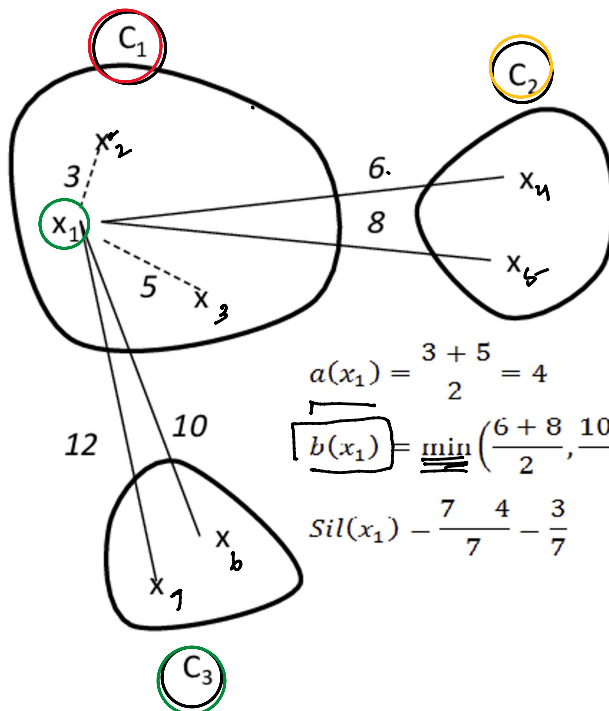
Silhouette Score  $(-1, +1)$  Clusters



a - within-cluster distance.

b - out of cluster distance.

Example of calculation of Silhouette score:



$$a(x_1) = \frac{3 + 5}{2} = 4$$

$$b(x_1) = \min\left(\frac{6 + 8}{2}, \frac{10 + 12}{2}\right) = 7$$

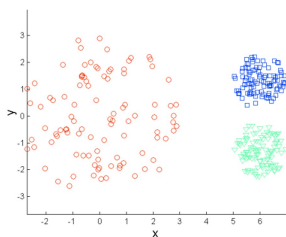
$$Sil(x_1) = \frac{7 - 4}{7} = \frac{3}{7}$$

$$a(x_1) \rightarrow \{c_1\}$$

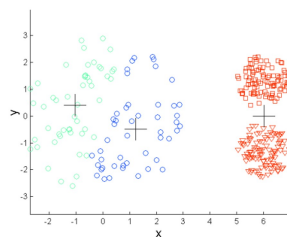
$$a = \frac{3 + 5}{2} = 4.$$

$$\min\left(\frac{6 + 8}{2}, \frac{10 + 12}{2}\right) = 7$$

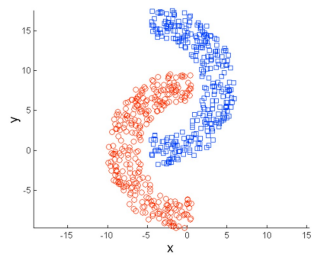
$$Sil = \frac{(b - a)}{\max(a, b)} = \frac{7 - 4}{7} = \left(\frac{3}{7}\right)$$



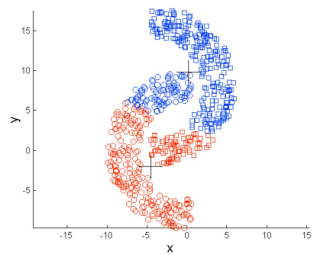
Original Points



K-means (3 Clusters)



Original Points



K-means (2 Clusters)