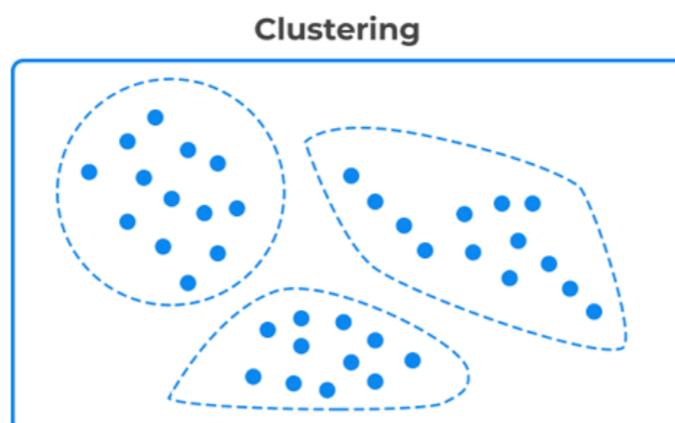
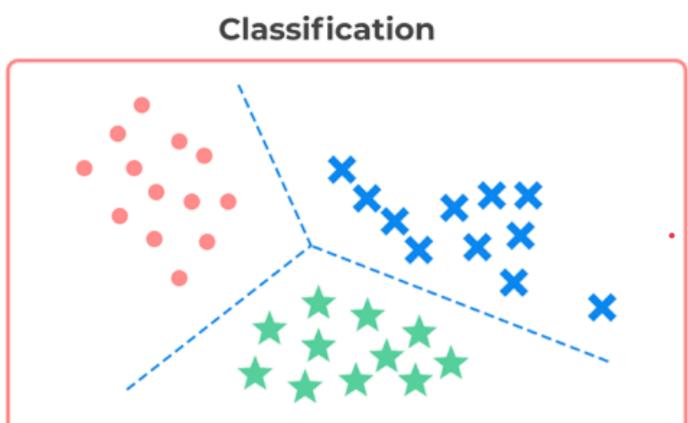


Ibrahim Radwan

20/09/2021

# Clustering

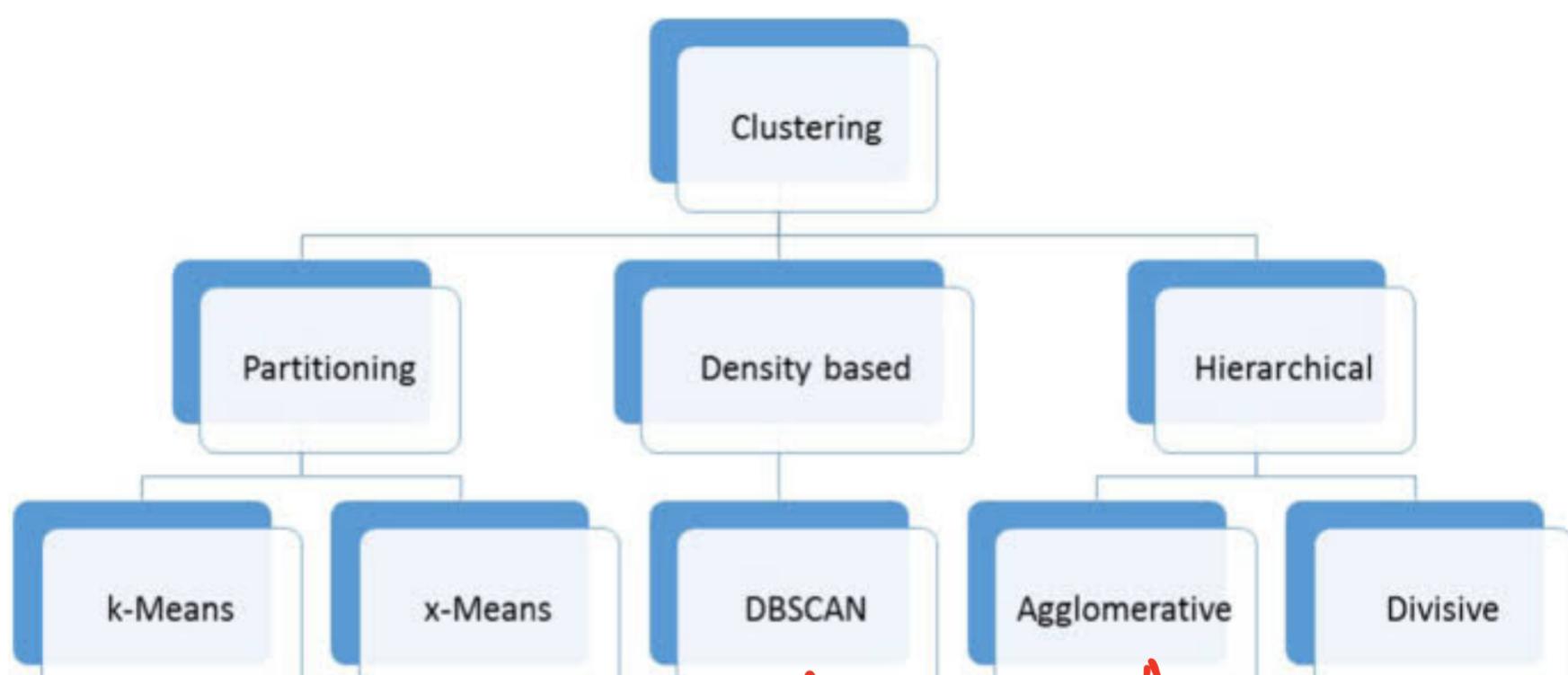
→ Grouping or Partitioning data based on their attributes



Label-based Attribute-based

## \* Types of clustering

$x \ y$   
 $s_1 =$   
 $s_2 =$   
 $s_3 =$   
 $\vdots$   
 $s_6 =$



\* Common ones

## → K-Means

- An iterative clustering algorithm
- Initialize: Pick  $K$  random points as cluster centers
- Alternate:
  1. Assign data points to closest cluster center
  2. Change the cluster center to the average of its assigned points
- Stop when no points' assignments change

General steps



$K = 2$

Number of clusters  $\rightarrow$  K-Means  
 Centroids of clusters  $\rightarrow$

# K-Means Mathematically

## K-means Solution

A proposal: Minimise the *distortion function*, i.e., the sum of the squared distances of each data point to its closest vector  $\mu_k$ .

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

↳ objective function

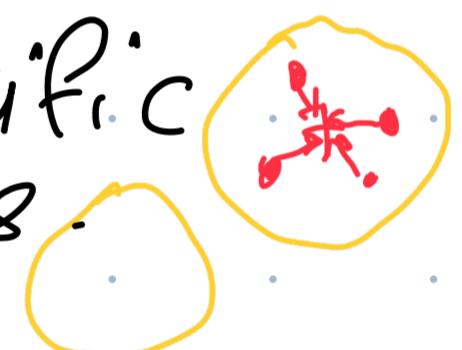
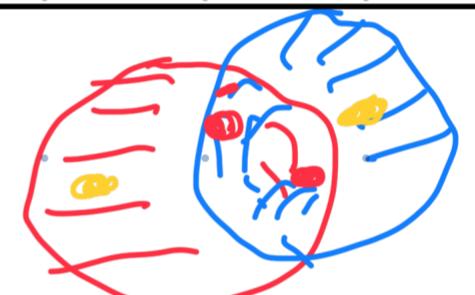
Equivalent to Expectation step.

Equivalent to Maximization Step.

- Given  $K$ , randomly select  $\mu_{k=1, \dots, K}$
- Minimise  $J$  with respect to  $r_{nk}$ , keeping the  $\mu_k$  fixed.
- Minimise  $J$  with respect to  $\mu_k$ , keeping the  $r_{nk}$  fixed.
- Repeat steps 2 (Expectation) and 3 (Maximisation) steps until convergence, that is,  $\Delta J < \epsilon$ .

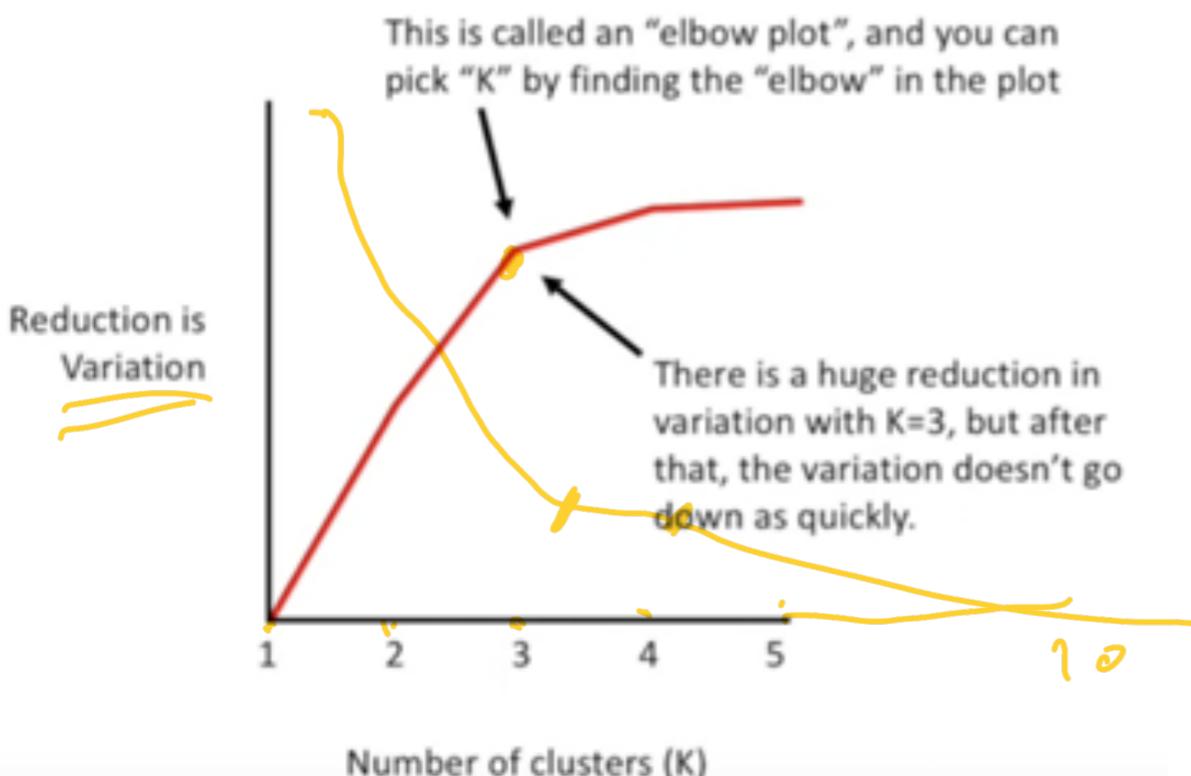
## Initialization of K-Means

- Randomly
- Manually by selecting specific positions to the centroids
- Can be done smartly  $\Rightarrow$  KMeans++



## How to find the best value for $K$ ??

### Elbow Method



Try different values of  $K$  & choose the point that shows significant reduction in variation

- Calculate WCSS for different  $K$ s
- Plot
- Identify

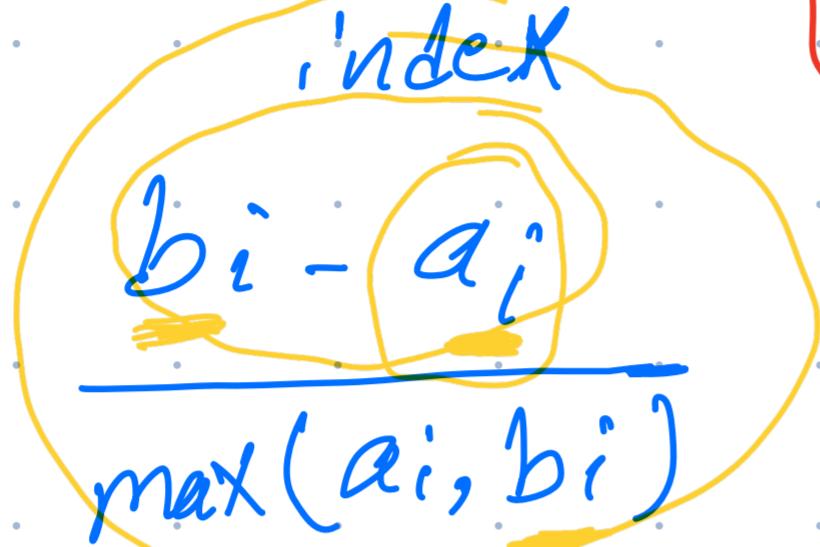
# How to Evaluate the clustering outcome?

## Silhouette Score

$SSI \Leftrightarrow \text{Silhouette score}$

$$SSI_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

index



The main idea is to assign each Point with a Score, that measures the Quality of being belonging to a specific cluster and far enough from other clusters

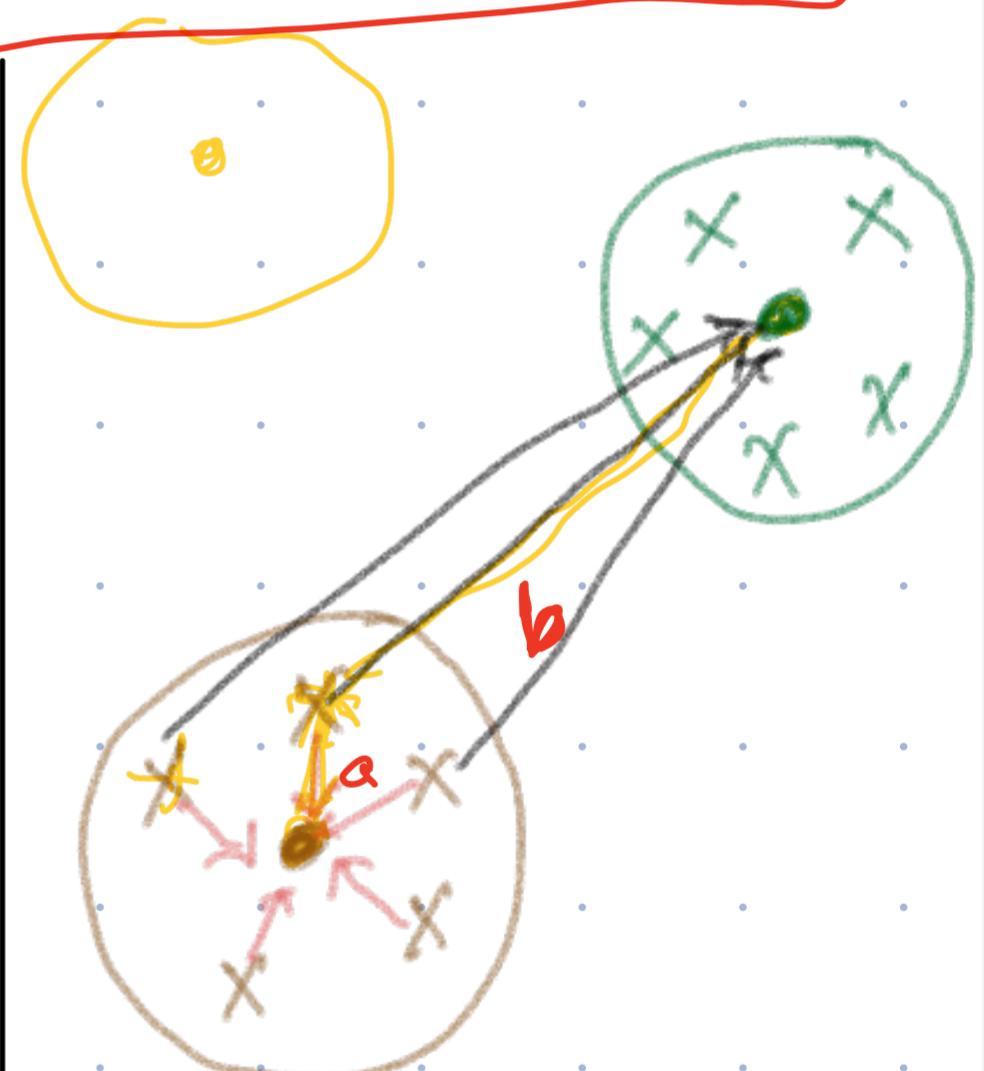
$a_i$ : distance between point  $i$  and the corresponding centroid.

$b_i$ : distance between point  $i$  and the closest centroid that this point does not belong to.

$SSI_i$ : can be between  
-1 as 1

very bad

very good



Silhouette Score  
can be used to determine the optimal number of  $K$ , how??