

"k-means and k-means++"

Supervised \rightarrow Label \rightarrow Target / Dependent variable

Unsupervised \rightarrow No label

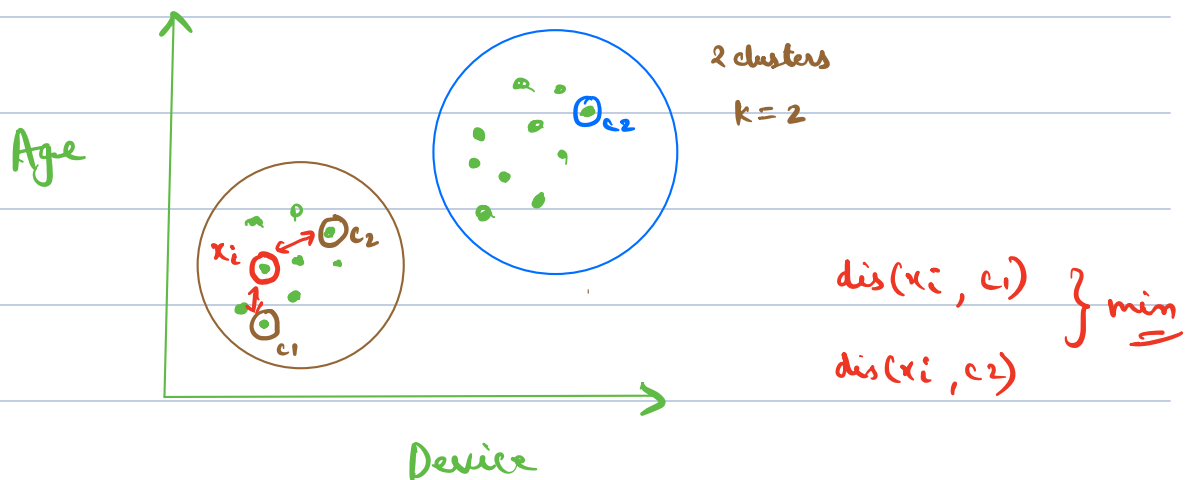
\hookrightarrow "CLUSTERING"

"JIONHOTSTAR"

Person	Age	Number of views	Duration (mins)	Device
A	20	4	180	Mobile
\rightarrow B	48	10	360	TV
C	30	1	20	iPad
D	25	15	10	Mobile

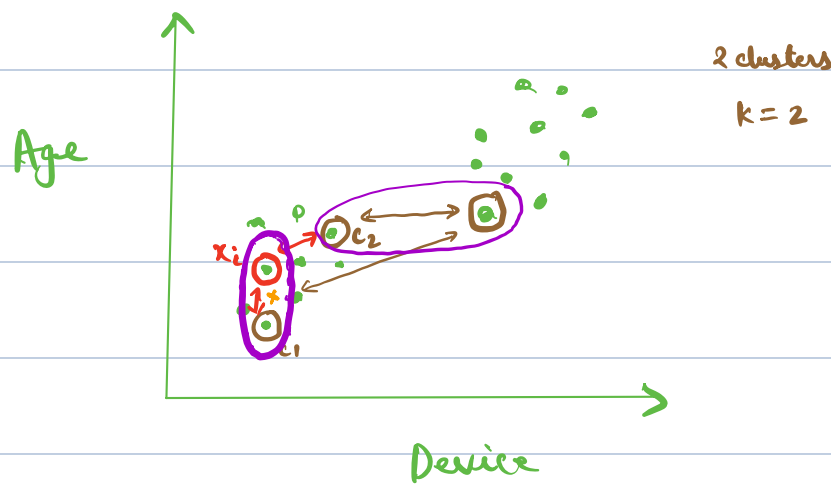
2 categories / clusters

- ① Young viewers and mobile users
- ② Senior adults and TV



Steps of k-means:

- ① Initialize 'k' centroids.
- ② Calculate distance of each point from all centroids.
- ③ Assign that data point to nearest centroid.
- ④ Re-calculate the new centroid.
- ⑤ Repeat "Steps 2 to 4"
- ⑥ Till the point of convergence [same centroids are getting calculated]



• EVALUATION METRICS :

① Silhouette Score

② Elbow Curve

↳ WCSS

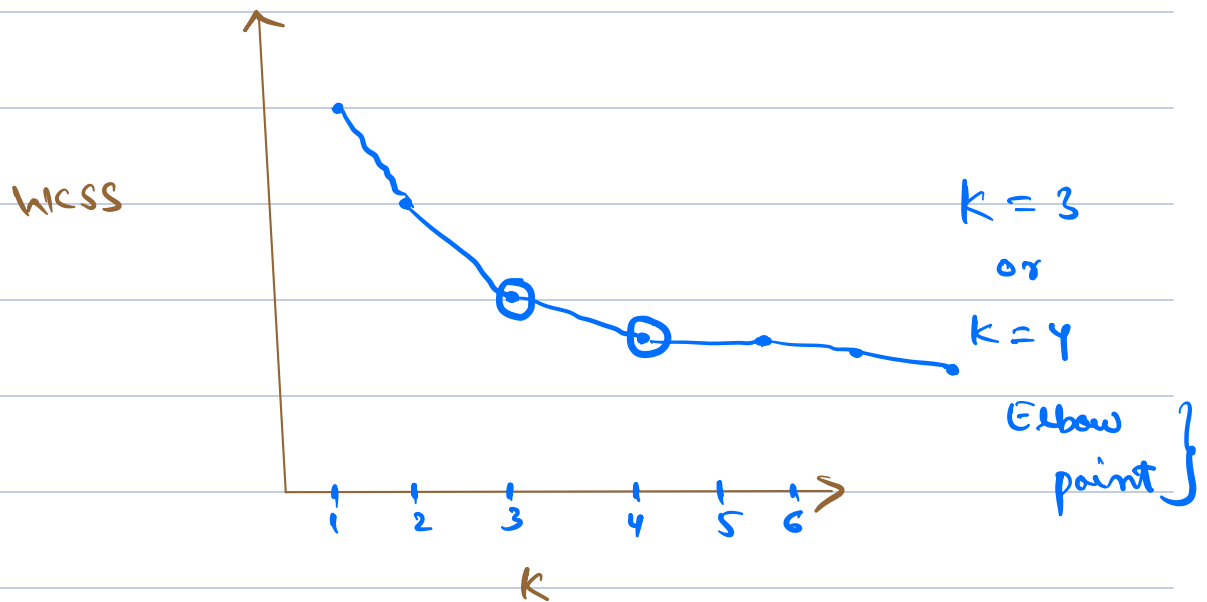
"Within Clusters Sum of Squares"

3 clusters

↳ For each cluster → intra-cluster distances (sum)

↳ sum of all Summed Up Distances of all clusters.

Python \rightarrow inertia_



k = [1 to 10]

WCSS = [Model1, Model2, ..., Model10]

To decide the value of 'k':

① Silhouette Score \rightarrow [-1, +1]

② Elbow Curve "Best Model"

③ Domain Knowledge \hookrightarrow SS as close to +1

Limitations of kmeans:

① Not robust to outliers

② Scaling of data.

③ Because of random centroids initialization, may not get correct clusters.

④ Mostly gives spherical clusters

"K-means++"



Another case of "Kmeans" where:

- instead of initializing random data points as centroids, we take farthest points as centroids.

Categories:

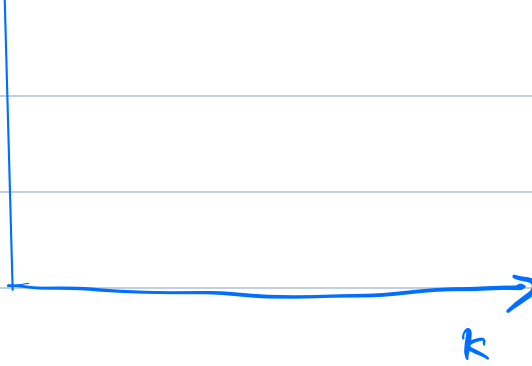
- ① Low spenders, medium, high
- ② Frequent visitors v/s non-frequent
- ③ High discount users v/s low
- ④ New users v/s old users

$$K = [1, 2, 3, \dots, 9] \rightarrow SS = [SS_1, SS_2, \dots, SS_9]$$

$$\rightarrow K\text{-models} = [\text{Model}_1, \text{Model}_2, \dots, \text{Model}_9]$$

$$\rightarrow \text{inertia} = [WSS_1, WSS_2, \dots, WSS_9]$$

WSS ↑



① Visualization $\rightarrow k=3$

② Elbow Curve $\rightarrow k=3$ or 4

③ Silhouette Score $\rightarrow k=2$ or $k=3$ or $k=4$



+ Domain Knowledge
" " $k=3$ "