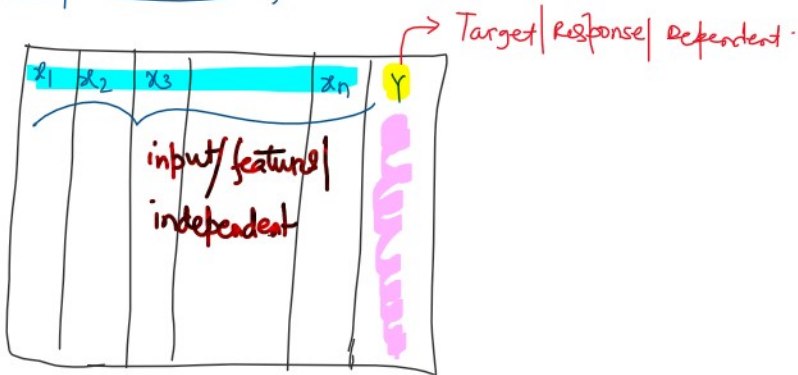


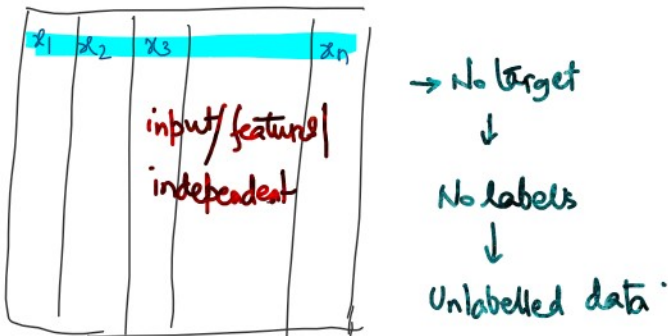
# K-Means Clustering

07 April 2024 19:59

## Supervised learning



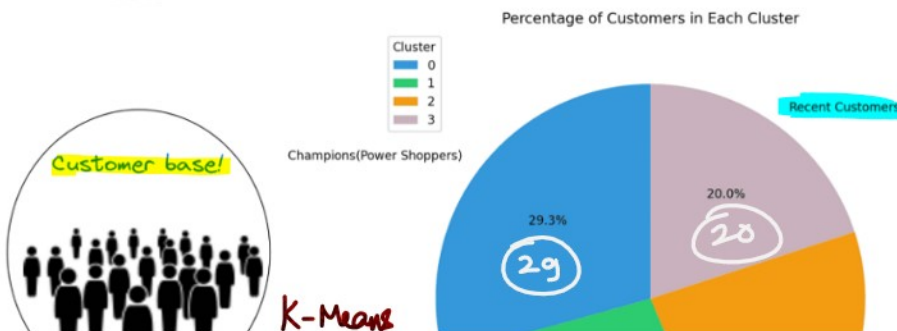
## Unsupervised learning

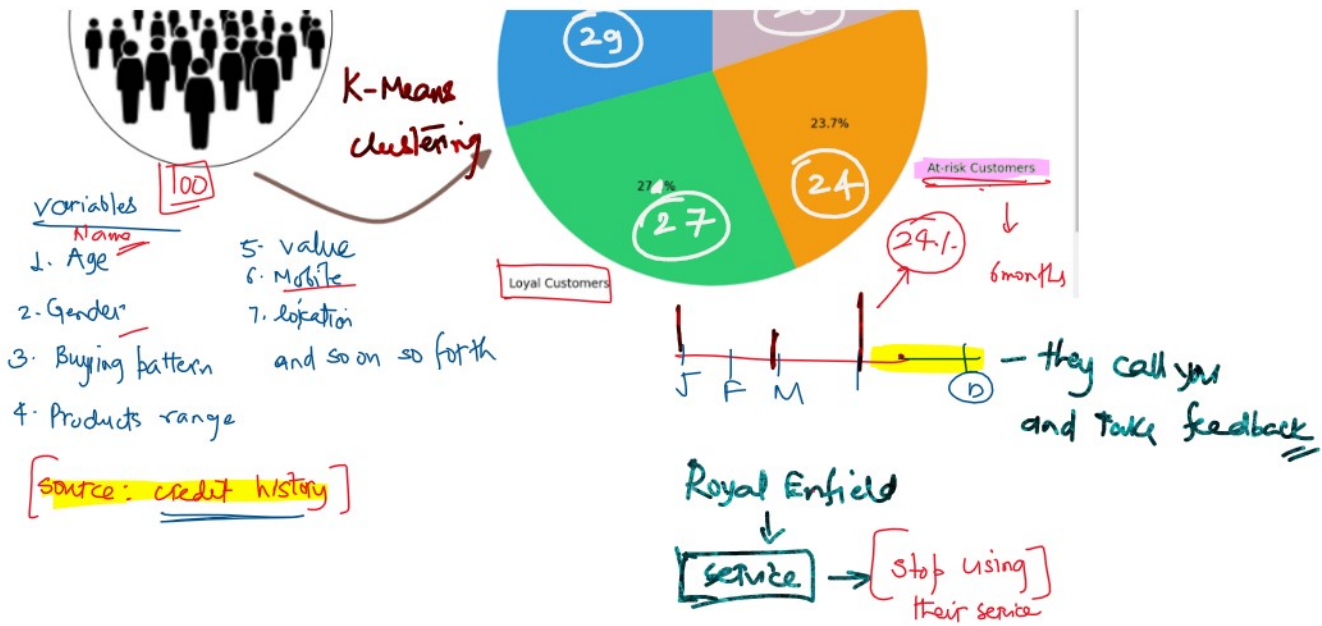


In unsupervised learning, there is **no output variable** to **guide the learning process**, and **data is explored** by algorithms to find the patterns.   
relationships between  $x_i$ s (features only)

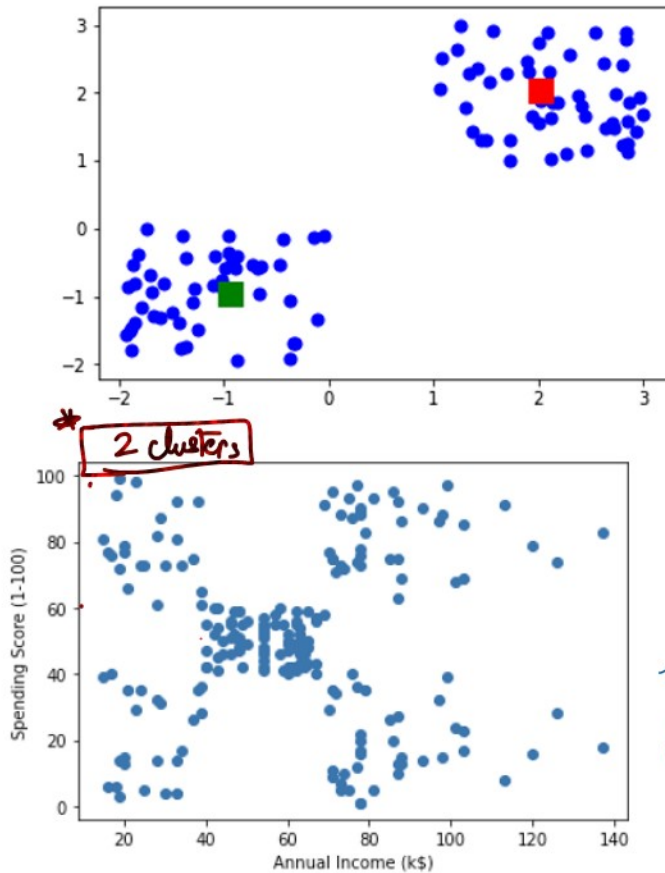
Given that the data has no labels, the algorithm identifies similarities on data points and groups them into clusters/segments.

## Customer Segmentation





## \* K-Means Clustering

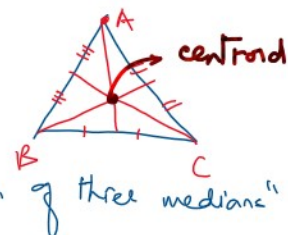


is it easy ??

No

Q: What is centroid?

Geometry

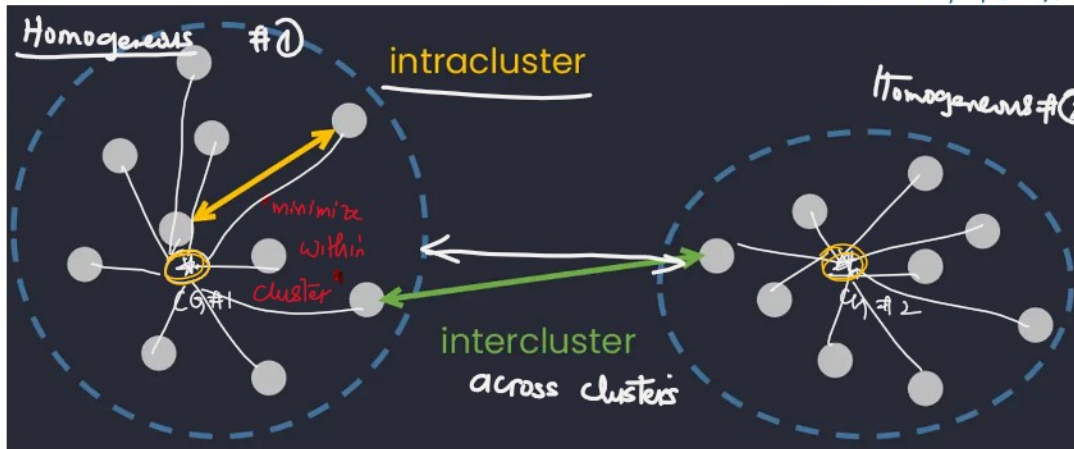
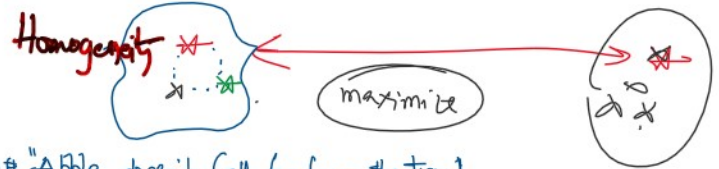


In K-Means clustering technique, each cluster is represented by its center (called a centroid), which corresponds to arithmetic mean of the

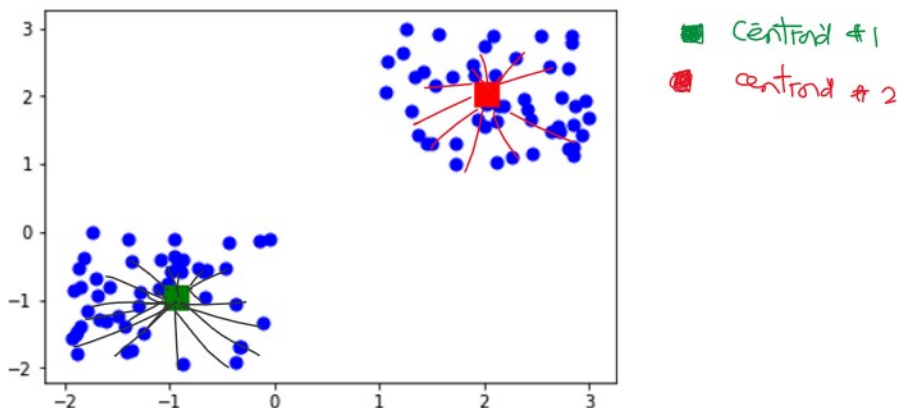
Intersection of three medians (called a centroid), which corresponds to arithmetic mean of the data points assigned to the cluster.

"Intersection of three medians"

A centroid is a single data point within the cluster which represents the center and might not be necessarily be a member of the dataset

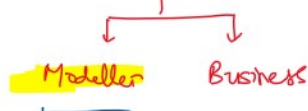


K-Means clustering algorithm tries to minimize the distance of the points in a cluster with their respective centroids.



## K-Means Algorithm steps

1. Specify the no. of clusters 'K' - user input



"Scree Plot" aka



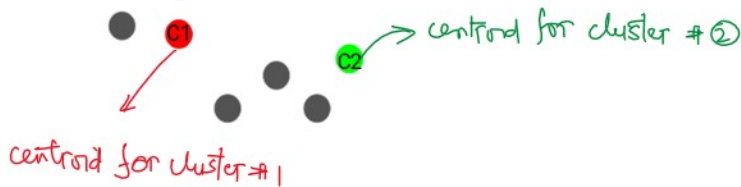
2. Initialize centroids by first shuffling the dataset and then randomly select ' $K$ ' data points for the centroids without replacement.

100 000

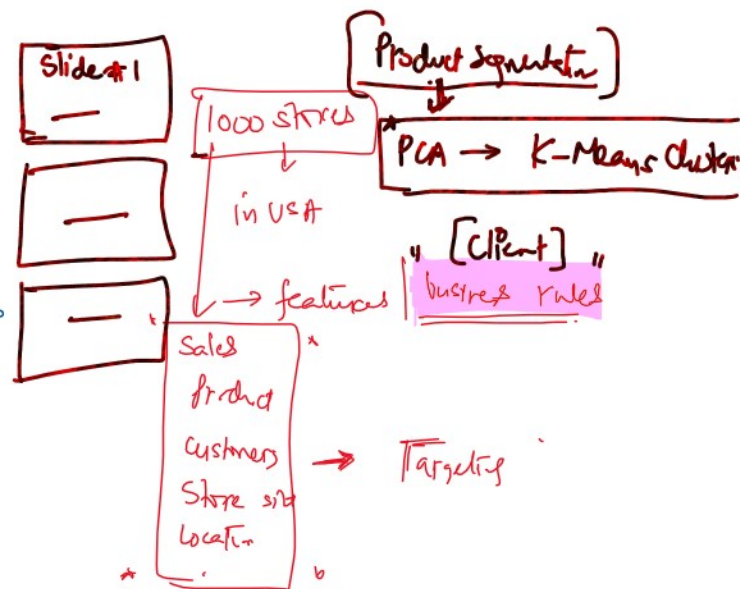
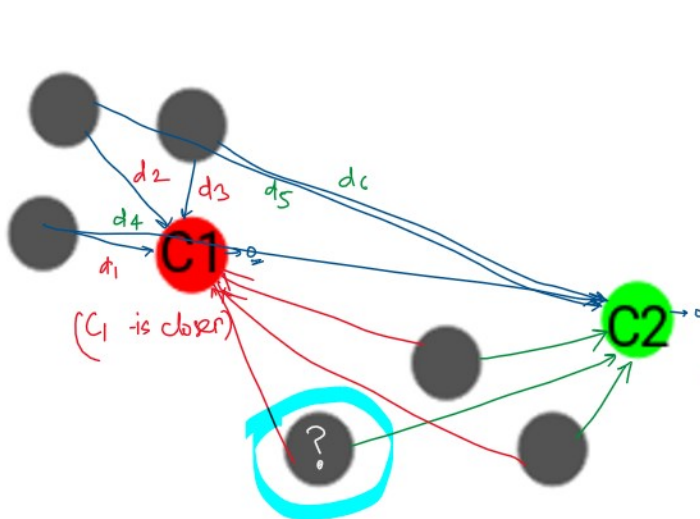
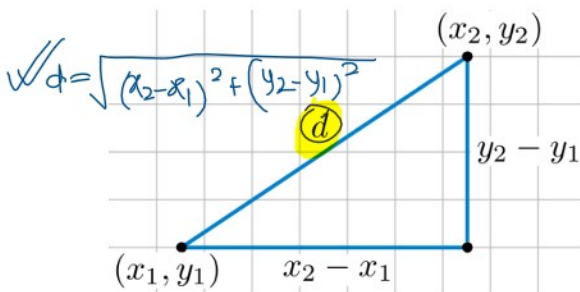
row ④

Total #8

$K=2$



Euclidean distance:

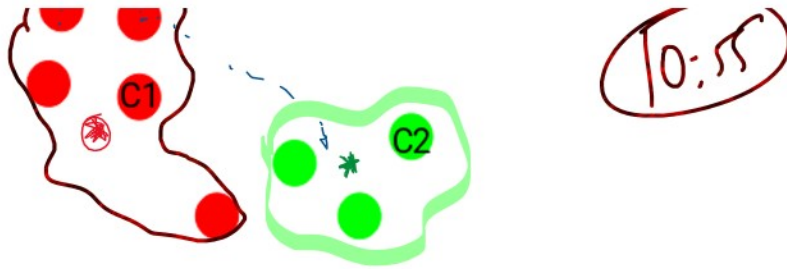


3. Assign all the points to the closest cluster centroid

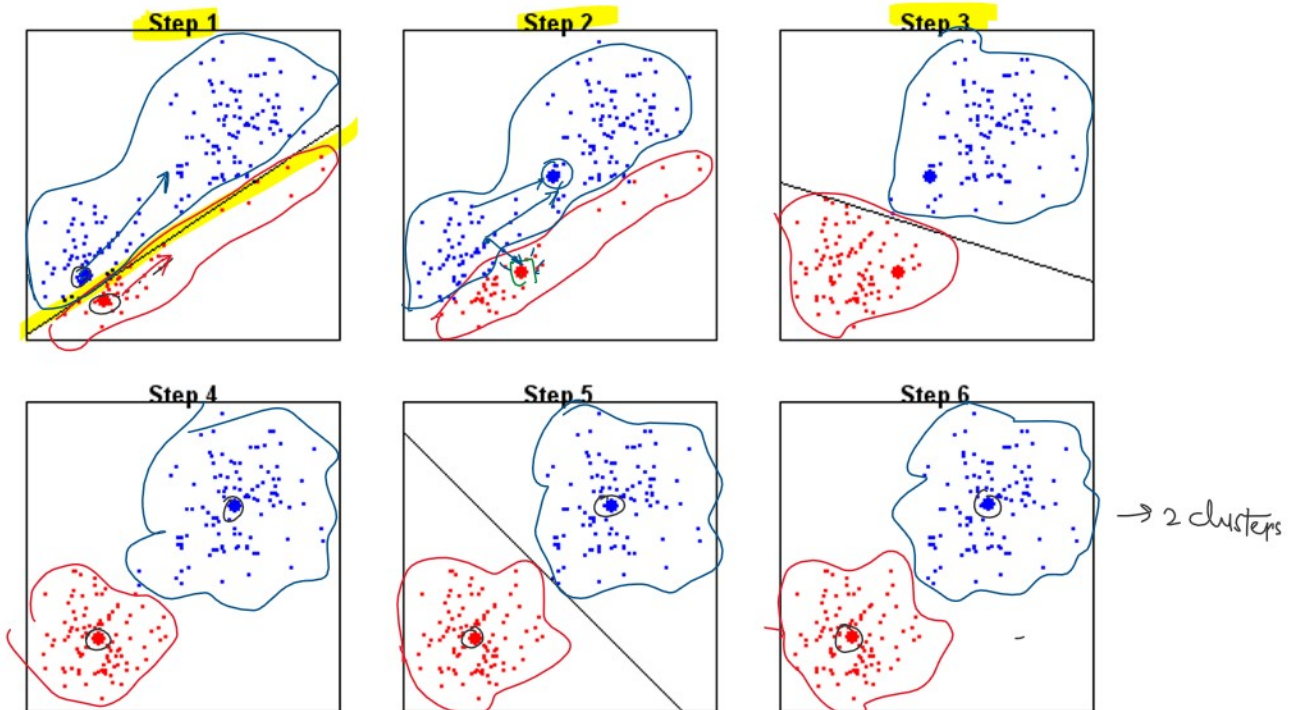
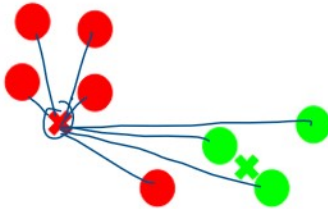


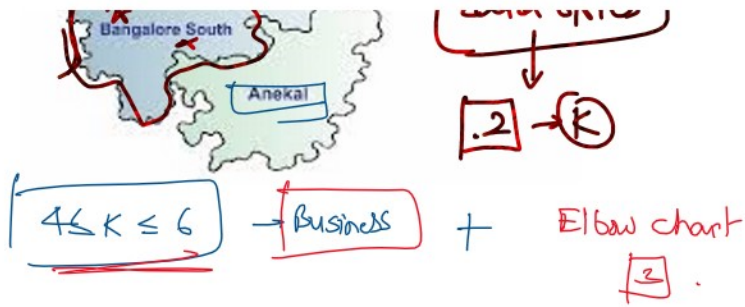
10:55



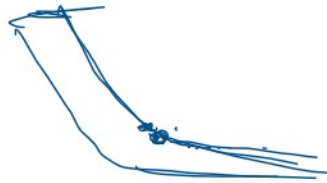
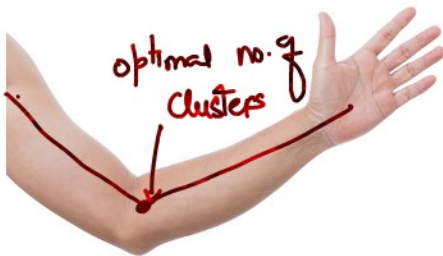


4. Recompute the centroids of newly formed clusters.



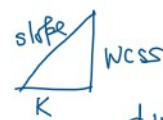
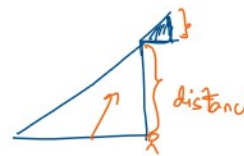
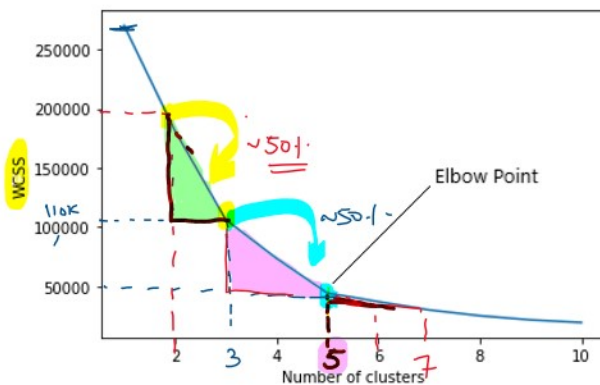


## # Elbow chart



\* Elbow method is a heuristic method to determine the optimal number of clusters ( $K$ ) in K-Means clustering

\* It involves plotting the within-cluster sum of squared distance (WCSS) against the no. of clusters and identifying the "elbow" point where the rate of decrease in WCSS slows down.

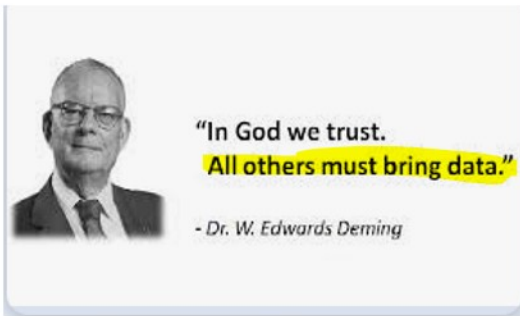


$$\frac{dWCSS}{dK} \downarrow \downarrow \downarrow$$

$$\frac{dy}{dx} = \frac{y}{x}$$

$K=5 \rightarrow$  no. of optimal clusters

Testing different no. of clusters and measuring the resulting WCSS and choosing the 'K' value at which an increase in 'K' (say  $K=5$  to  $K=6$ ) will cause a very small decrease in WCSS. one more cluster.



~~In 2024~~  
~~→ Gut feeling~~  
~~→ Business Acumen~~



Task # Hierarchical clustering

- Agglomerative
- Divisive.