

Unsupervised Learning  
 ↓  
 Data with input cols only

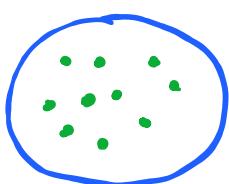
Amazon Customer's Data

Name	Age	City	Gender	Purchases	Amount spent
0	-	-	-	-	-
1	-	-	-	-	-
2	-	-	-	-	-
3	-	-	-	-	-
4	-	-	-	-	-
5	-	-	-	-	-

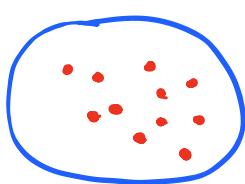
↓  
 Unsupervised Algorithm

They will learn two things :

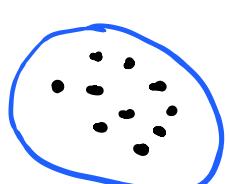
- How many groups/clusters exist in our data ?
- Which customer/datapoint belongs to which cluster.



Low  
spenders



Medium  
spenders



High  
spenders

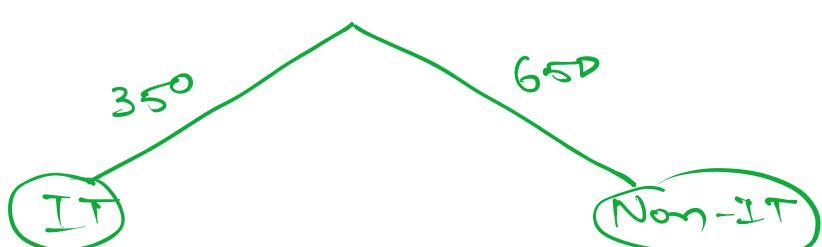
# Intellipaat (1000)

Name	Age	Gender	Edo qualif	Domain	Exp Level	cluster numbers
→ 0	-	-	-	IT / Non-IT	-	C <sub>1</sub>
→ 1	-	-	-	-	-	C <sub>4</sub>
→ 2	-	-	-	-	-	C <sub>2</sub>
→ 3	-	-	-	-	-	C <sub>1</sub>
→ 4	-	-	-	-	-	C <sub>3</sub>
→ 5	-	-	-	-	-	C <sub>2</sub>
→ 6	-	-	-	-	-	C <sub>1</sub>
→ 7	-	-	-	-	-	C <sub>1</sub>

Gender

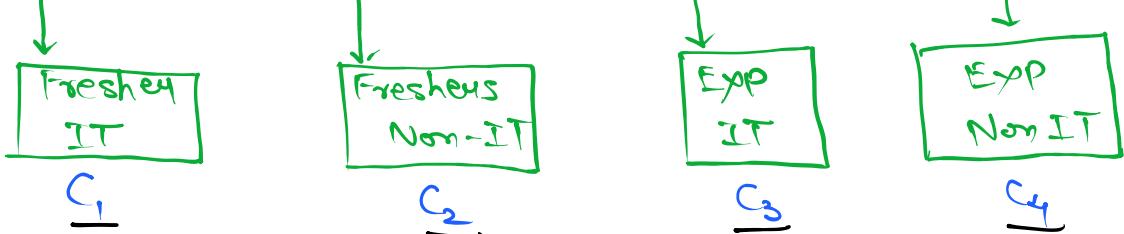


Domain



Domain + Exp Level

<1 years



- ① Targeted Marketing
- ② Recommendations
- ③ Offers, coupons, discounts.

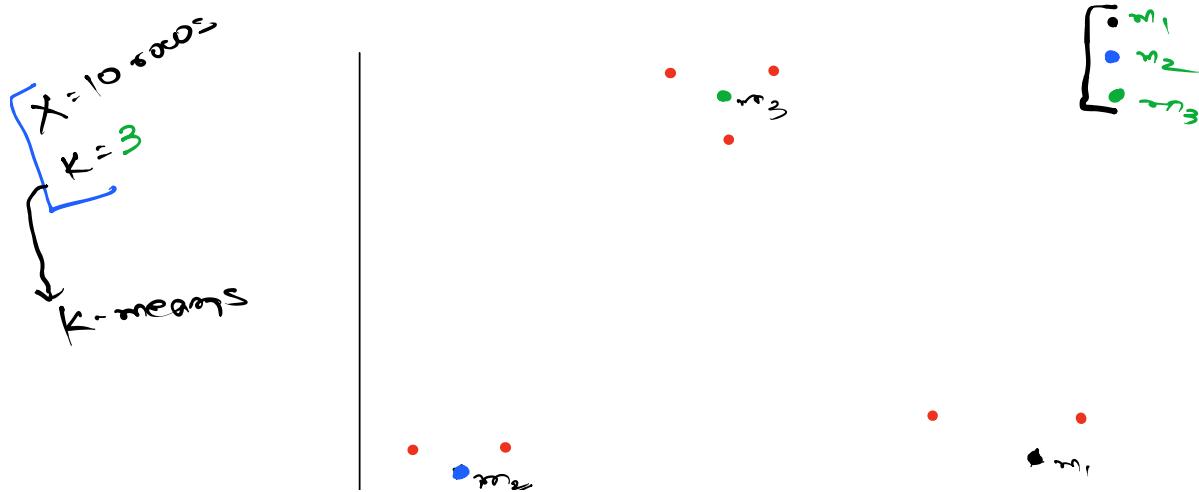
## Unsupervised Learning Algorithms →

- ① K-Means Clustering
- ② Hierarchical Clustering
- ③ DBScan Clustering
- ④ Gaussian Mixture Models (GMMs)

### K-Means Clustering

Two requirements of K-Means :

- ① Data (input cols)  $\rightarrow X$
- ② K-value : No of clusters to be identified.

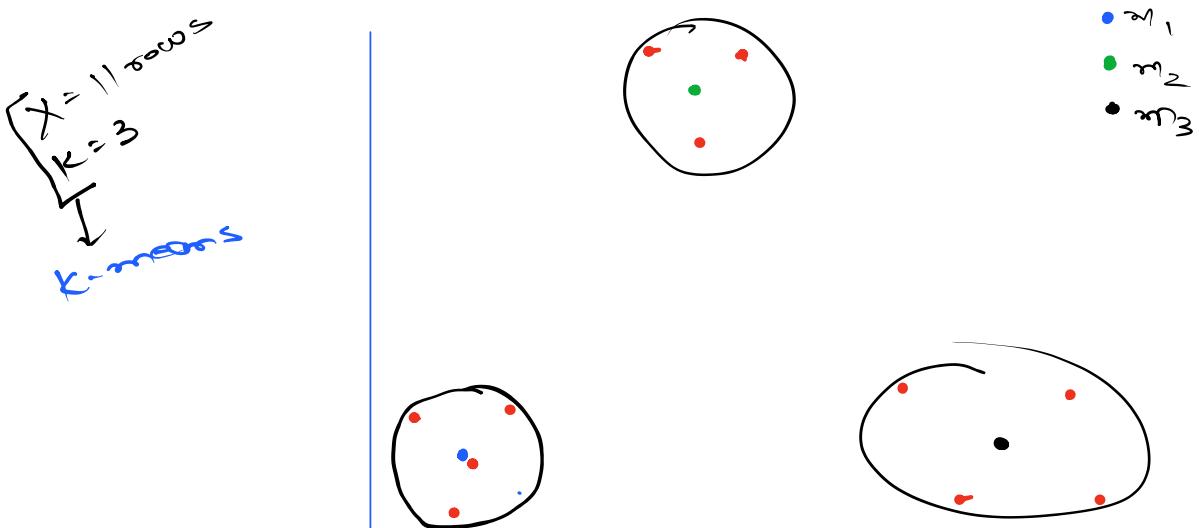




- ① Create centroids/cluster means on the data randomly.
- ② Calculate the distance b/w each datapoint to each centroid.
- ③ Put together the datapoints with their closest centroids in clusters.

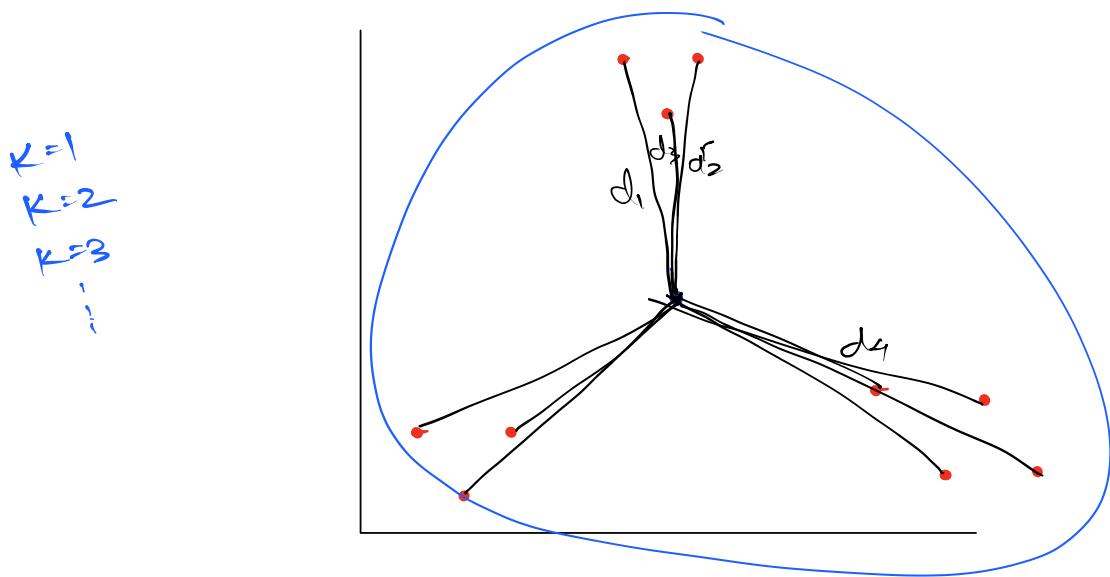
④

⑤



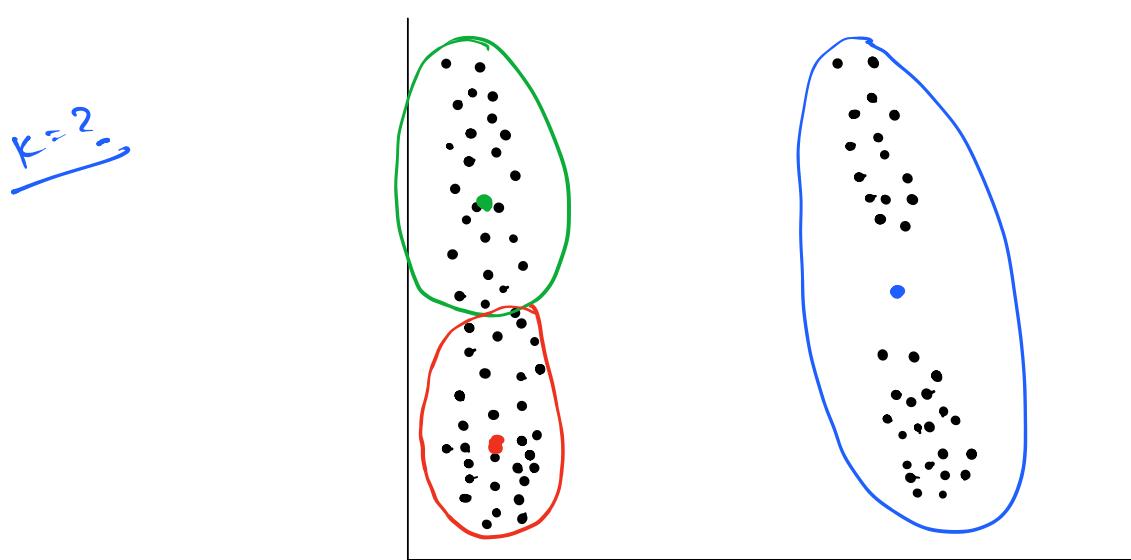
Finding optimal k-value using WCSS/elbow method:

WCSS: Within cluster Sum of squared distances.  
or  
SSD



$$\text{WCSS}(k=1) = d_1^2 + d_2^2 + d_3^2 + \dots + d_{10}^2$$
$$= 2023$$

Initialization Trap :



To overcome this trap, we now have updated version of K-means algorithm which is known as 'K-means++'.

| ...      . . .      ... |  
K-means++

