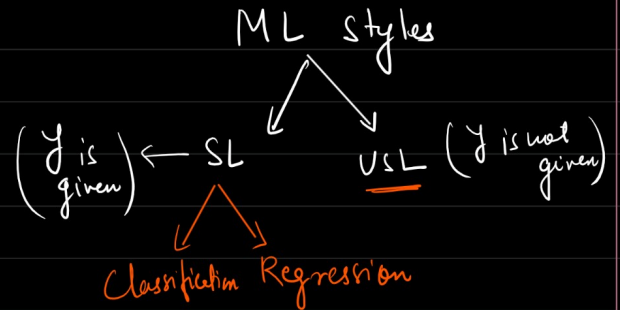


# Introduction to Unsupervised learning

\* Target variable is Not given (y)

\* USL groups / segments your data  
(Finds pattern in the data)



## Examples Zara store

Customer purchase amount	Customer salary in lakhs
--------------------------	--------------------------

2500

10

↑  
Cust Amount

SL PA ↓  
①

①  
S ↑ PA ↑

③  
SL PA ↓

②  
S ↑ PA ↓

Clothes mix up

Shop keeper → will try to touch the cloth and understand the pattern, different recks

To increase revenue

→ Group 2 customers should be focused more.

→ Group 1 → loyalty card & Extra discount.

\* depends on you how you want interpret the groups and bring business.

## Motivation

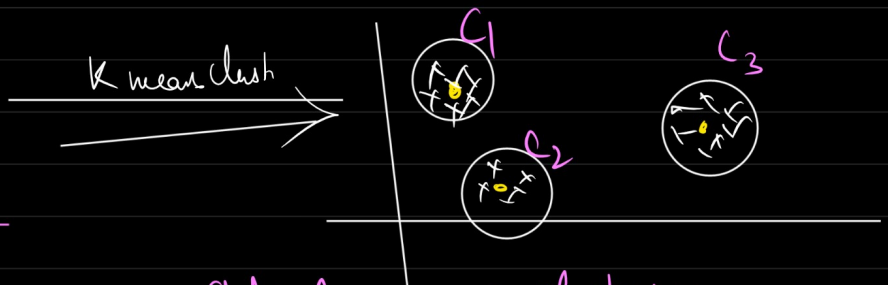
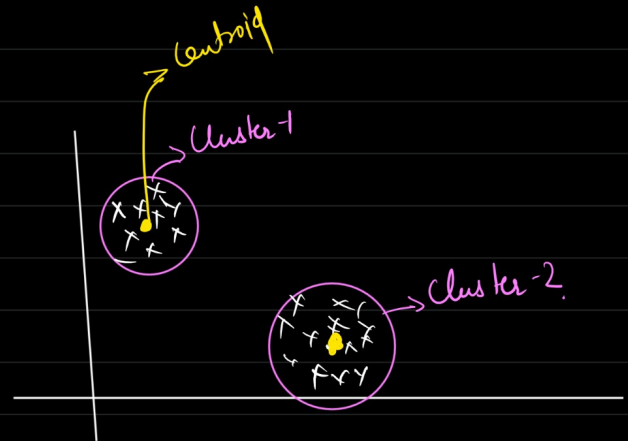
- Need to launch a campaign based on income / sales.
- Different ways of fraud / money laundering.
- Grouping of images / patterns of specific image
- Document / Article Analysis.
- Cohort Analysis.

\* In higher dimensions, you cannot identify groups / patterns manually. Therefore you need USL.

# USL

- ① K-means clustering
- ② Hierarchical clustering
- ③ DBSCAN

## \* K means clustering



## Steps of K means clustering

Step-1 initialize Centroid (K)

- ↳ Can be random dp
- ↳ Centroid of all the dp
- ↳ A random new point

Step-2 → Points nearer to the centroid will be labelled as that group

Step-3 Re-calculate the centroid  
Again repeat  $\hookrightarrow$  Avg.  
Step 2 untill Centroid doesn't change.



Exp	Salary
2	5
3	6
5	7
$\left( \frac{2+3+5}{3}, \frac{7+6+5}{3} \right)$	

\* Distance

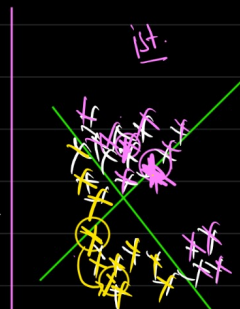
① Euclidean distance

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

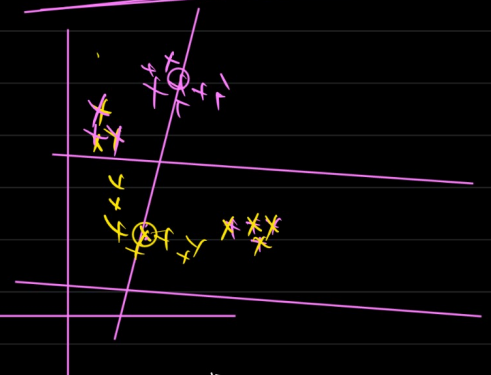
② Manhattan distance



$$\text{Manhattan distance} = |x_2 - x_1| + |y_2 - y_1|$$

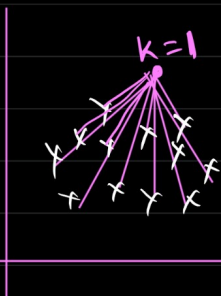


⇒ after recalculating Centroid.



# How to select to k value

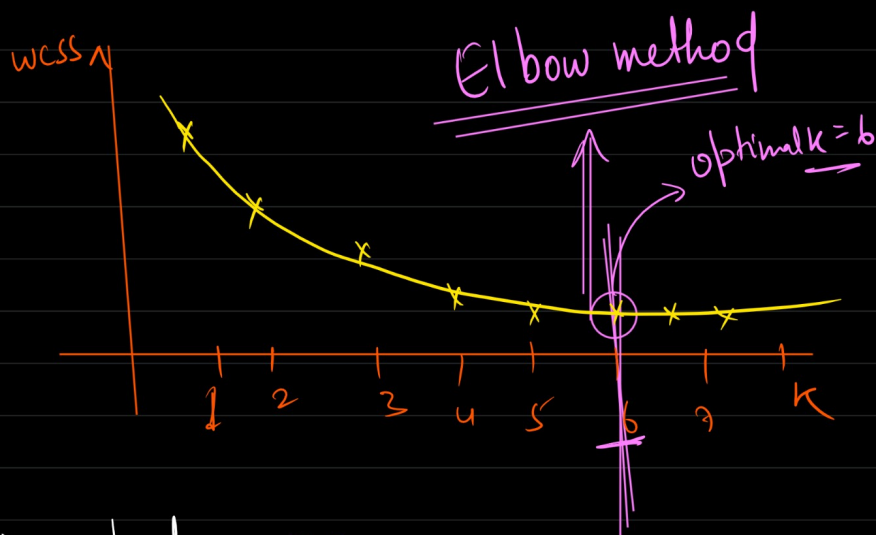
WCSS  $\rightarrow$  within cluster sum of squares distance



$$WCSS = \sum_{l=1}^n \left( \text{distance b/w points to nearest Centroid} \right)^2$$

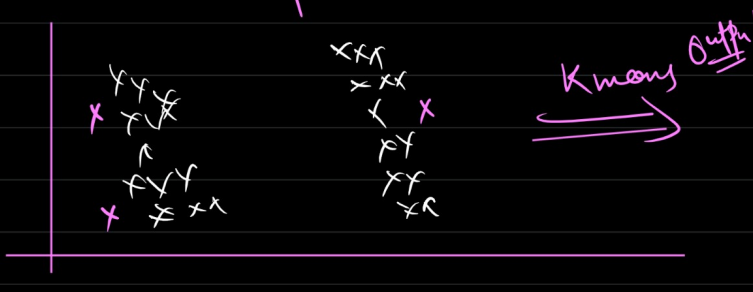
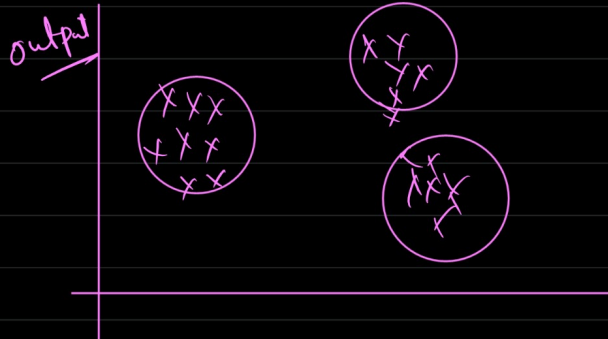
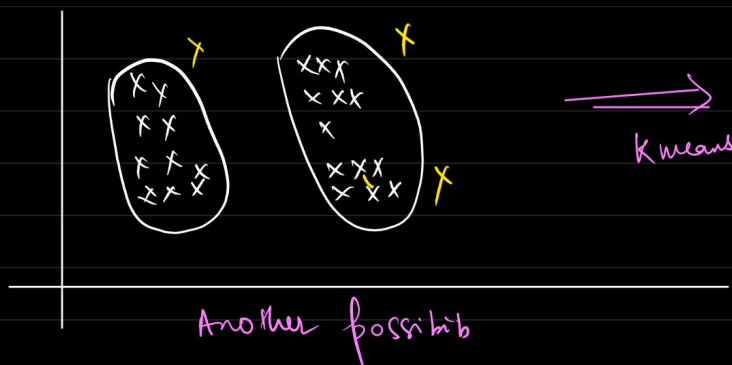


$\Rightarrow$  for  $k=2 \rightarrow$  As you increase  $k$ ,  
WCSS decreases  
 $\rightarrow$  And at optimal  $k$ ,  
WCSS stops changes.



## \* Random Initialization trap

$k=3$



\* The final cluster depends on the initialization of  $k$ .



$k$  means ++



initialize the  $k$ 's as much far as possible

