* PCA.

$\begin{cases} * \text{ USL} \\ * \text{ Anomaly detection} \\ * \text{ Time series.} \end{cases}$

USL → K means.
→ Hierarchical
→ DBSCAN

Anomaly detection → Isolation forest
→ DBSCAN
→ Local outlier fac..

Time series → Smoothing models
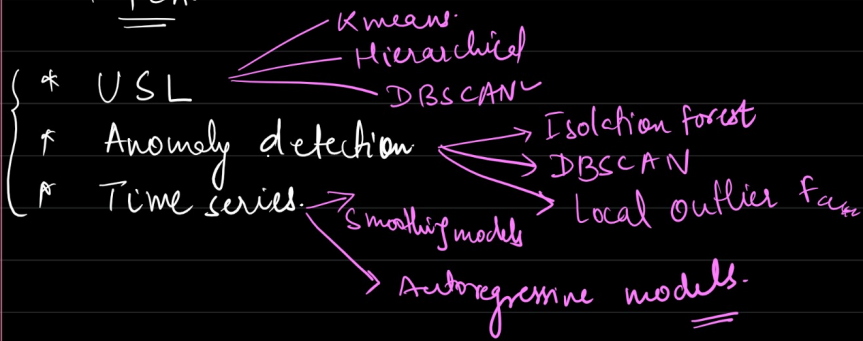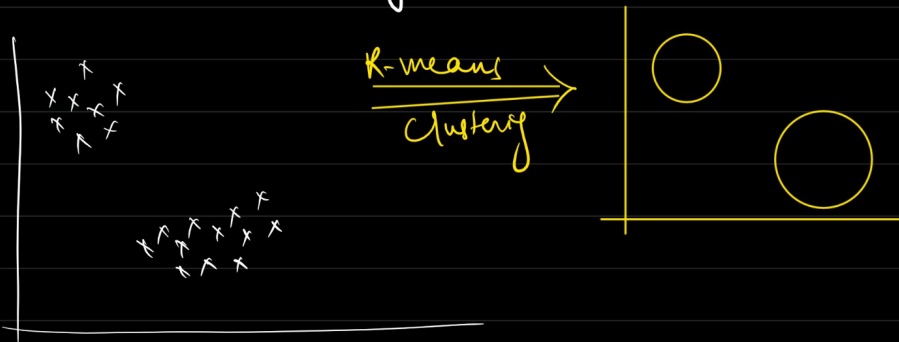→ Autoregressive models.

## USL

→ y (target variable is not there.
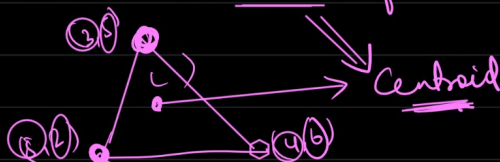→ finds groups/Pattern in the data.

## ① K- means clustering



K- means
↓        ↓
K        means (mean/Average)

$f_1$  $f_2$  $f_3$

→ (5   3   2)
→ (2   4   3)
→ (6   1   5)
─────────────
  13/3  8/3  10/3

Average → Arithmetic center of the data.



→ Centroid.

$$x_{cent} = \frac{5+3+4}{3}$$

$$y_{cent} = \frac{2+5+6}{3}$$

K-mean cluster

What ??

Steps of K-means clustering
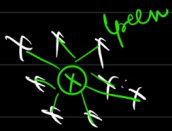
(i) initialise centroid (K)

→ Can be a random dp
→ Centroid of all dps
→ A random new dp's

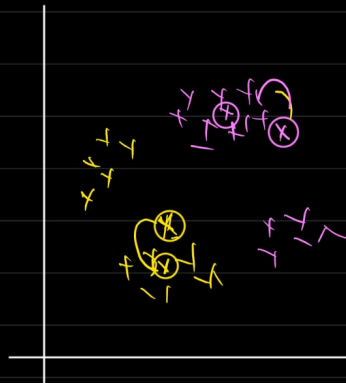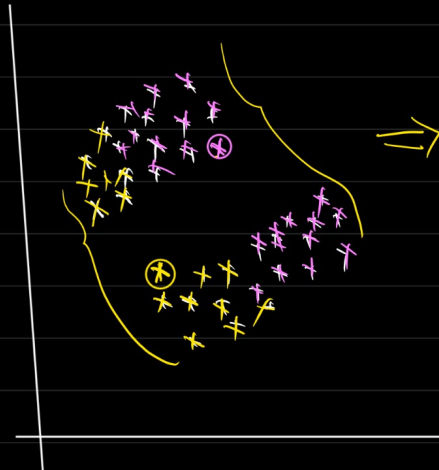(ii) Points nearest to the centroid will be labelled as that group.

(iii) Recalculate the centroid

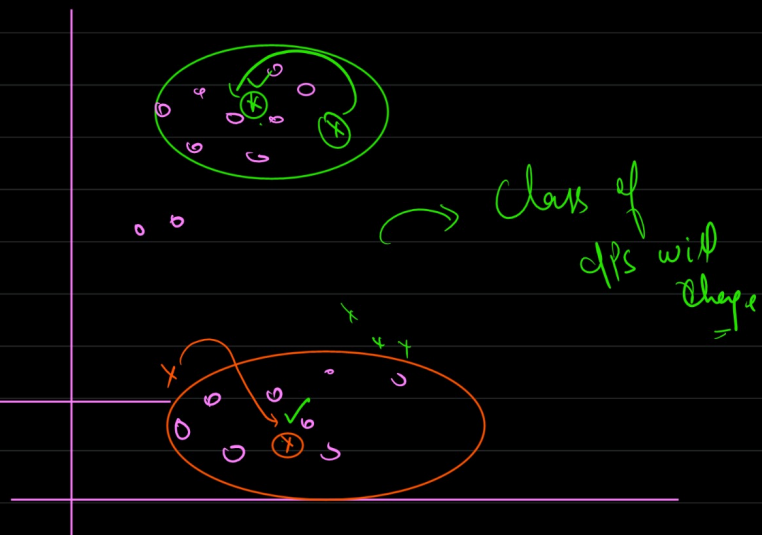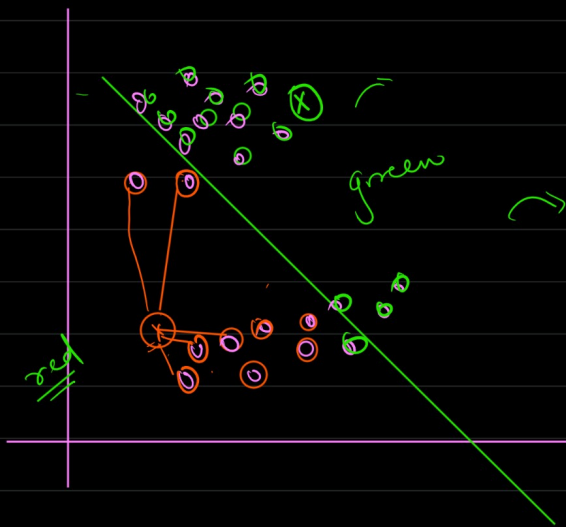→ Step 2 will be repeated untill centroid doesn't change.

⇒ label will change
⇓
New centroid.

Expenditure

green

Pink

B

Salary

green

red

Class of
d/ps will
change

distance
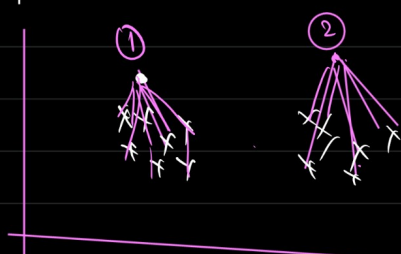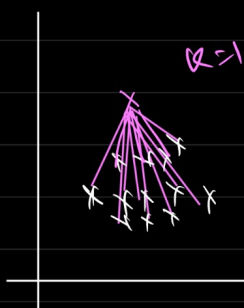
① Euclidean $\sqrt{(x_2-x_1)^2 + (y_2-y_1)^2}$

② Manhattan distance $= |x_2-x_1| + |y_2-y_1|$

* How to decide K.

WCSS → Within Cluster Sum of Square distance.

$$WCSS = \sum_{i=1}^{n} \left( \begin{array}{c} \text{distance b/w} \\ \text{Points to} \\ \text{nearest} \\ \text{Centroid} \end{array} \right)^2$$

K=1

① ② $WCSS_1 + WCSS_2$

K=2

As you increase K, WCSS decreases.

→ At optimal K, WCSS will not change.

WCSS

Elbow method.

optimal.

1   2   3   4   5   6   7   8   K.

if all the
dP are classified ⇐ { if class label not
to the nearest          Changing
Centroid                  ⇓
Class.              Centroid not change

*  if WCS s not
   Changing.
       ⇓
   All the dp's are
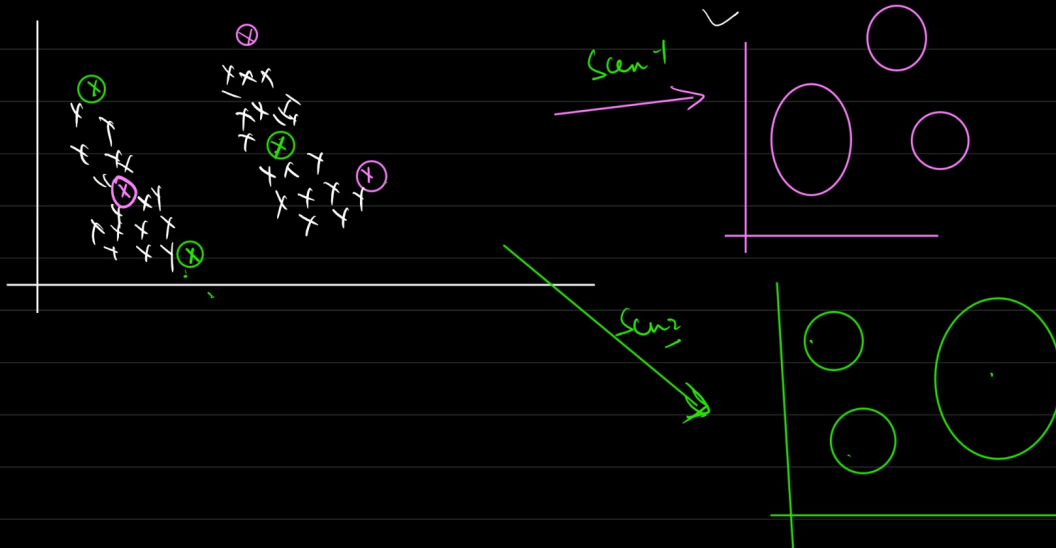   correctly marked to
   the nearest cluster class.
         ⇓
      optimal k.

Scen-1



6 units

WCSS

Sen-2



→ 3 unit.

How k is decided :-

✓ → elbow method.
  → business team.

* Random initialization k- trap



Scen-1 →

Scen-2

* initialize the k's as much far as possible
$\Downarrow$
K means ++