Any doubts ?)

Stats,
Prob,
Sampling,

Hypothesis test,

ML,

LR,

LR

Python Implementation!

Box plot in Python?
↙         ↓
Pandas?    function?

df. Series
↓
?

Heatmap?

Evaluation of Models?
Sklearn.metrics

Standard Scaler?
Sklearn. preprocessing

Prob. of Prediction! → which method
            of that ML model
model. Predict_proba ( features )
                        ↙    ↘
                     x train   x test

→ Clustering → Create Clusters / group ) segments → Data Mining

→ No Target required

Set of Similar objects → rows

=1    100 Rows

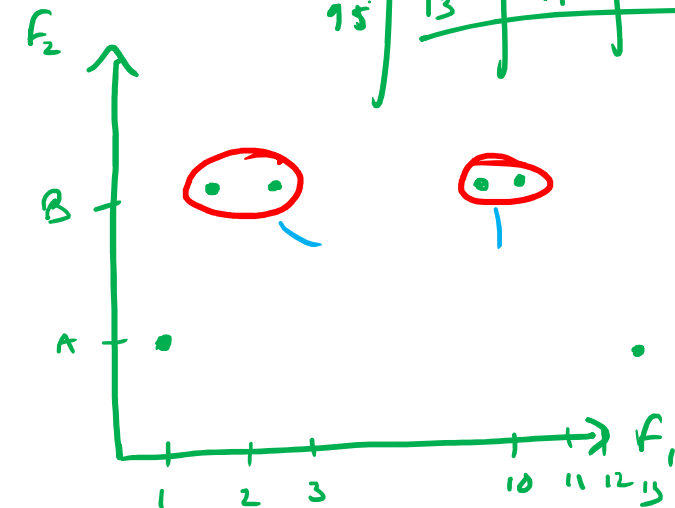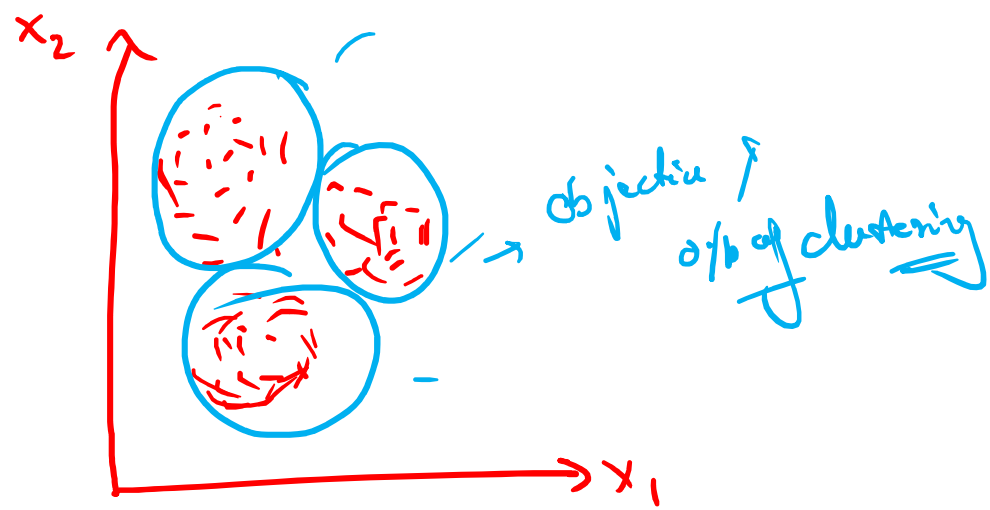Features → Clustering ML Algo → Clusters

L
Data

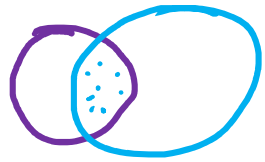No plotting required

→ Quality
Evolution

2500, 7580

| | $F_1$ | L $F_2$ | | |
|---|---|---|---|---|
| 90 | 1 | 4 | | |
| 91 | 10 | B | | |
| 92 | 2 | B | | |
| 93 | 1, | B | | |
| 94 | 3 | B | | |
| 95 | 13 | A | | |



Objective
O/p of clustering

$\underline{I}$

Incorrect cluster

$\underline{II}$

Bad Quality

$\underline{III}$

Good Clustering

Cluster → Collection of objects., which are "Perceived" to be similar b/w them
& dissimilar to objects belonging to diff cluster.

$\rightarrow$ K - Means

$\downarrow$

arbitrary Number

$K \Rightarrow \underline{1}$ , $(N) \rightarrow$ $\boxed{K = \{3, 20\}}$

$\downarrow$

Not so many clusters!

$\downarrow$

$1^{0.0}$

$1^{n}$

Statistically a healthy cluster

can have about $(5\%) - (35\%)$ of the total data

$(20) \rightarrow (3)$

$(1\%)$

$(50\%) \times$

100 Rows $\rightarrow$ Styles

$\downarrow$ $\searrow$

$100$ $1000$ clusters

1 Row in each cluster

$\dfrac{25}{\cancel{100}}$ $K = 4$ 

$\dfrac{100}{4}$ $(5\%)$

$2\%$

On every data $\rightarrow$ 20 clusters $\rightarrow$

$\downarrow$

Data $\rightarrow$

K = 1
K = 2
$\boxed{K = 3}$
K = 4
$\boxed{K = 5}$
$\vdots$
$K = 20$

$\rightarrow$ we will select the best value of $(K)$ on that dataset

$\downarrow$

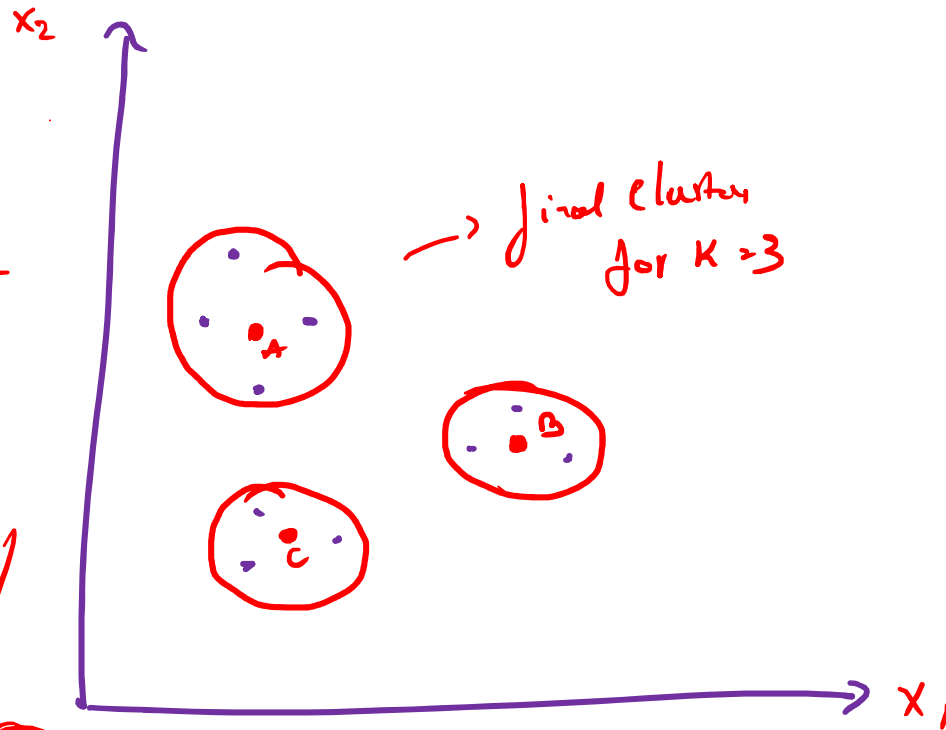final output

→ K-Mean Algorithm

K = 3

③ Groups of shortest distance will be center

④ The centers will be recalculated / shifted

⑤ Repeat from step ②
Till in consecutive 2 Iterations you get **same** centers



→ final cluster for K = 3

(axes: $X_2$ vertical, $X_1$ horizontal, clusters A, B, C)

① will Add K **Random points** to your data
These are known as centers / Centroids / Mean

② from each data point, it will calculate the distance to **each** added center

| | A | B | C |
|---|---|---|---|
| → 1 | 3 | 5 | 4 |
| → 2 | 4 | 2 | 3 |
| → 3 | — | — | — |
| | — | — | — |
| 5 | — | — | — |
| 2 | — | — | — |
| | — | — | — |