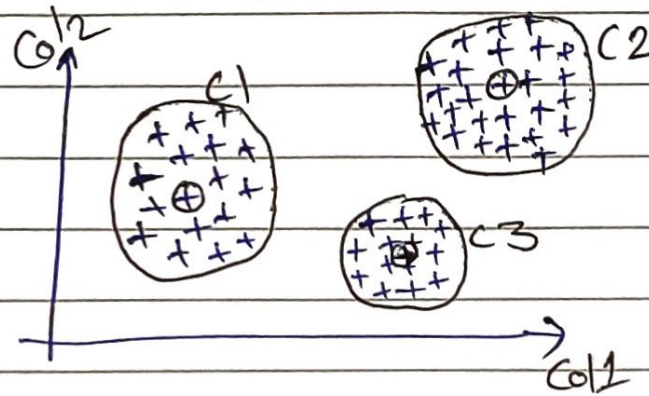
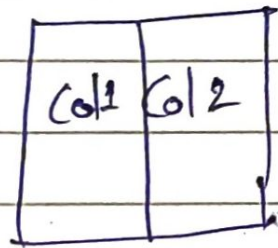


K-Means Clustering

Group them into "K" clusters

"K" no. of clusters.

For every cluster it will assign Centroid (to each cluster)



$C_1, C_2, C_3 \rightarrow$ Clusters

Set of points in $C_1 \rightarrow S_1$

" " " " " $C_2 \rightarrow S_2$

" " " " " $C_3 \rightarrow S_3$

*) $D = S_1 \cup S_2 \cup S_3 \rightarrow$

Dataset

This gives my Complete Data

*) $S_1 \cap S_2 = \phi$

$S_2 \cap S_3 = \phi$

$S_1 \cap S_3 = \phi$

} properly clustered

\rightarrow each data point belongs to only 1 cluster

\rightarrow not even single datapoint should belongs to 2 clusters

\downarrow mutually exclusive.

K-Means

↳ K-Centroids

↳ K-Set of points (S_1, S_2, \dots, S_K)

↳ $S_1 \cup S_2 \cup S_3 \cup S_4 \dots \cup S_K = D$ - Dataset

↳ Mutually exclusive sets (S_1, S_2, \dots, S_K)

$$S_1 \cap S_2 = \phi$$

$$S_2 \cap S_3 = \phi$$

$$S_1 \cap S_3 = \phi$$

Q) How to find Centroid.

Ans:- $C_i = \frac{\text{Sum of Datapoints}}{\text{\# of Datapoints}} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$

nothing but Mean

$$C_i = \frac{\sum_{j=1}^n x_j}{n}$$

We can also use Median as well

$$C_i = \frac{1}{n} \sum_{x_j \in S_i} x_j$$

datapoints in Set 1.

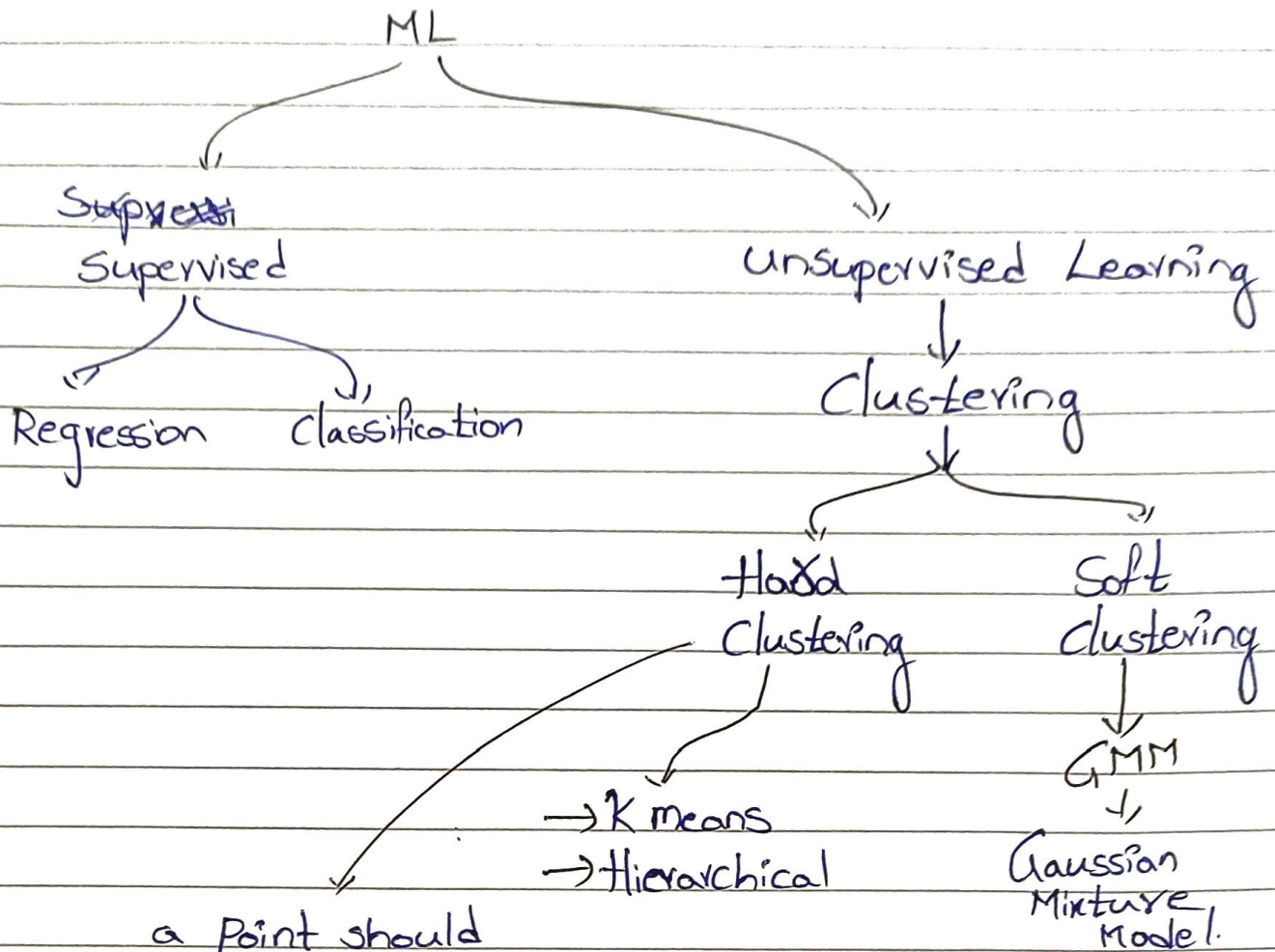
Size of S_i

we are using Mean

so, it will be impacted by outlier.

$$n = |S_i|$$

$$C_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$

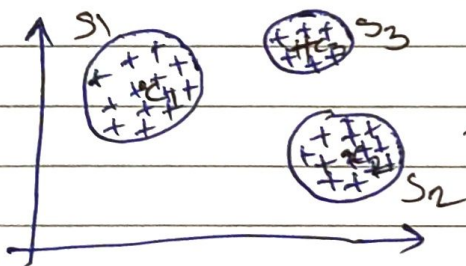


a point should belong to only 1 cluster.

↳ It means intersection b/w two clusters = \emptyset

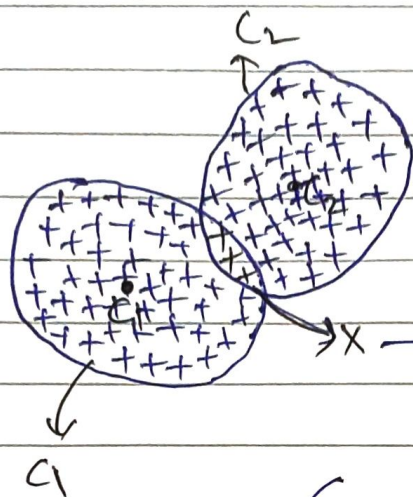
$$\left. \begin{aligned} \hookrightarrow S_1 \cap S_2 &= \emptyset \\ \hookrightarrow S_2 \cap S_3 &= \emptyset \\ \hookrightarrow S_1 \cap S_3 &= \emptyset \end{aligned} \right\} \rightarrow \text{Mutually exclusive}$$

Then it is known as Hard Clustering



→ Hard clustering.

Soft clustering! \rightarrow (Ex!-GMM)



- \rightarrow Two clusters are overlapping
- \rightarrow It has 3 intersection points
- \rightarrow 3 data points are having/ belongs to both clusters.

$$P(x) \begin{cases} \rightarrow C_1 \rightarrow 0.6 \\ \rightarrow C_2 \rightarrow 0.4 \end{cases}$$

(The probability of the data point that belongs to C_1 (centroid 1) (or) C_2 (centroid 2) and

\downarrow
which is having high probability
 \downarrow

Then that data point belongs to high probability Centroid.

K-Means

\hookrightarrow Hyperparameter $\rightarrow K \rightarrow$ (Total no. of clusters)

(Note! In Geo-metric Algorithms

\hookrightarrow Linear, Logistic, SVM..... \rightarrow In all these algorithms

we are looking at optimization eqn.

K-means \rightarrow distance \rightarrow Euclidean distance (In order to perform for K-means)

\rightarrow As of now there is no optimization equation for K-means

(Actually there are optimization eqn for K-means but the issue is it will take very high computational power for optimization equation.
 \hookrightarrow As of now research is going on but not yet proposed

\downarrow
Because of this reason (no optimization eqn for K-means) we will use Lloyd's Algorithm.

\downarrow
That's why K-means is also known as Lloyd's Alg

Task! \rightarrow Need to form K-clusters such that intracluster distance is low and intercluster distance is high

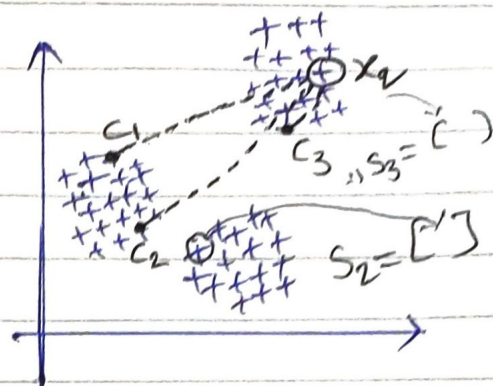
1st step! \rightarrow Randomly pick 'K' data points from the Dataset (D) and call them as $c_1, c_2, c_3, \dots, c_k$ (Centroids)

(This is Initialization step)

→ Initialization Step

$c_1, c_2, c_3 \rightarrow$ These 3 Centroids are randomly chosen.

(This are not the real Centroids, So, we need to update the Centroid)



Step 1: \rightarrow

For each point x_i in the Dataset \rightarrow Select the nearest the Centroid C_j

\rightarrow Then add x_i to Set S_j .

How to select is calculate distance from x_q to C_1 and x_q to C_2 and x_q to C_3

\rightarrow Get the Centroid which is having minimum distance

\rightarrow In the above example $x_q \in C_3$ (Centroid 3)

\rightarrow So, it belongs to Set 3 (S_3)

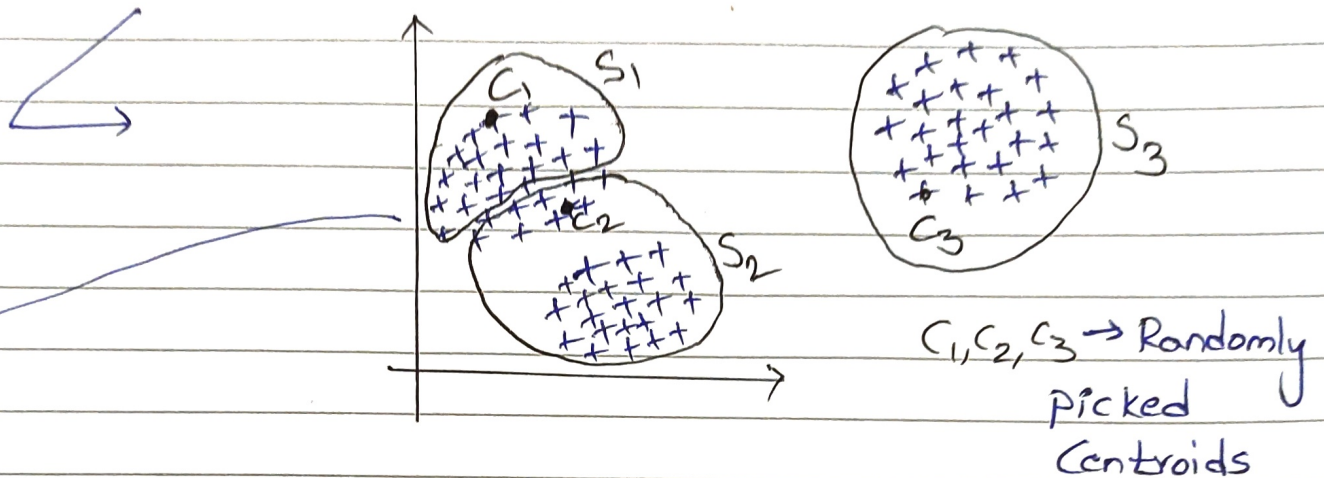
\rightarrow In this way we will calculate distance from every point in the dataset to centroids and we will store it into the sets (it belongs to which set then
 \rightarrow means storing it in respective set

In Simple:-

For each point x_i in the Dataset

↳ select the nearest centroid C_j

↳ Add x_i to set S_j



(This sets are formed according to the distance (nearest Centroid))

(This is called Assignment Step \rightarrow becoz assigning each datapoint to corresponding set a/c to distance nearest Centroid)

Step-3:-

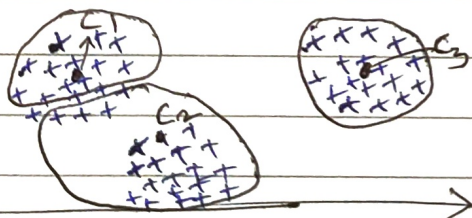
↳ Recalculate C_j (centroid) for each Set S_j .

(This step is called Re Computation of Centroid.)

$$C_j = \frac{1}{|S_j|} \sum_{x_i \in S_j} x_i \rightarrow \text{Sum of datapoints.}$$

↳ no. of data points

Centroids are updated



\Rightarrow Centroids positions are changed.

Step-4! →

Repeat Step-2 & Step-3 until Convergence.

↓
Centroids won't change.
(i.e. sets won't change)

(and these clusters are mutually exclusive)

$$\left\{ \begin{array}{l} S_1 \cup S_2 \cup S_3 = D \\ S_1 \cap S_2 = \emptyset \\ S_2 \cap S_3 = \emptyset \\ S_1 \cap S_3 = \emptyset \end{array} \right.$$