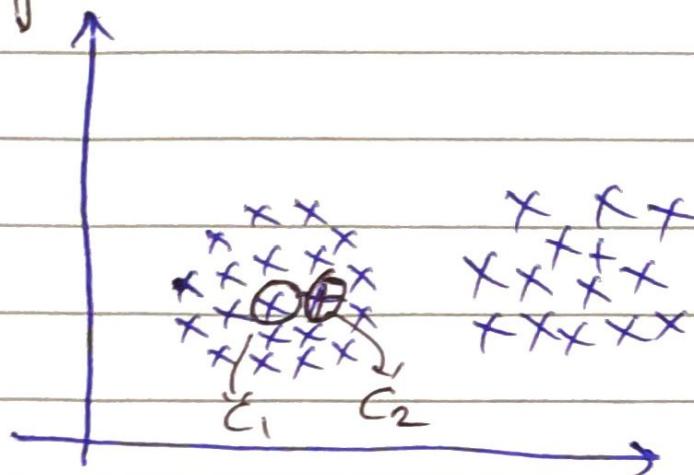


Q) What if all the selected Centroids (at initialization- randomly picking the centroids at start) are next to each other?



For large datasets the probability of getting nearest (boss) centroids next to each other is very very less ≈ 0

→ But for small datasets then probability of getting Centroids next to each other is won't be 0.

↳ chances are there for small datasets

→ In real time, we have millions of data then it won't have problem (little bit chances are there)

Q) So, is there any better initialization of Centroids (i.e step 1)
A) { In order to
It can't solve the issue of Getting centroids next to }
each other

→ Yes, there is a better initialization
↓
K-Means++ initialization.

In K-means.

Step-1: Random Initialization

↳ Decide the value of K & pick " K " points from the Dataset D & call them as $C_1, C_2, C_3, \dots C_k$

↳ So, Rather than randomly picking, we will replace this step with K-Means++ initialization.

↓
→ Move Computational heavy task

K-Means

Step-1. → Replace with K-Means++ initialization

Step-2

Step-3

Step-4

} Same

K-Means ++

↪ It is nothing but K-Means algorithm by improving

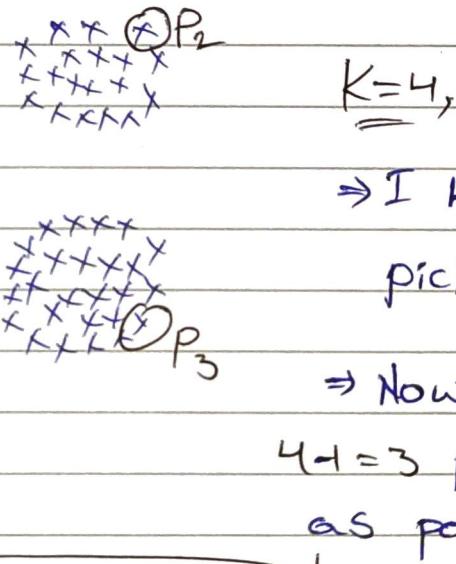
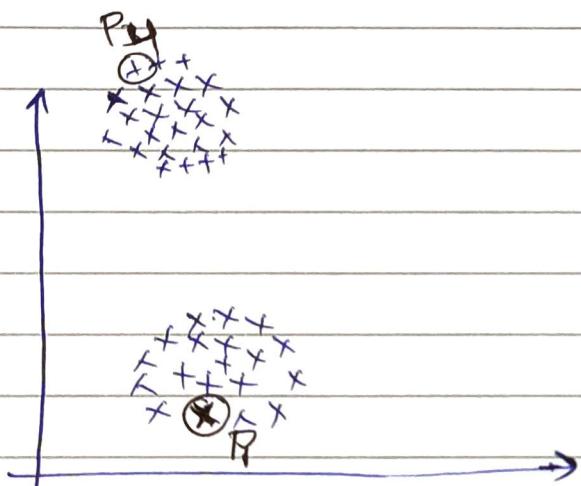
Not selecting (doing) 1st step.

↪ Replacing with K-Means ++.

K-Means:

Step-1: Decide the value of K and randomly pick K -Centroids c_1, c_2, \dots, c_K .

K-Means++
Step-1: Decide the value of K and randomly pick ~~K~~ 1 Centroid
→ And pick $K-1$ Centroids as far as possible.



→ then i will pick $P_2 \rightarrow$ it is farthest having large distance b/w them ($P_1 \& P_2$).

- Now I will pick P_3 as it is having high distance from $P_1 \& P_2$.
- In the same way we will pick P_4 as it is having high distance from $P_1, P_2 \& P_3$.
- Once we pick $P_1, P_2, P_3 \& P_4$ (Let's say $k=4$) then our next steps will be done quickly (step-2, 3, 4), as
 - the Centroids will Converge quickly (sets won't change).
(Because in the above example we got the centroids easily as far as possible then easy form the sets)
 - But Getting $P_1 \& P_2 \& P_3 \& P_4 \rightarrow$ takes very high computation after that it won't take much time to converge; as it is easily clustered (In the above example points are present in each cluster)

Q) How do you decide best value of "K"

Ans:

For this we need to calculate WCSS

\downarrow
Within cluster sum

Clusters: →

of squares.

(a) Intracluster distance to be low

(b) Intercluster distance to be high

∴ So, for good clusters, Intra cluster (distance with the Centroid inside a cluster) should be low.

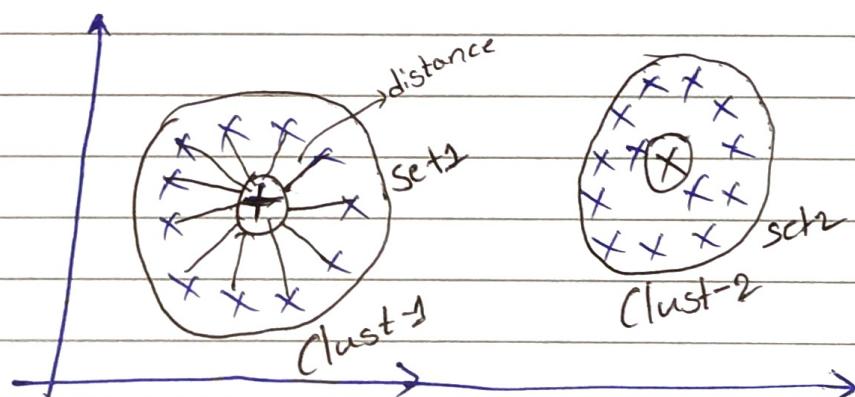
∴ So, i am calculating distance from centroid to point



$$\hookrightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \Rightarrow \|x - c\|^2$$

(x_1, y_1) (x_2, y_2)

⇒ we can represent distance as $\|x - c\|^2$



for clust-1 (centroid-1)

$$\sum_{x \in S_1} \|x - c_1\|^2$$

for clust-2 (centroid-2)

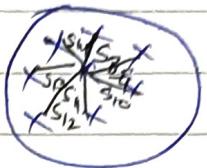
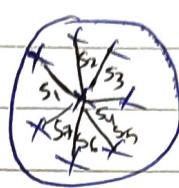
$$\sum_{x \in S_2} \|x - c_2\|^2$$

\Rightarrow Now, I will sum up intracenter distance for the 2 centroids. (And it has to be low then only my clusters are good)

$$\rightarrow \sum_{x \in S_1} \|x - c_1\|^2 + \sum_{x \in S_2} \|x - c_2\|^2$$

$$\rightarrow \sum_{i=1}^K \sum_{x \in S_i} \|x - c_i\|^2$$

$$\Rightarrow (s_1 + s_2 + s_3 + \dots + s_7 + s_8 + s_9 + \dots + s_{14})$$



Cluster-1 $\Rightarrow s_1, s_2, \dots, s_7 \rightarrow$ distances from centroid c_1 to points

And this sum has to be low.

Cluster-2 $\Rightarrow s_8, s_9, \dots, s_{14} \rightarrow$ distances from centroid c_2 to points

$$c_1^*, c_2^*, \dots, c_K^* = \arg \min_{c_1, c_2, \dots, c_K} \left\{ \sum_{i=1}^K \sum_{x \in S_i} \|x - c_i\|^2 \right\}$$

Optimization Eqn.

K-Means \Rightarrow it is also known as Lloyd's Algorithm (4 steps)

Q) Why can't we apply Gradient Descent on optimization Eqn so that we can get Centroids.

Ans)

It's because this optimization Eqn is of one of

the hardest equation and it is really hard to solve. That's the only reason we are going with the Lloyd's Algorithm and KMeans++ Initialization.

→ Even the libraries like Sklearn, ... (anything which uses) uses Lloyd's algorithm & ~~so~~ it won't use Gradient Descent for KMeans++

(Any other library which are used for KMeans++, then it must uses Lloyd's Alg rather than Gradient Descent.)

Note: And good part about this optimization equation is

that we can leverage this part to find out K.

$$c_1^*, c_2^*, \dots, c_k^* = \arg \min_{c_1, c_2, \dots, c_k} \sum_{i=1}^K \sum_{x \in S_i} \|x - c_i\|^2$$

→ (This is intra Cluster Distance.)

K-Means++ :

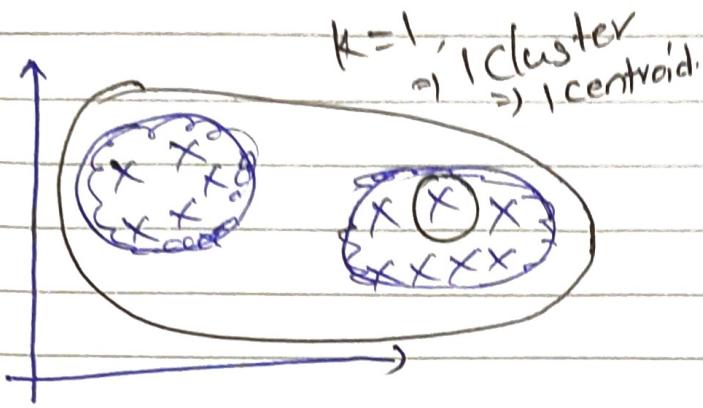
a) Decide the value of K.

Step-1: → K-Means++ initialization

Step-2: → Assignment (assigning each data point to Corresponding Set)

Step-3: → Recomputation (of centroids)

Step-4: → ~~Converge~~ Repeat 2 & 3 until Convergence



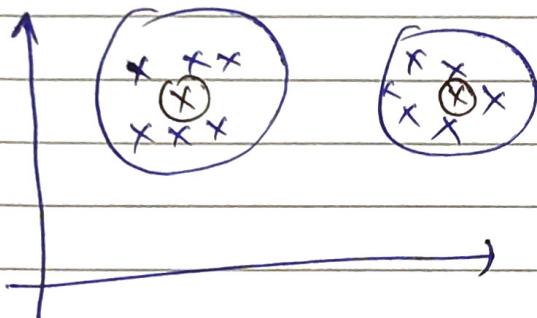
Let's take $K=1$,

and do the

steps-1, 2, 3, 4.

→ Calculate Intracluster distance
(d_1)

$K=2$

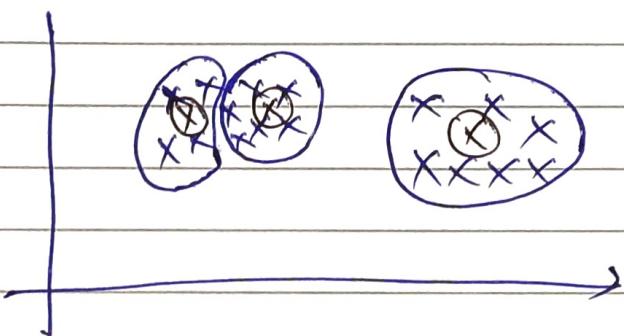


$K=2$, (

Now, do the step-1, 2, 3, 4)

→ Cal Intracluster distance (d_2)
↳ Here $d_1 > d_2$

$K=3$



$K=3$

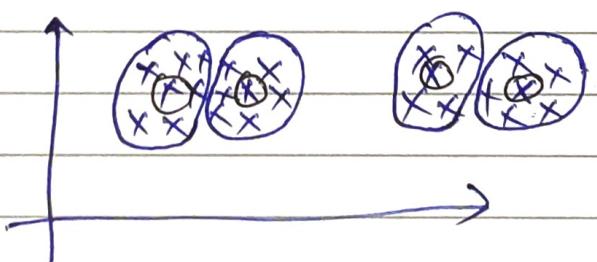
Now, do the step-1, 2, 3, 4

→ Cal

Intracluster distance (d_3)

↳ Here $d_1 > d_2 > d_3$

$K=4$



$K=4$,

Now do the step-1, 2, 3, 4

→ Cal

Intracluster distance (d_4)

↳ Here

$d_1 > d_2 > d_3 > d_4$

$$d_1 > d_2 > d_3 > d_4$$

→ My ideal K is 2

but My Equation is saying $K=4$ is good

because of Very less Intracluster distance
(d_4)

↳ It's because we are overfitting the data
(clusters ($K=4$) are overly fitted)
(It's like Tight genes part)

→ and for $K=1$,

My Equation is saying d_1 is very high

So, $K=1$ is not good.

→ 1 cluster for all the datapoints → so d_1 is very high
(for this example)

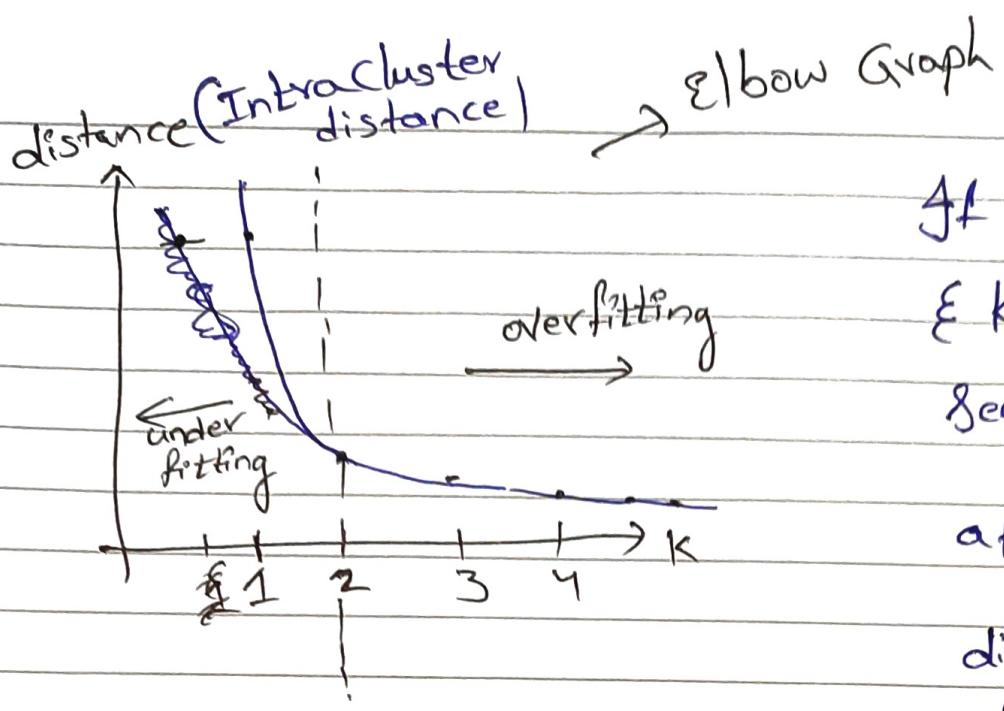


It like loose genes part



Underfitting.

→ "K" is the hyperparameter.



If we plot the distances ϵk then we can see sudden drop at $k=2$. After that very small difference in the next $k (k>2)$

at $k=2$

↳ Intracluster distance will fall very high

for $k > 3$ ↳ Intracluster distance will gradually reduce (slowly)

→ So in this way we use to get optimal no. of clusters.