

Kmeans clustering:

original data: ~~AAAA BBBB OOO~~

Let's see if computer can make these 3 cluster

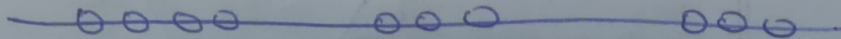


Step 1

Select the no of cluster you want to identify in your data.

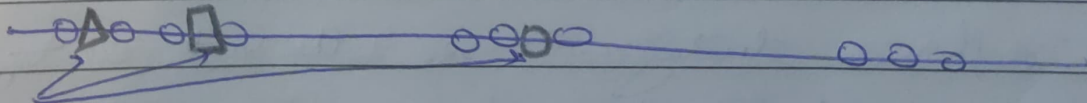
This is K in Kmeans clustering

In this case $K = 3$



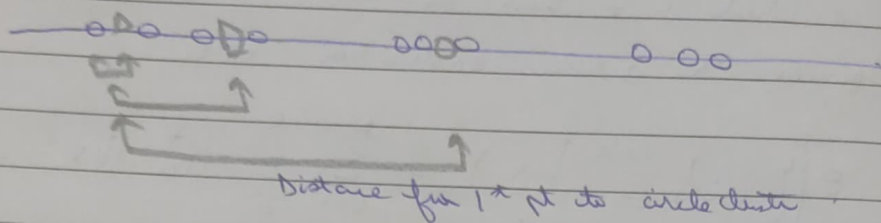
Step 2

Randomly select 3 pts (distinct)

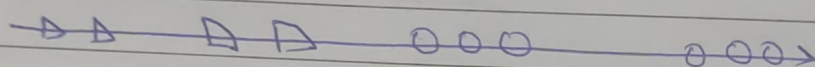


Initial Cluster

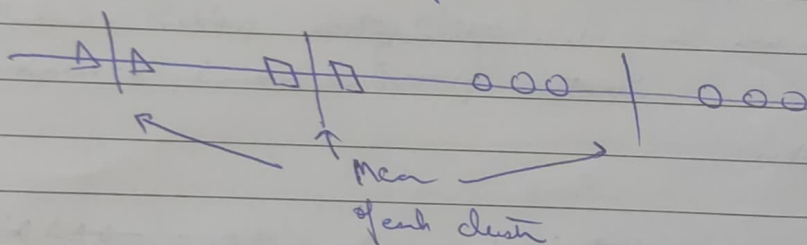
Step 3: Measure the distance b/w 1st pt at three initial clusters



Step 4: Assign 1st pt to nearest cluster



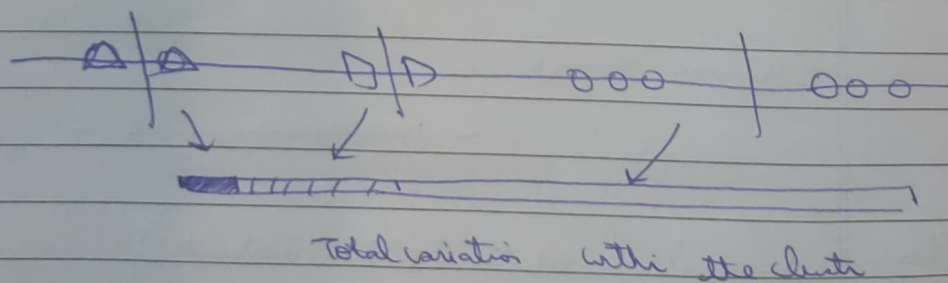
Step 5: Calculate mean of each cluster



The we repeat what we just did (measure all clusters) using the mean values.

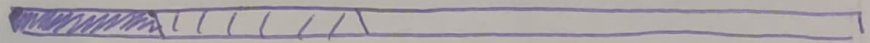
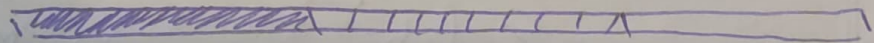
In this case clustering won't change much

We can assess the quality of clustering by adding up the variation within each cluster



Since K means clustering can't see the best clustering, its only option is to keep track of these clusters, all their total variances, and do the whole thing over again with diff starting pts

So we start again and then cluster all the remaining pts, calculate mean of each cluster, re-cluster based on new means. It repeats until clusters no longer change. Now that the data are clustered, we sum the variation within each cluster.

1st attempt2nd attempt3rd attempt

At this pt, K means clustering knows that 2nd attempt is best clustering so far. But it does not know if it's best overall, so it will do a few more clusters (it does as many times as we tell) and come back and return that one if it is still the best.

You can find by using elbow plot.

Q) $X = \{2, 3, 4, 10, 11, 12, 20, 25, 30\}$
 Perform K means clustering for $K=2$

n) $K=2 \therefore K_1, K_2$ - 2 clusters

Step 1 Randomly split (2 means) $m_1=4, m_2=12$

$$\begin{matrix} (2-4) \times 2 = 2 \\ (2-12) \times 2 = 10 \end{matrix} \quad K_1 = \{2, 3, 4\}$$

$$K_2 = \{10, 11, 12, 20, 25, 30\}$$

Step 2 $m_1 = \frac{2+3+4}{3}$

$$m_2 = \frac{10+11+12+20+25+30}{6}$$

$$= 3$$

$$= 18$$

$$K_1 = \{2, 3, 4, 10\}$$

$$K_2 = \{11, 12, 20, 25, 30\}$$

Step 3 $m_1 = \frac{2+3+4+10}{4}$

$$m_2 = \frac{11+12+20+25+30}{5}$$

$$= 4.75$$

$$m_2 = 19.6$$

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$K_2 = \{20, 25, 30\}$$

Step 4 $m_1 = \frac{2+3+4+10+11+12}{7}$

$$m_2 = \frac{20+25+30}{3} = 25$$

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$K_2 = \{20, 25, 30\}$$

$$\boxed{m_1 = 7}$$

$$\boxed{m_2 = 25}$$

Since it giving same value for cluster m_1 and m_2 ,
 we have to stop.

$$K_1 = \{2, 3, 4, 10, 11, 12\}$$

$$K_2 = \{20, 25, 30\}$$

$$m_1 = 7$$

$$m_2 = 25$$