

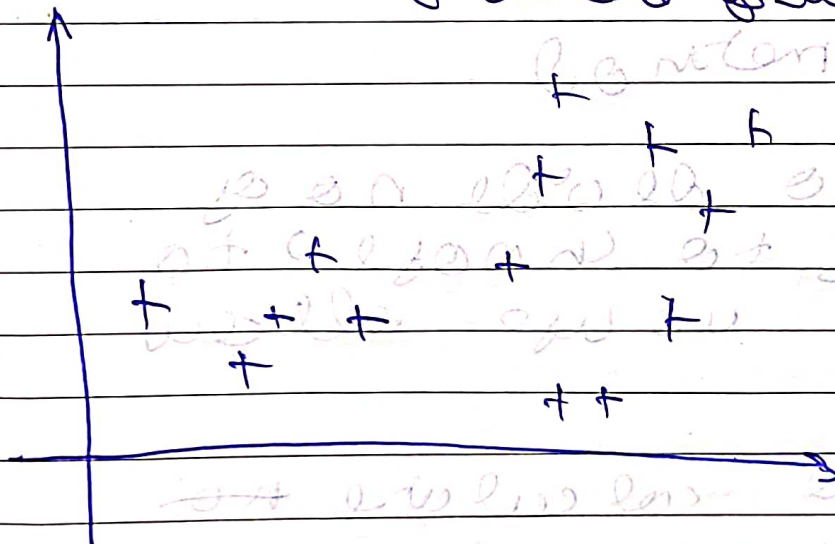
→ Clustering models

* This is part of unsupervised learning in this model is not given labelled data but it is given unlabelled data and model has to group it.

(model doesn't have any reference)

→ k-means clustering

(we say we are given this)



(we will decide how many clusters we want)

Then we will find

For every place a centroid (number depend on no of clusters)

Then we assign data to groups Based on the centroid

then we calculate center of mass of data points (including centroid) and then we move the centroid there and repeat the process

until this does not change anything we do this

→ Elbow method

How to decide no of clusters to choose to do this we use elbow method

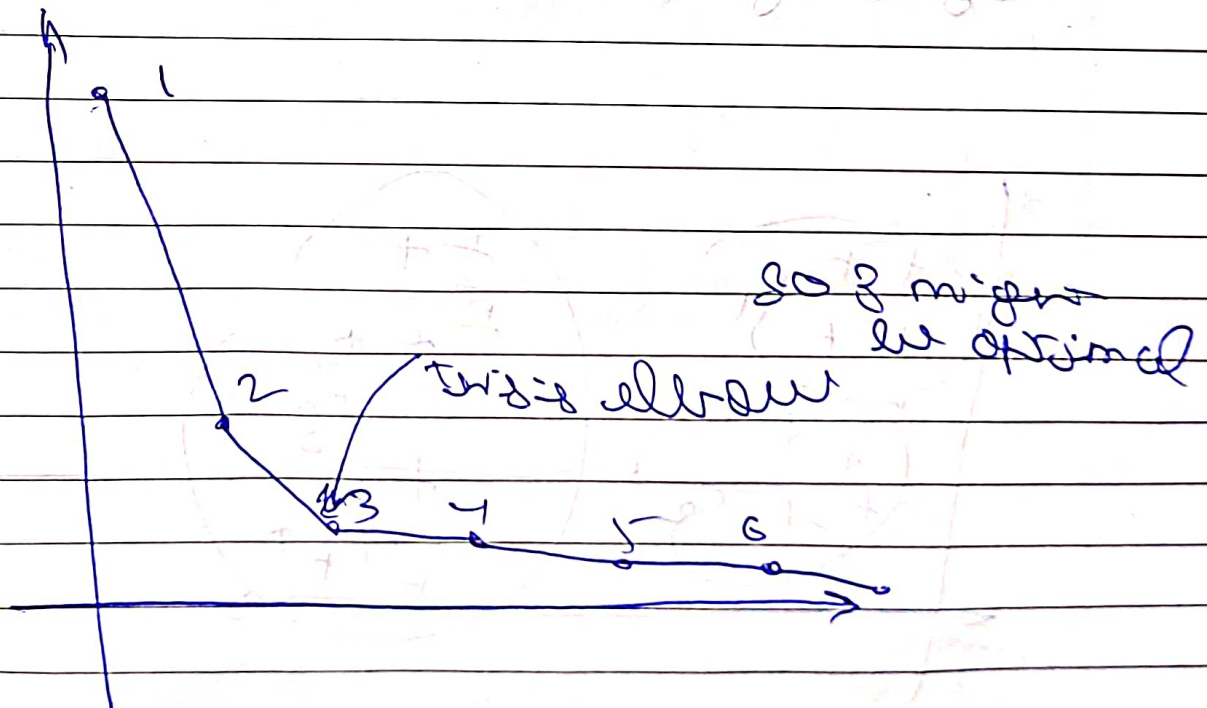
we calculate the ~~of clusters by~~
 we use elbow method

$$WCSS = \sum_{P_i \in \text{cluster } 1} \text{dist}(P_i, c_1)^2 + \sum_{P_i \in \text{cluster } 2} \text{dist}(P_i, c_2)^2 + \dots$$

$P_i \in$
 cluster 1
 $P_i \in$
 cluster 2
 +

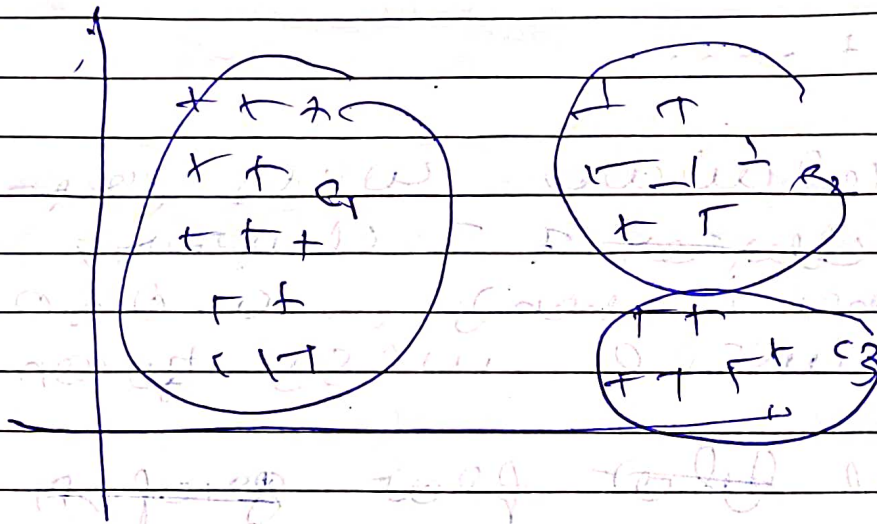
we calculate WCSS for
 cluster 1 & 2 cluster,
 then 2 then 3 ... and so
 on until WCSS becomes

then we plot for graph

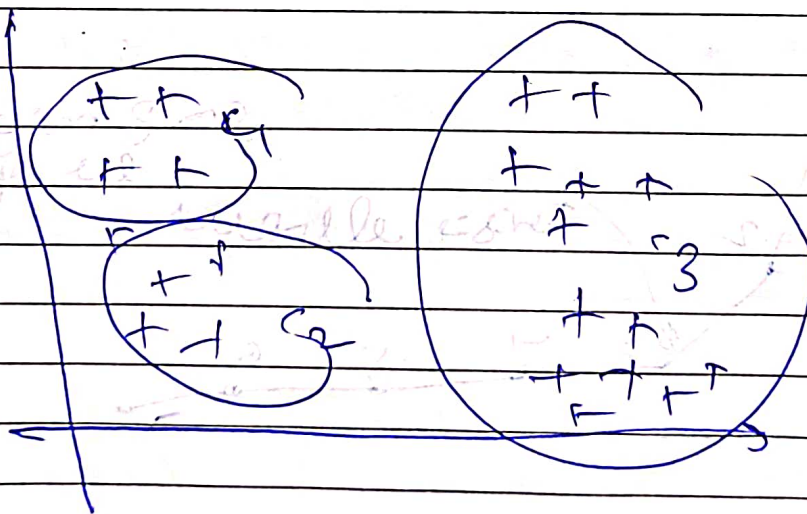


→ $k \rightarrow$ means $++$

you have data set



But if centroids are initial
of 3rd cl



we have different
results

Random Initialization

In k means++

- 1) first we choose k centroids at Random
- 2) we compute distance (D) to the nearest centroid already selected
- we choose the next centroid using weighted Random selection (weighted by D^2)
- 3) Repeat 2 and 3 until all the centroids are selected
- 4) apply k -means

→ Hierarchical clustering

There are 2 types

1] Agglomerative (Bottom up)

2] Divisive (covered in sheets)

1] first make each data point a single cluster
↓
N clusters

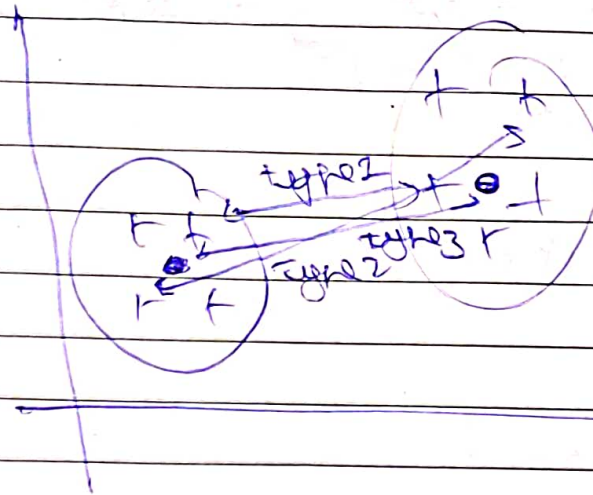
2] take two closest data points and make them one cluster
N-1 clusters

3] take two closest clusters and make them one N-2 cluster
↓
Repeat steps until only one cluster

(FIN)

[we can use any distance
But prefer Euclidean]

How to measure distance between 2 clusters



T1 = dist
between
cluster
pairing

T2 distance between furthest points

T3 - distance between centroids

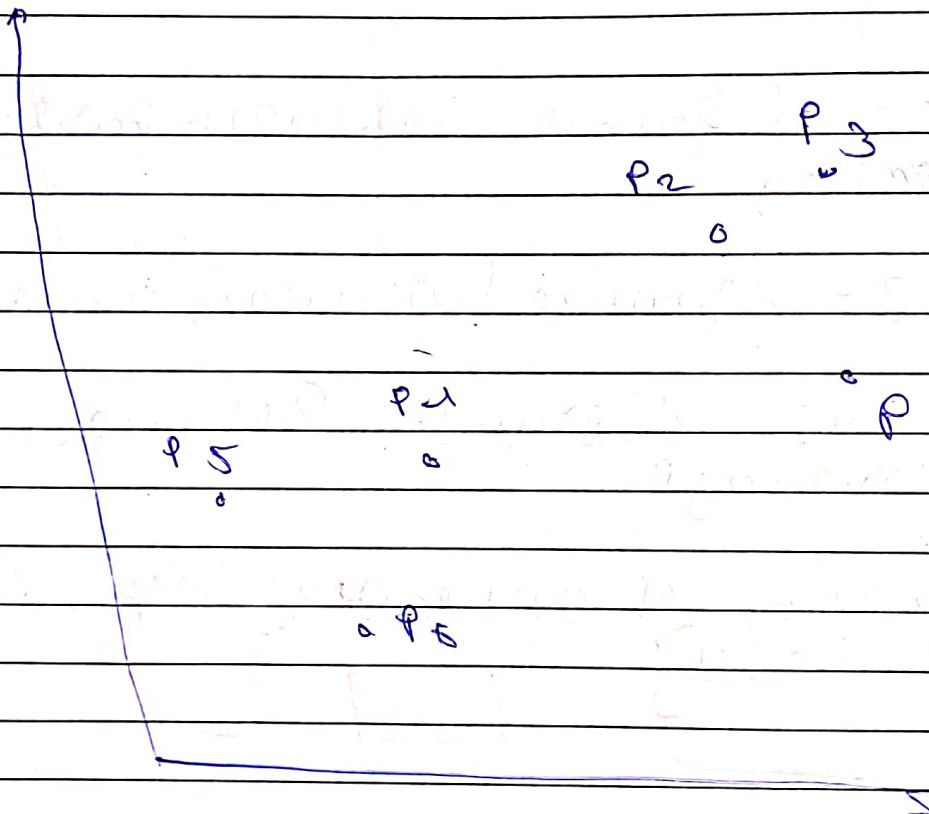
or T4 - distance between average

[we can choose any depending on data]

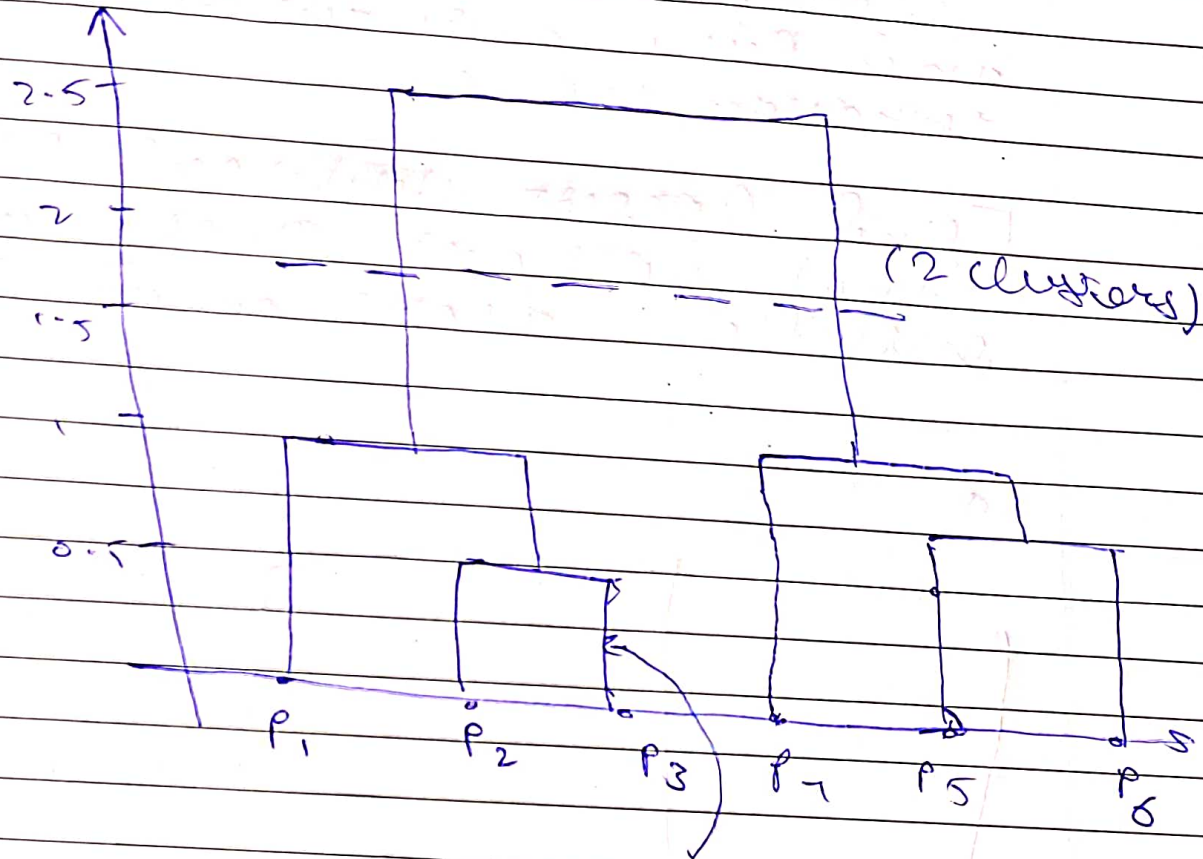
→ Dead & greedy

In T-IC the computer maintains the memory of how the algorithm was performed

Suppose we have



Dendrogram



Representing
the
dissimilarity
(distance between
points)

We decide a threshold of
dissimilarity

(no of lines threshold
crosses are the no of
clusters)

we find the largest distance
 (largest distance) and
 and then cut through them
 through them

Find largest vertical line
 that will not cross any
 intersect horizontal line

29

