

$$VIF = \frac{1}{1-R^2} \Rightarrow [1, \infty)$$

$(0, 1)$ ↙

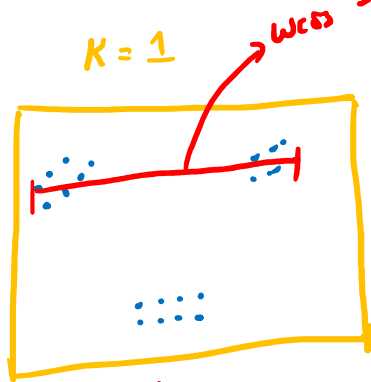
$\geq 5, \Rightarrow$ high corr

< 5 , No or weak corr

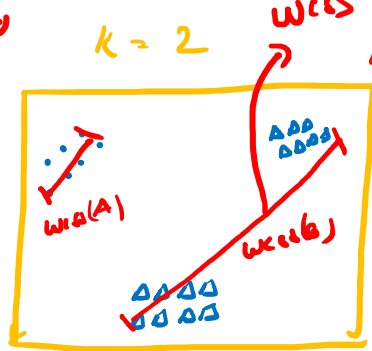
K-means?

→ Choosing the best value of K → $\{1, 2, \dots\}$

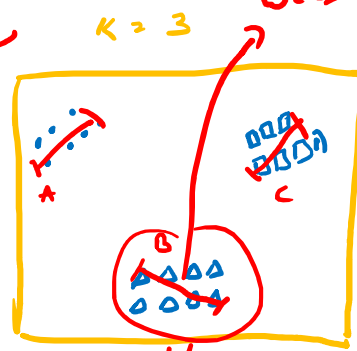
Determined using wcss method / elbow method



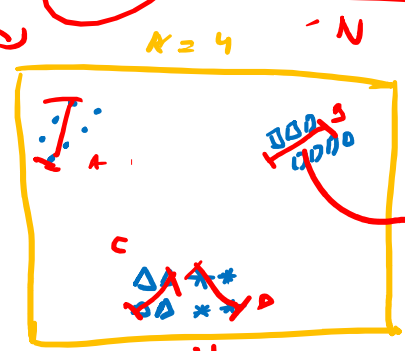
1 - cluster



2 - cluster



3 - clusters



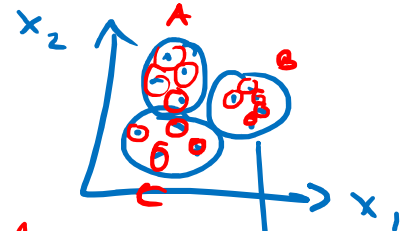
4 - clusters

(K=5)

for multiple values we will choose the $\max(WCSS_1, WCSS_2, \dots)$

$WCSS / WSS$
(within-cluster sum of sq)

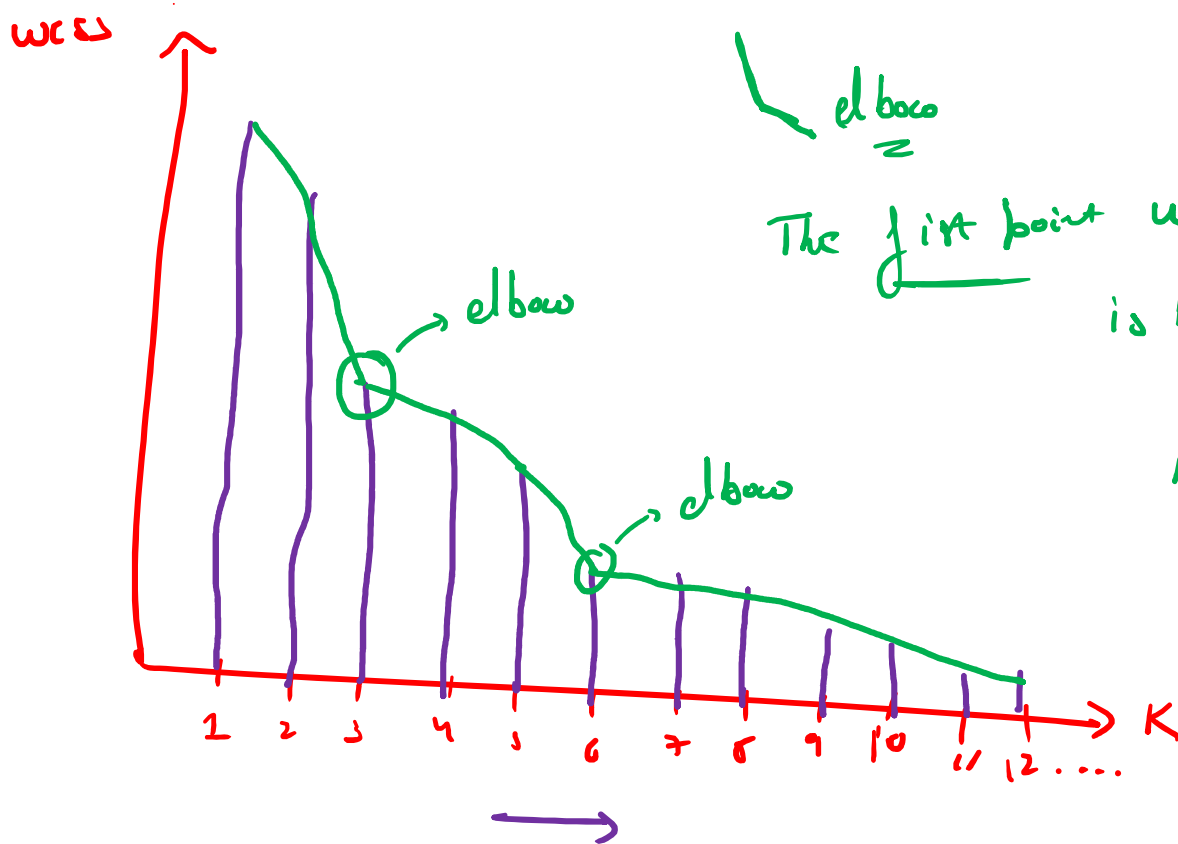
within the cluster



Spread centroid

$$SS = \frac{\sum (x_i - \bar{x})^2}{N}$$

Plotting the value of w_{css} vs K



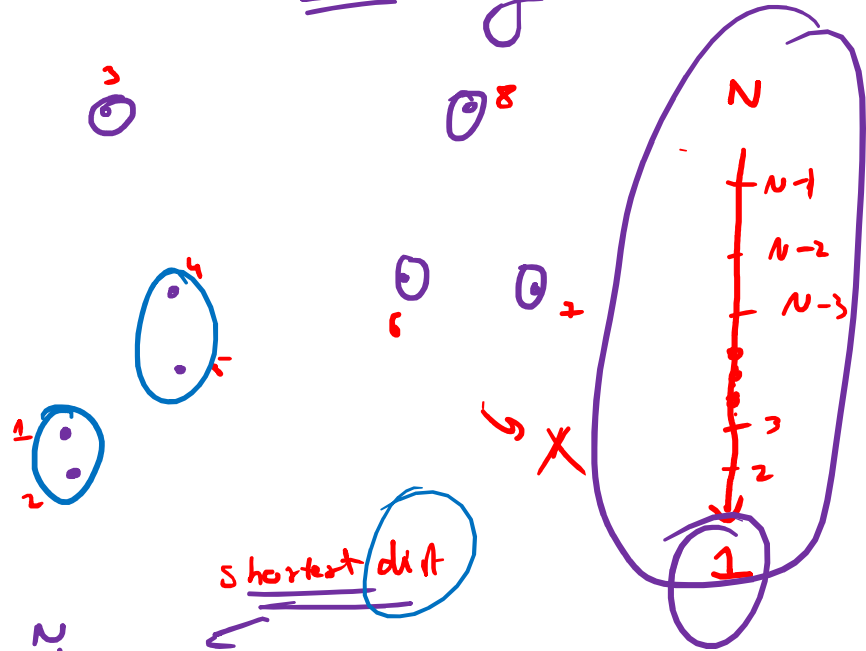
elbow

The first point where the elbow is formed is the best value of K

In case, the elbows are formed close to each other

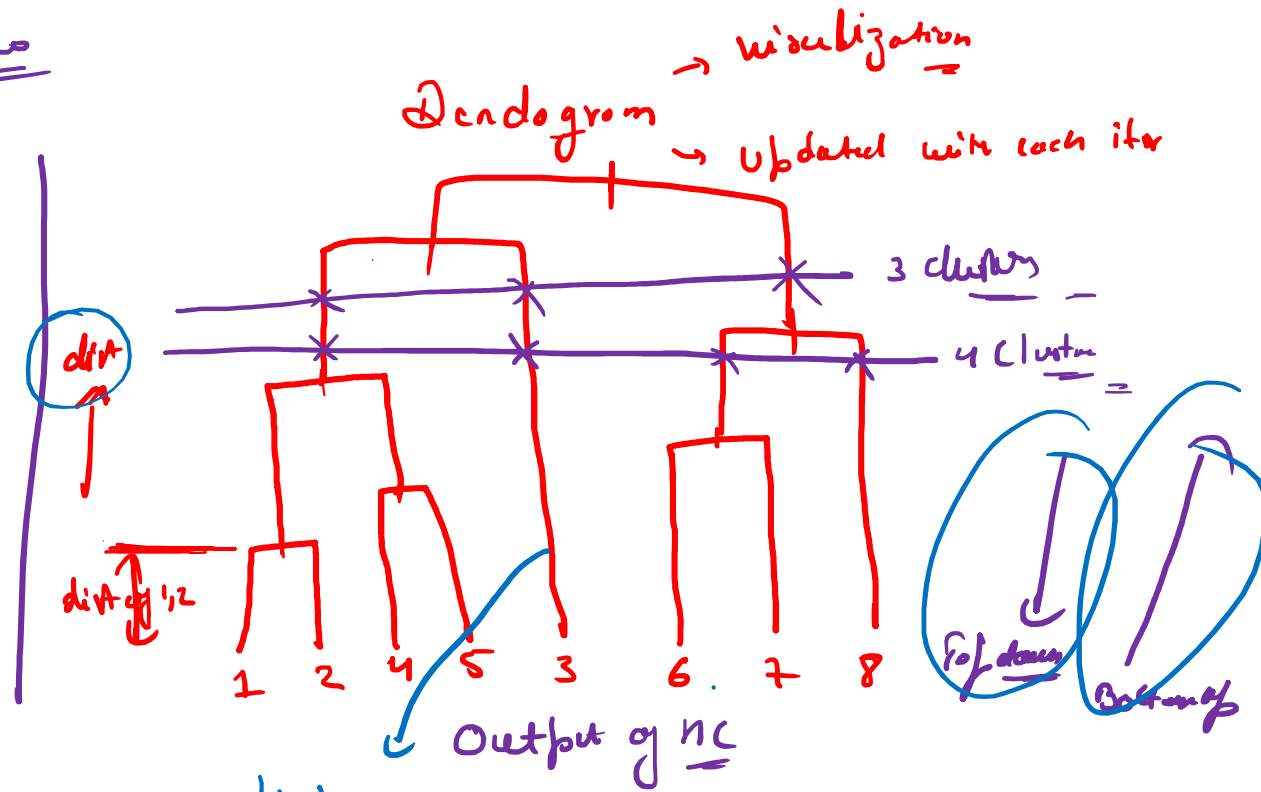
then business needs to decide on what one to keep
 \downarrow
 SME

→ Hierarchical Clustering (HC) → Flow



\downarrow Bottom up → Agglomeration HC

\uparrow Top up → Divisive HC (farthest first)



Linkage
linking the two groups together using certain calculations

→ Difference in HC, K-means

① HC cannot be used for a large dataset
→ too much time
↳ $O(n^2)$
Dendrogram will be too complex
↓
difficult to analyze

K-means can be used for a large dataset → $O(n)$ → faster than HC

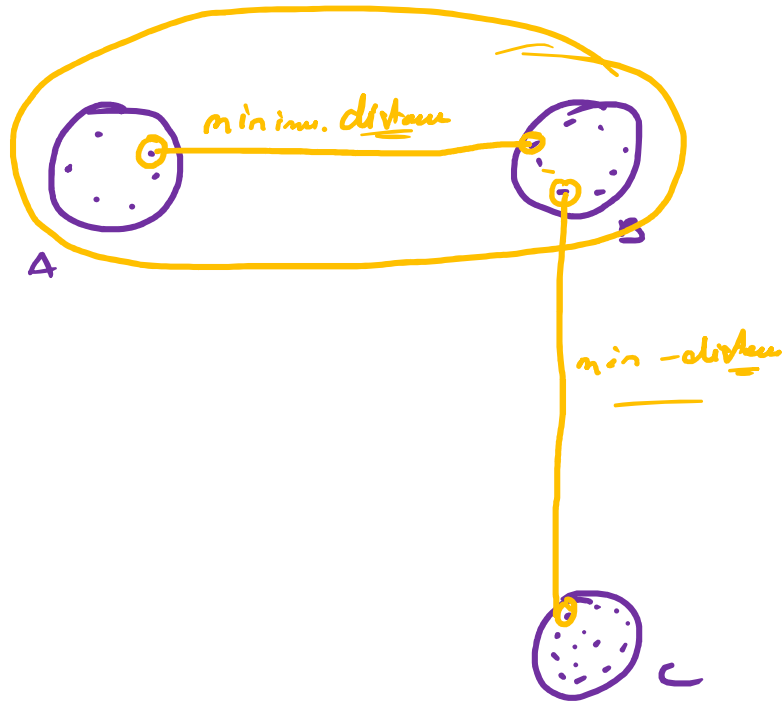
(17)
1002

② HC you can reproduce the same clusters on same data every time

In K-means you cannot ^{can} reproduce the result, because of the random centroid allocation

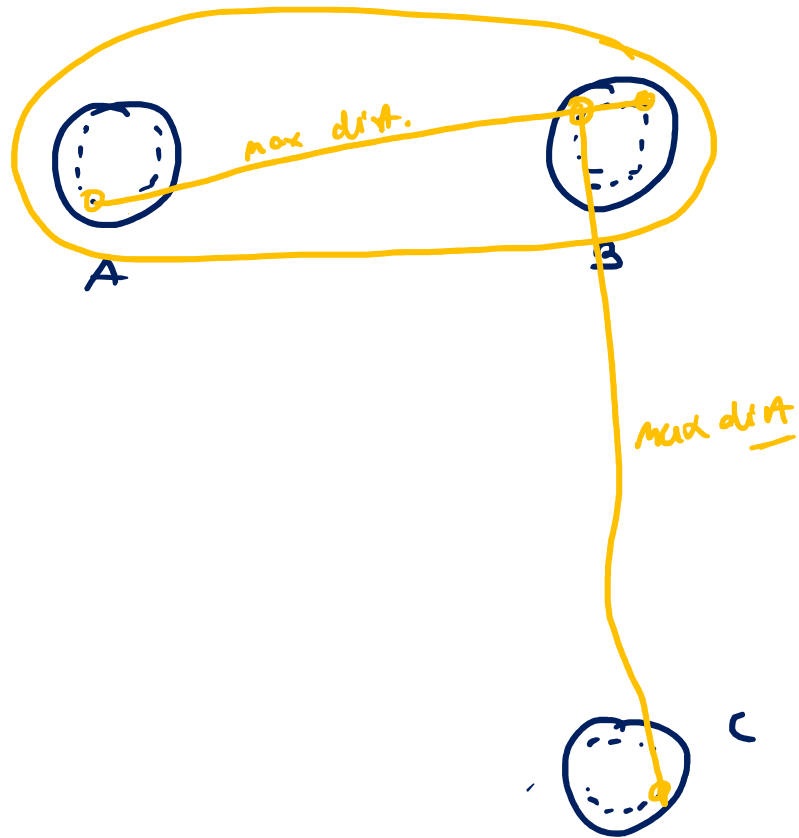
→ Types of linkages

① Single linkage → closest minimum distance



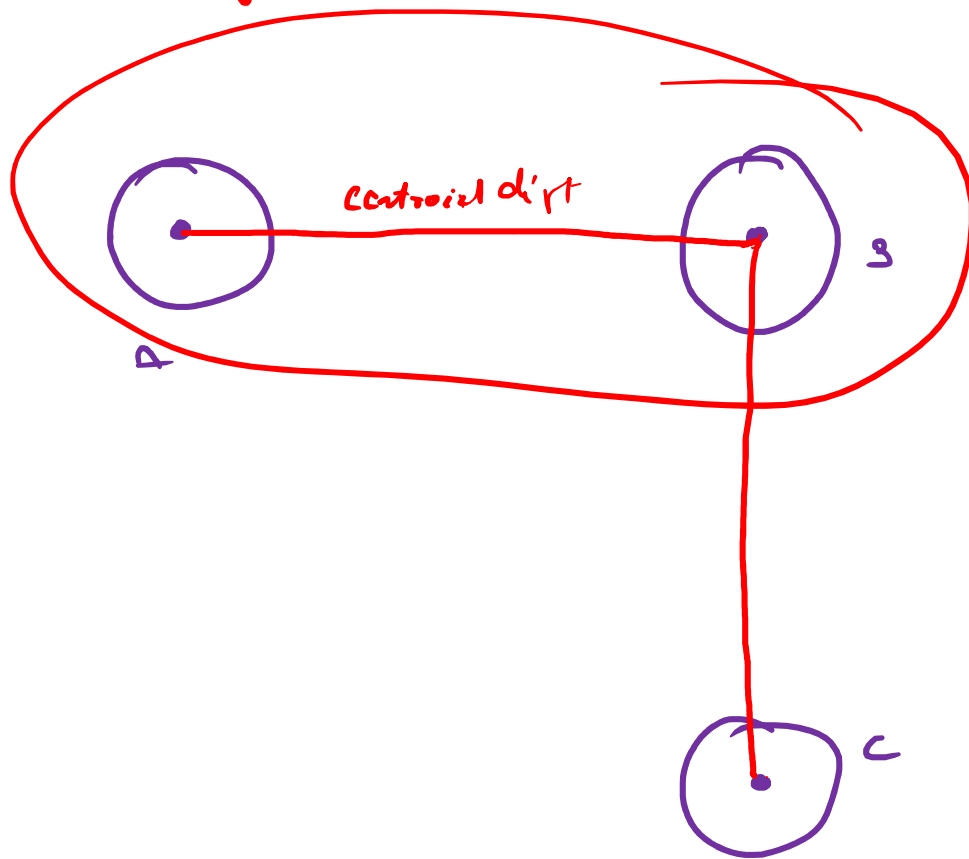
Any one
can be
chosen
randomly
in case of
tie

② Complete linkage → closest Maximum distance



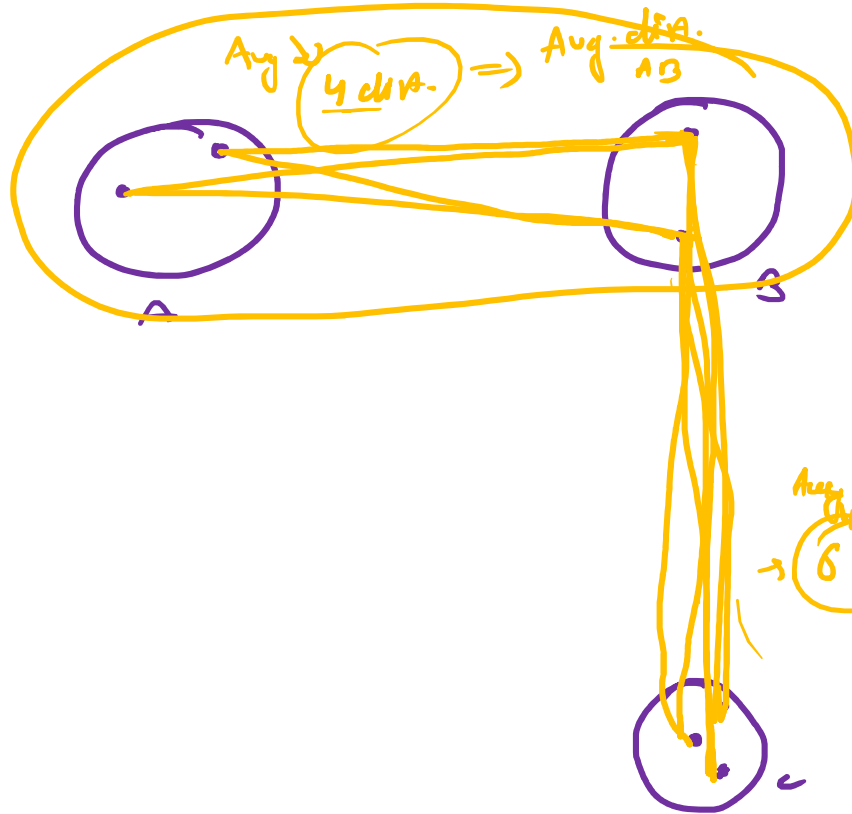
1	2	mean
+	3	-
-	2	-

③ Centroid linkage \rightarrow downward Centroid distance



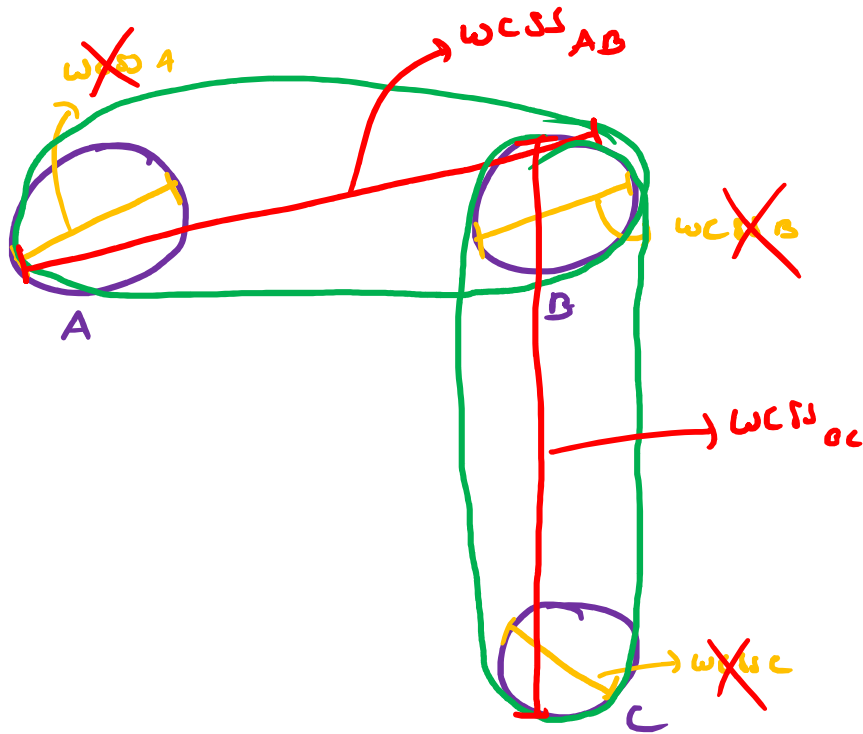
→ Average linkage → down Avg distance

Avg of Every possible distance from all the points of two clusters



	Avg
AB	—
BC	—

⑤ Ward linkage \Rightarrow down wcss/ess value \rightarrow most used



down ($wcss_{AB}, wcss_{BC}$)

group having the down wcss will be grouped

for Example:

for linkage code:

→ Quality of clusters → Silhouette Score

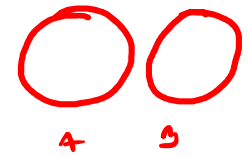
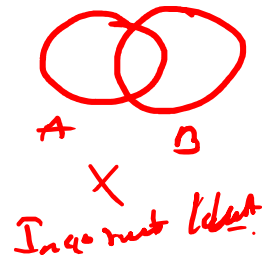
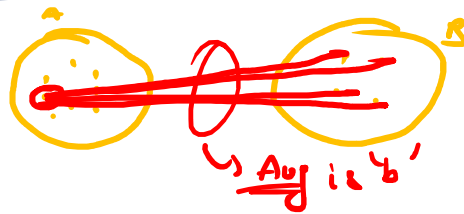
$$S = \frac{b - a}{\max(a, b)}$$

How much closely packed

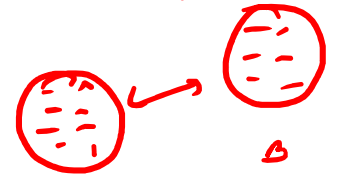
a = Mean distance b/w a randomly selected point of the cluster & all other points in the same cluster



b = Mean dist. b/w the same selected pt. of the cluster & all the points in the nearest neighbour cluster



Bad quality



good quality

Best value of 'a' = 0

$$S = \frac{b - 0}{\max(b, 0)} = \frac{b}{b} = 1$$

Best quality

$$b = a$$

$$S = \frac{b - a}{\max(b, a)} = \frac{0}{\max(b, a)} = 0$$

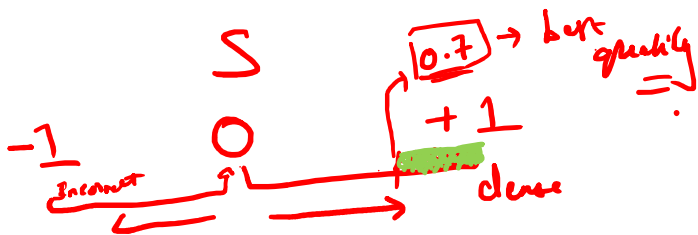
Bad quality

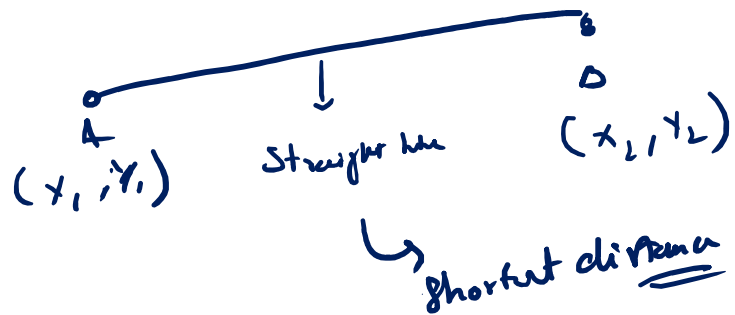
$a > b$

↳ overlapping clusters

$$b = 0$$

$$S = -1 \rightarrow \text{Incorrect cluster}$$





→

$$AB = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

\Downarrow
Euclidean dist.