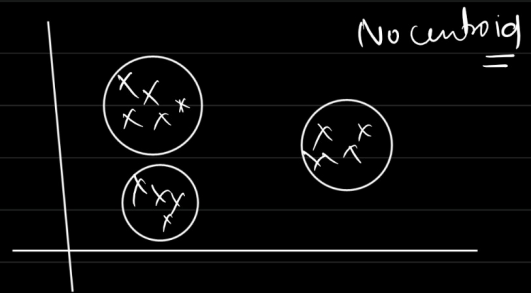
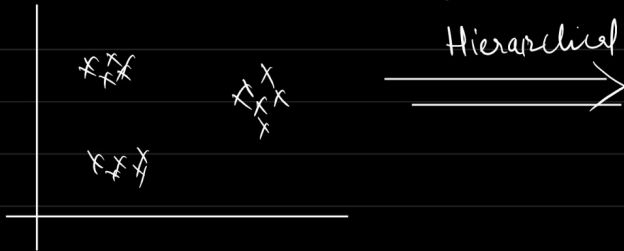


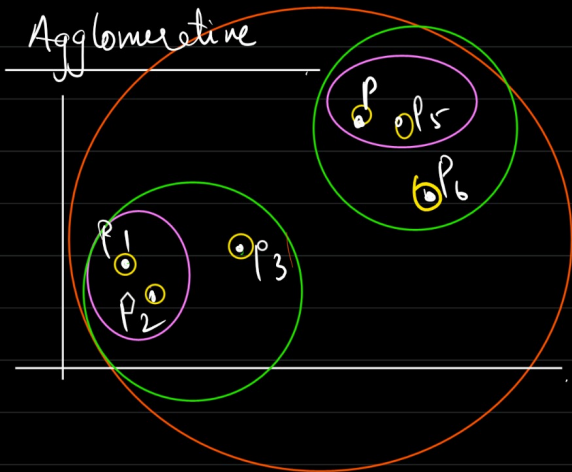
* Hierarchical clustering



* Hierarchical clustering.

- ① Agglomerative Clustering (combining)
- ② Divisive Clustering (dividing)

* Agglomerative



Steps (Agglomerative)

- ① Each point is a cluster in its own
- ② Find the nearest point and create a new cluster.
- ③ Keep on doing step 2 until we get a single cluster.

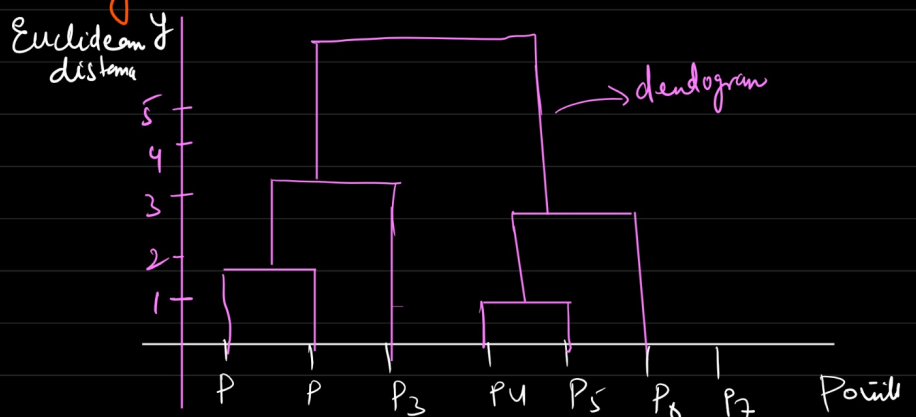
→ Euclidean dist

→ Manhattan distance

→ Cosine Similarity → To calculate distance between two vector (Categorical data / strip variable)

Still how to select K ?

Using dendrogram



To determine K

→ threshold on Euclidean distance. (can be tricky job \Rightarrow business team)

→ trick \rightarrow we can take longest vertical line of dendrogram where none of the horizontal line of dendrogram is passing.

① Using threshold on Euclidean distance

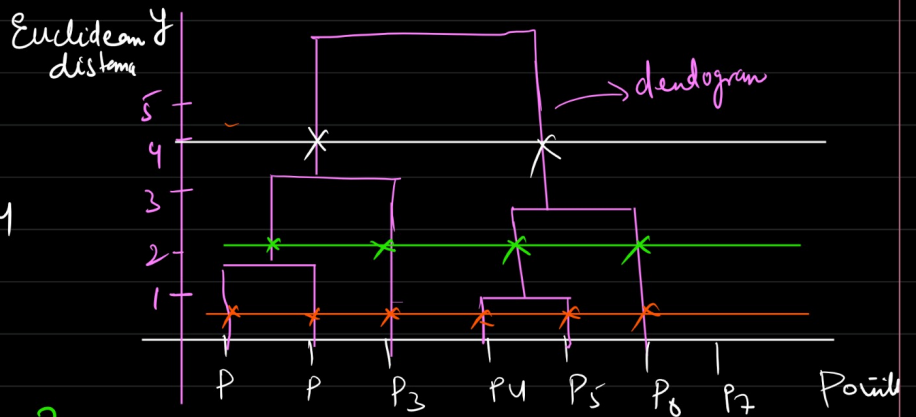
→ Say threshold = 4

$$\Downarrow \\ k=2$$

→ Say threshold is 2
 $k=2$

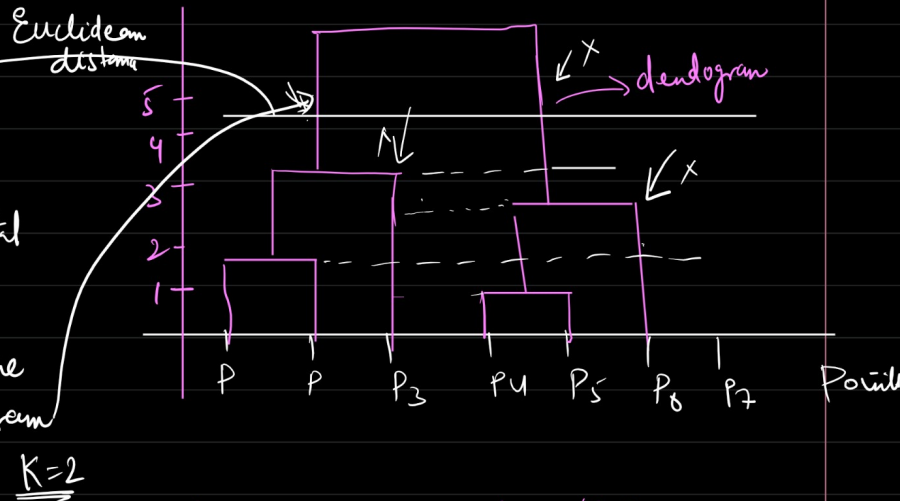
→ say threshold is 1

$k=6$ (all the d.p.'s is cluster in its own)



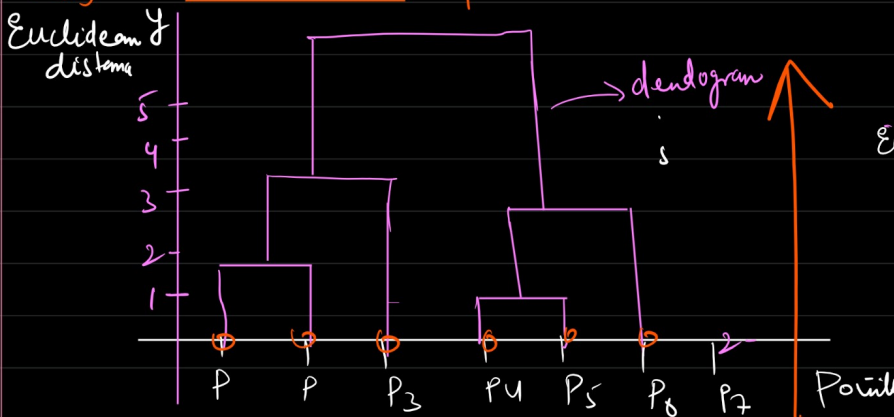
②

Apart from this horizontal line no any other horizontal line is cutting this vertical line of dendrogram



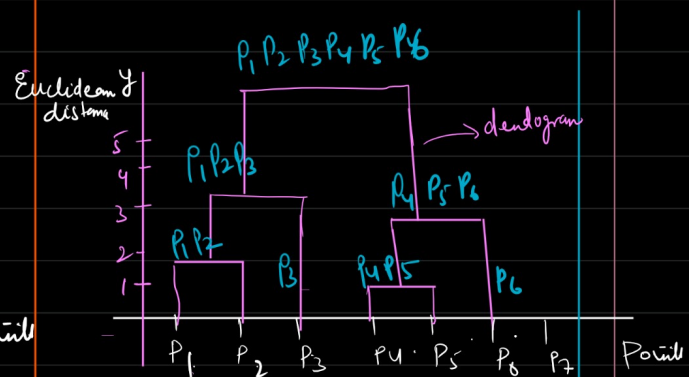
Conclusion → Select the lowest vertical line such that no horizontal line passes through it.

Agglomerative clustering



→ Here every d.p.'s are getting combined

Divisive clustering



Keep diving until every individual d.p. becomes cluster in its own.

Dividing

* K Means vs Hierarchical (scalability & flexibility)

① Size of data → Huge → K Means
→ Small → Hierarchical

② K means — Numerical data (Euclidean | Manhattan distance)
(only)

Hierarchical Clustering — Variety of data (Euclidean, Manhattan, cosine similarity)

③ K means → Centroid(k) → Elbow method

Hierarchical — No centroid.

