# * Hierarchical Clustering *

↳ Unsupervised ML Algo, group unlabelled dataset into a cluster c/a Hierarchical Cluster Analysis.
↳ Develop tree-structured hierarchy of clusters, tree-shaped structure c/a dendrogram.

↳ Num. of clusters not predefined.

↳ 2 approaches:-

① Agglomerative: Bottom-up approach.
↳ All data points consider individual clusters & iteratively merge until 1 cluster left.

② Divisive: Top-down approach.
- Reverse of Agglomerative.
- Start with consider 1 cluster and iteratively divide till reach individual cluster.

↳ working of Hierarchical Clustering :-

① start, consider every data point separate individual cluster.
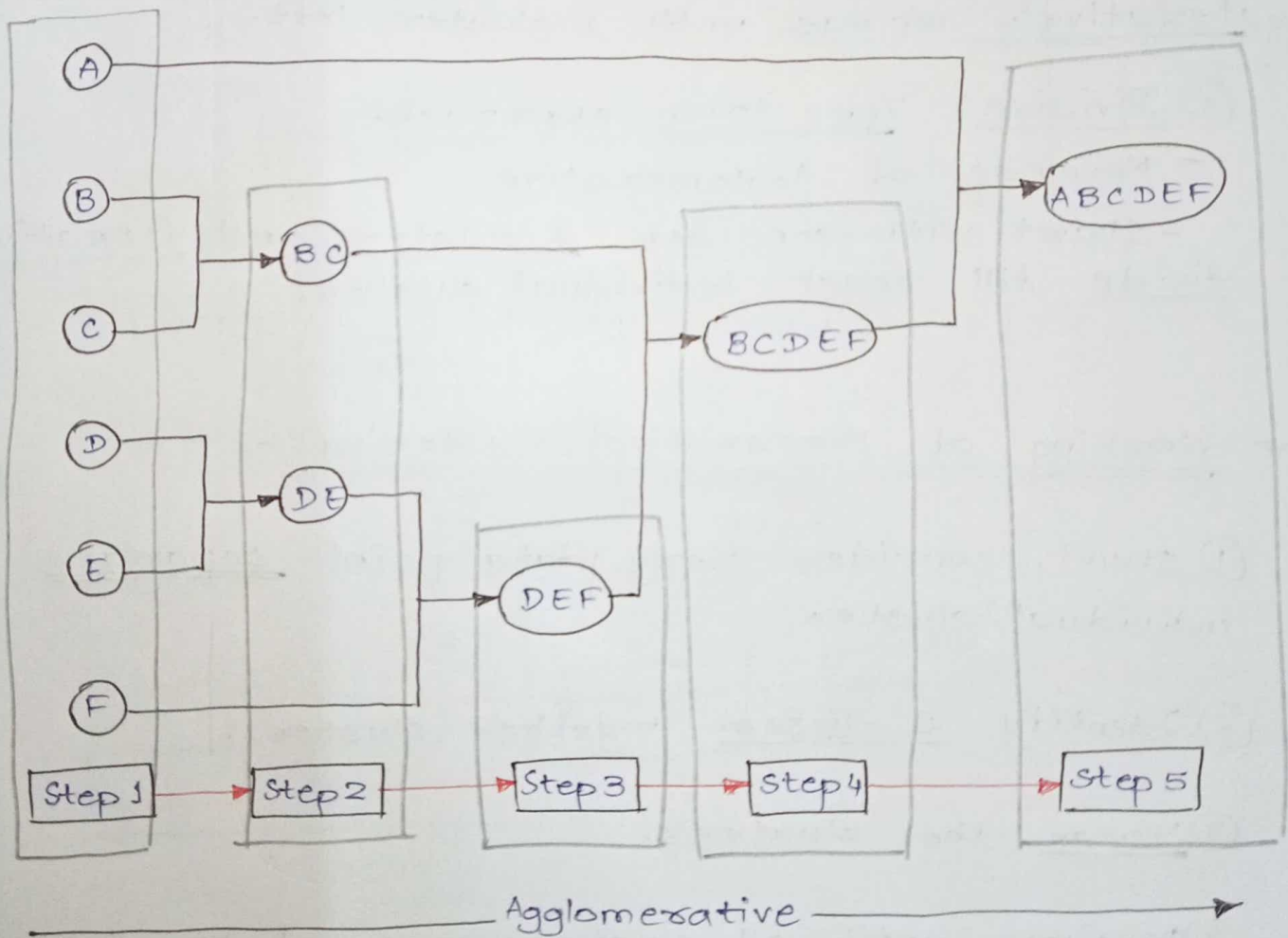
② Identify 2 closest together clusters.

③ Merge the clusters.

④ continue until all clusters merged together.

# 1) Agglomerative Clustering :- Bottom-up method.

① Consider every data point individual cluster.

② Calculate proximity matrix (similarity of cluster with all other clusters).

③ Merge highly similar / closest clusters.

④ Recalculate proximity matrix of each clusters.

⑤ Repeat steps 3 & 4, until only 1 cluster remain.

- Example:



Step 1 → Step 2 → Step 3 → Step 4 → Step 5
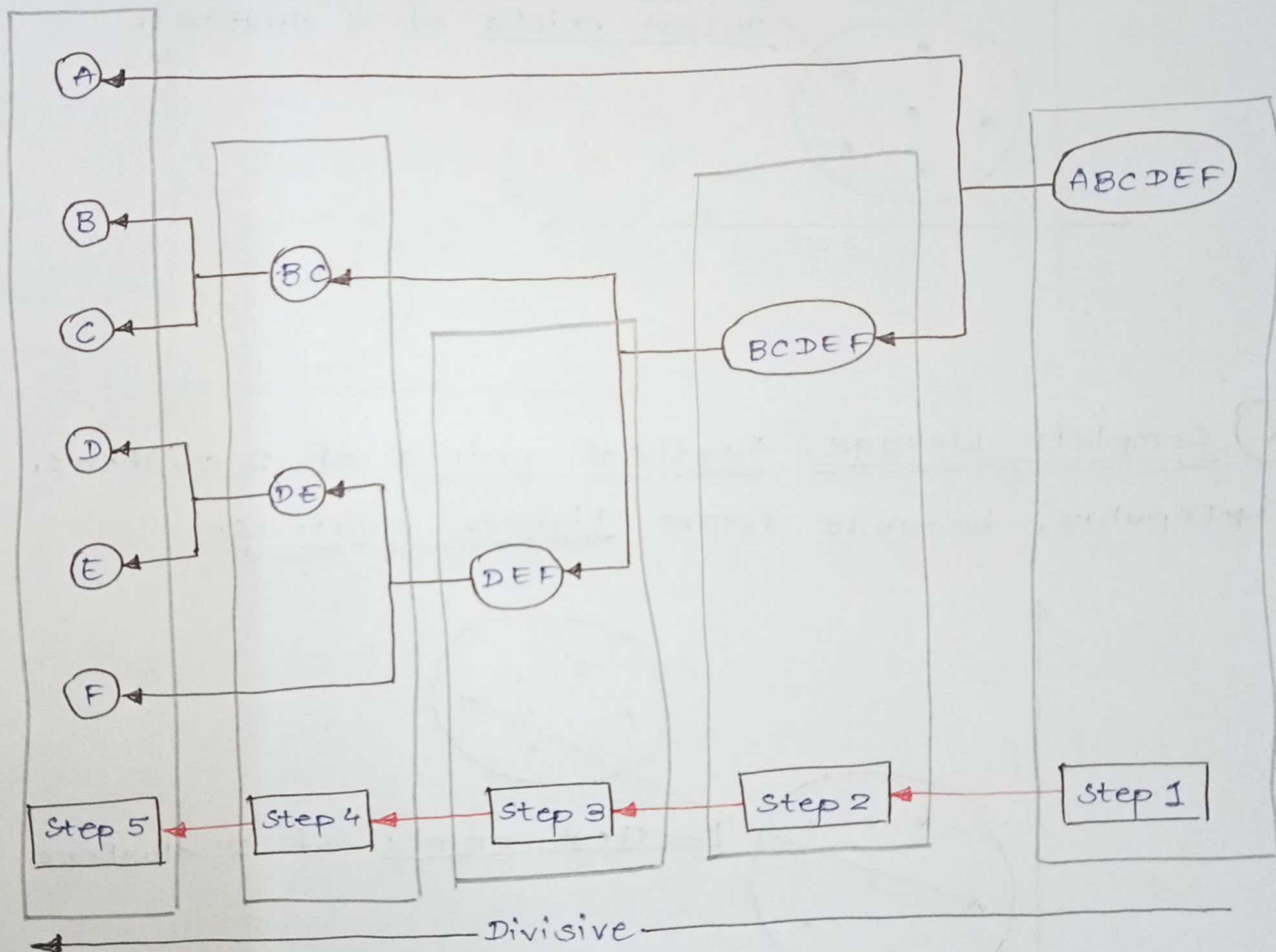
──────────── Agglomerative ────────────→

## 2) Divisive Clustering :- Top-down method.

↳consider all data points a single cluster.

↳In each iteration, separate data points from clusters which aren't comparable.
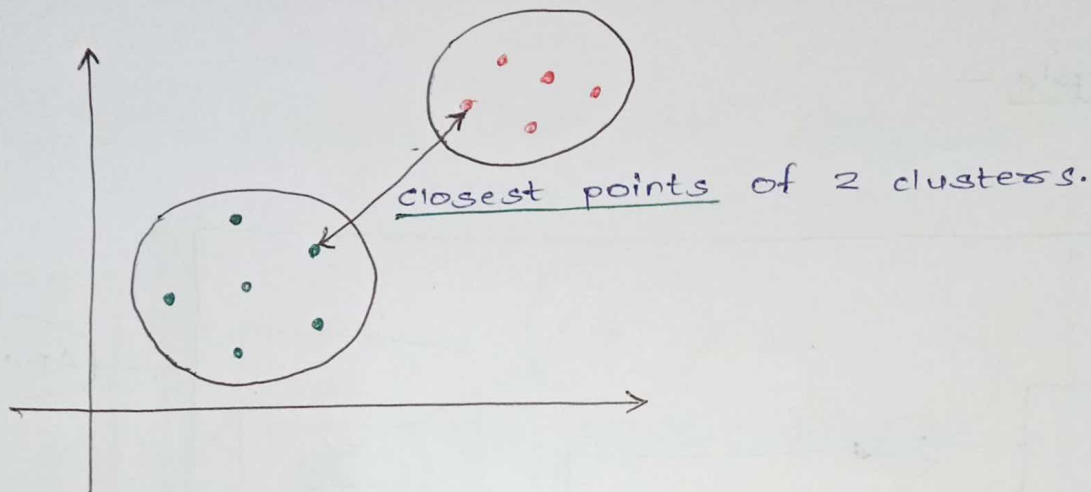
↳At the end get N-clusters.

- Example —



```
                                                      ABCDEF
                                          BCDEF
         BC
                         DEF
         DE


Step 5 ← Step 4 ← Step 3 ← Step 2 ← Step 1
```

Divisive

- <u>Measure Dist. bet$^n$ 2 clusters</u>:-
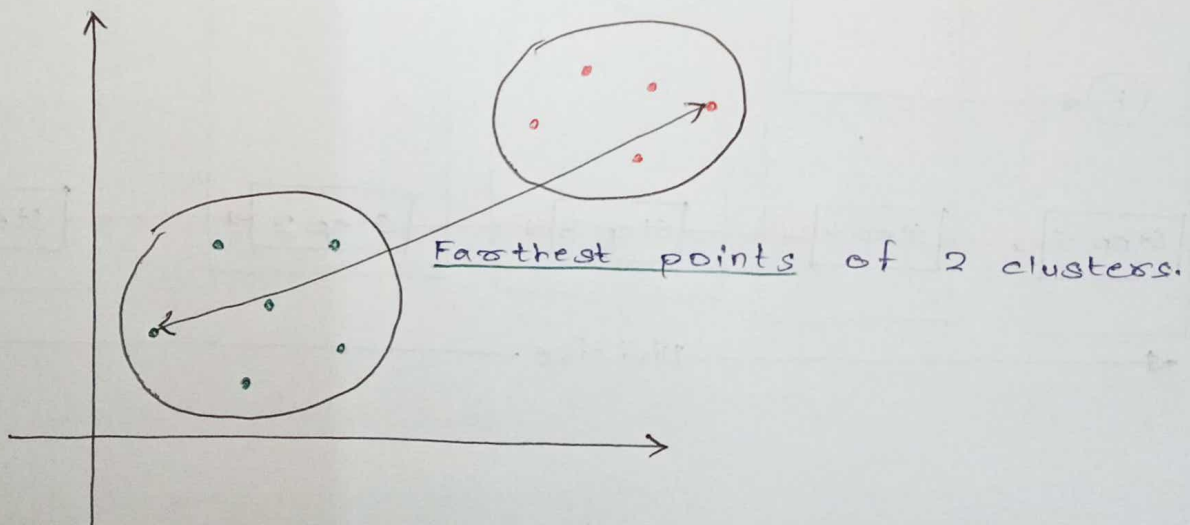
↳Find <u>closest distance</u> bet$^2$ 2 clusters.

↳Different methods c/a <u>Linkage Methods</u>.


① <u>Single Linkage</u>: shortest distance between the <u>closest points</u>.



Closest points of 2 clusters.


② <u>Complete Linkage</u>: <u>Farthest points</u> of 2 clusters.

↳Popular, because form <u>tighter clusters</u>.



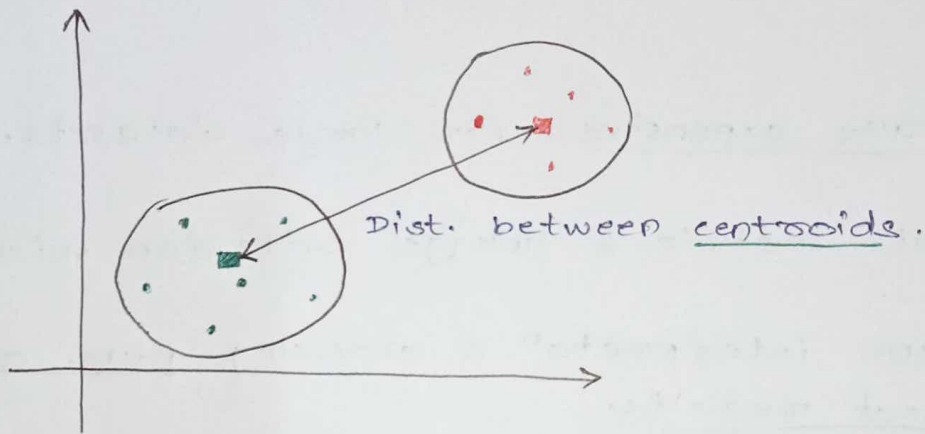<u>Farthest points</u> of 2 clusters.

③ Average Linkage :— calculate average distances.

↳ Add up dist. bet$^n$ each pair of datasets & then divide by total num. of data points.

$$\left(\begin{array}{c}\text{Average}\\ \text{Dist}\end{array}\right) = \frac{\Sigma \text{ dist. bet}^n \text{ each data pair}}{\text{Num. of data points}}$$

④ Centroid Linkage : Dist. bet$^n$ centroids.



Dist. between centroids.

- **Advantages:**

↳ Not required to predefine num. of clusters.

↳ Hierarchical cluster representation useful for multiple level data points.

↳ Work well for different dist. metrics & linkage criterion.

- **Limitations:**

↳ Computationally expensive for large datasets.

↳ Choose dist. metric & linkage criterion wisely.

↳ Dendrogram interpretaⁿ & choosing num. of clusters affect results.