

# EXPERIMENT NO. 8

**Aim:** Implementation of any one Hierarchical Clustering method

**Software Used:** Java/ Python

**Theory:**

A **Hierarchical clustering** method works via grouping data into a tree of clusters. Hierarchical clustering begins by treating every data points as a separate cluster. Then, it repeatedly executes the subsequent steps:

1. Identify the 2 clusters which can be closest together, and
2. Merge the 2 maximum comparable clusters. We need to continue these steps until all the clusters are merged together.

In Hierarchical Clustering, the aim is to produce a hierarchical series of nested clusters. A diagram called **Dendrogram** (A Dendrogram is a tree-like diagram that statistics the sequences of merges or splits) graphically represents this hierarchy and is an inverted tree that describes the order in which factors are merged (bottom-up view) or cluster are break up (top-down view).

There are two types of hierarchical clustering methods:

1. Agglomerative hierarchical clustering:

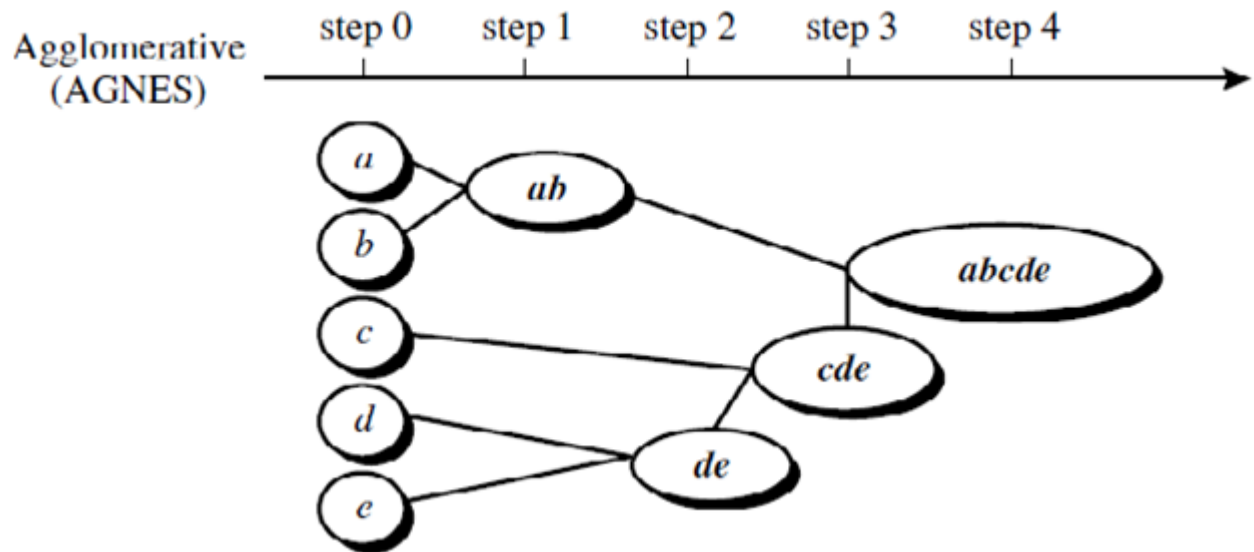
This bottom-up strategy starts by placing each object in its own cluster and then merges these atomic clusters into larger and larger clusters, until all of the objects are in a single cluster or until certain termination conditions are satisfied.

2. Divisive hierarchical clustering:

This top-down strategy does the reverse of agglomerative hierarchical clustering by starting with all objects in one cluster. It subdivides the cluster into smaller and smaller pieces, until each object forms a cluster on its own or until it satisfies certain termination conditions, such as a desired number of clusters is obtained or the diameter of each cluster is within a certain threshold.

**AGGLOMERATIVE HIERARCHICAL CLUSTERING:** - Figure shows the application of AGNES (AGglomerative NESTing), an agglomerative hierarchical clustering method to a data set of five objects(a, b, c, d, e).

- Initially, AGNES places each object into a cluster of its own.
- The clusters are then merged step-by-step according to some criterion.



### Agglomerative Algorithm: (AGNES)

Given

-a set of  $N$  objects to be clustered

-an  $N \times N$  distance matrix ,

The basic process of clustering is this:

**Step1:** Assign each object to a cluster so that for  $N$  objects we have  $N$  clusters each containing just one Object.

**Step2:** Let the distances between the clusters be the same as the distances between the objects they contain.

**Step3:** Find the most similar pair of clusters and merge them into a single cluster so that we now have one cluster less.

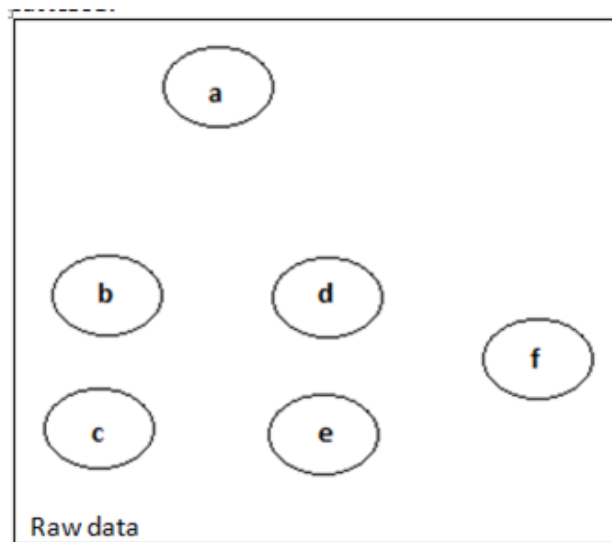
**Step4:** Compute distances between the new cluster and each of the old clusters.

**Step5:** Repeat steps 3 and 4 until all items are clustered into a single cluster of size  $N$ .

- Step 4 can be done in different ways and this distinguishes single and complete linkage.
  - > For complete-linkage algorithm:
    - clustering process is terminated when the maximum distance between nearest clusters exceeds an arbitrary threshold.
  - > For single-linkage algorithm:
    - clustering process is terminated when the minimum distance between nearest clusters exceeds an arbitrary threshold.

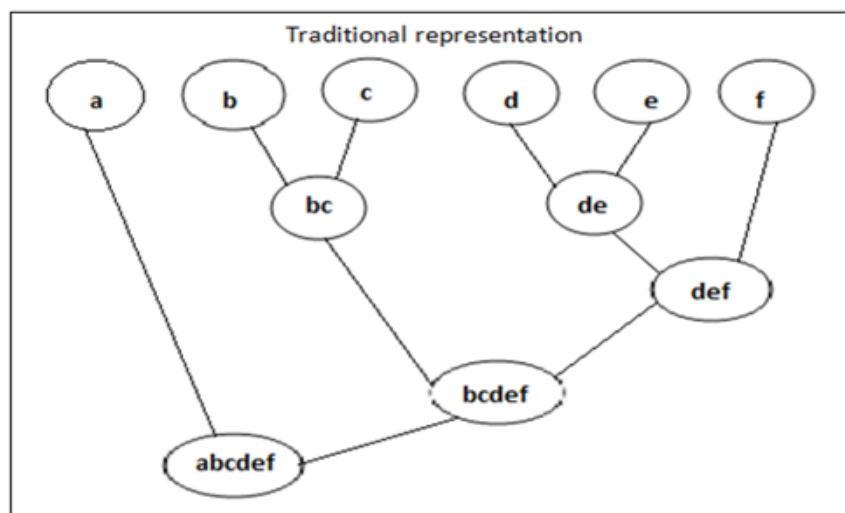
### EXAMPLE:

Suppose this data is to be clustered.



- In this example, cutting the tree after the second row of the dendrogram will yield clusters {a} {b c} {d e} {f}.
- Cutting the tree after the third row will yield clusters {a} {b c} {d e f}, which is a coarser clustering, with a smaller number but larger clusters.

The hierarchical clustering dendrogram would be as such:



In our example, we have six elements {a} {b} {c} {d} {e} and {f}.

The first step is to determine which elements to merge in a cluster.

Usually, we take the two closest elements, according to the chosen distance.

Then, as clustering progresses, rows and columns are merged as the clusters are merged and the distances updated. Suppose we have merged the two closest elements b and c, we now have the following clusters {a}, {b, c}, {d}, {e} and {f}, and want to merge them further.

To do that, we need to take the distance between {a} and {b c}, and therefore define the distance between two clusters. Usually the distance between two clusters A and B is one of the following:

- The maximum distance between elements of each cluster (also called complete-linkage clustering):  

$$\max \{d(x,y):x \in A,y \in B\}$$
- The minimum distance between elements of each cluster (also called single-linkage clustering):

$$\min \{d(x,y):x\in A,y\in B\}$$

- The mean distance between elements of each cluster (also called average linkage clustering):

$$\frac{1}{|\mathcal{A}| \cdot |\mathcal{B}|} \sum_{x \in \mathcal{A}} \sum_{y \in \mathcal{B}} d(x, y)$$

Each agglomeration occurs at a greater distance between clusters than the previous agglomeration, and one can decide to stop clustering either when the clusters are too far apart to be merged (distance criterion) or when there is a sufficiently small number of clusters (number criterion).

## PROGRAM:

```
[1] import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
```

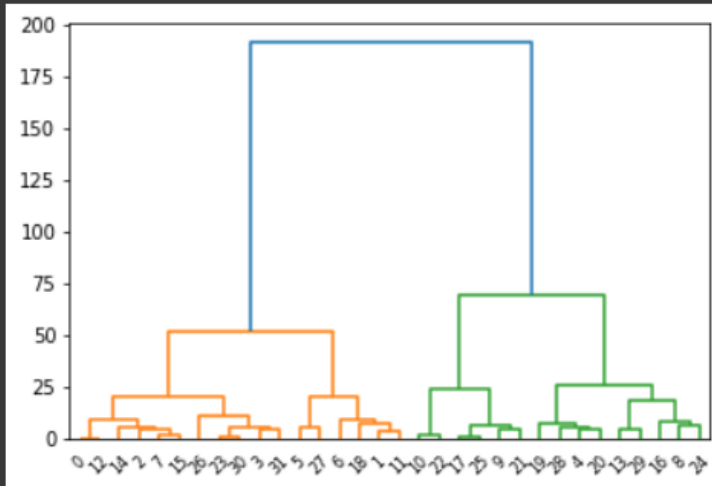
```
[2] from google.colab import files
uploaded = files.upload()
```

```
import io
df = pd.read_csv(io.BytesIO(uploaded['index_dim.csv']))
print(df)
```

1	BEN2	1	28	66	1007	20
2	BEN3	1	23	79	1012	21
3	BEN4	1	20	91	1012	9
4	BOM1	7	28	42	1018	11
5	BOM2	8	34	63	1011	7
6	BOM3	6	26	76	1004	9
7	BOM4	6	32	82	1011	14
8	DEL1	1	20	30	1020	6
9	DEL2	1	29	16	1006	11
10	DEL3	1	40	24	992	19
11	DEL4	1	30	70	1006	11
12	HYD1	1	20	83	1016	15
13	HYD2	1	34	38	1003	21
14	HYD3	1	26	82	1007	26
15	HYD4	1	25	84	1011	11
16	JAI1	1	19	37	1016	12
17	JAI2	1	36	15	1001	8
18	JAI3	1	35	68	1000	8
19	JAI4	1	18	50	1016	14
20	KAN1	1	17	46	1020	7
21	KAN2	1	28	12	1004	15
22	KAN3	1	40	26	992	15
23	KAN4	1	27	89	1007	7
24	NAG1	1	22	30	1013	13
25	NAG2	1	40	16	1000	11
26	NAG3	1	29	89	1000	12
27	NAG4	1	23	58	1014	6
28	PUN1	1	23	44	1014	6

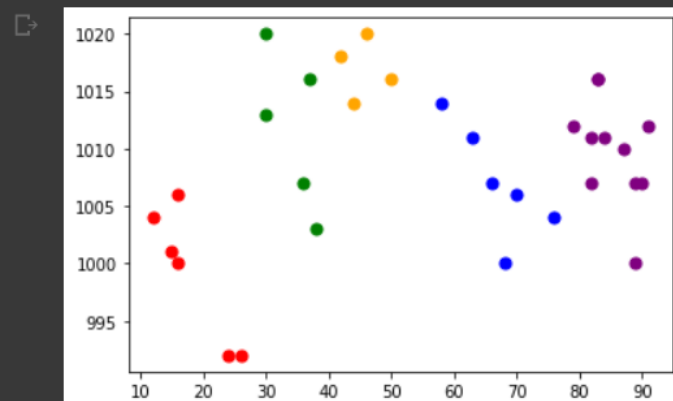
```
[4] import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
from sklearn.cluster import AgglomerativeClustering
import scipy.cluster.hierarchy as sch
```

```
[5] Adataset = pd.read_csv('index_dim.csv')
X = dataset.iloc[:, [3, 4]].values
dendrogram = sch.dendrogram(sch.linkage(X, method='ward'))
```



```
[6] model = AgglomerativeClustering(n_clusters=5, affinity='euclidean', linkage='ward')
model.fit(X)
labels = model.labels_
```

```
[7] plt.scatter(X[labels==0, 0], X[labels==0, 1], s=50, marker='o', color='red')
plt.scatter(X[labels==1, 0], X[labels==1, 1], s=50, marker='o', color='blue')
plt.scatter(X[labels==2, 0], X[labels==2, 1], s=50, marker='o', color='green')
plt.scatter(X[labels==3, 0], X[labels==3, 1], s=50, marker='o', color='purple')
plt.scatter(X[labels==4, 0], X[labels==4, 1], s=50, marker='o', color='orange')
plt.show()
```



**CONCLUSION:**

The different clustering algorithms of data mining were studied and one among them named Agglomerative clustering algorithm was implemented using Python. The need for clustering algorithm was recognized and understood.

**SIGN AND REMARK**

<b>R1</b> <b>(3 M)</b>	<b>R2</b> <b>(3 M)</b>	<b>R3</b> <b>(3 M)</b>	<b>R4</b> <b>(3 M)</b>	<b>R5</b> <b>(3 M)</b>	<b>Total</b>	<b>Sign</b>

**DATE**