

Name :- Aniket Kumar

Sn.No :- 30

Roll No :- PF45

PRN No :- 1032171203.

AML LAB ASSIGNMENT 06

Problem Statement :- Implementation and Comparison of various clustering techniques: k-means, Hierarchical.

Theory :-

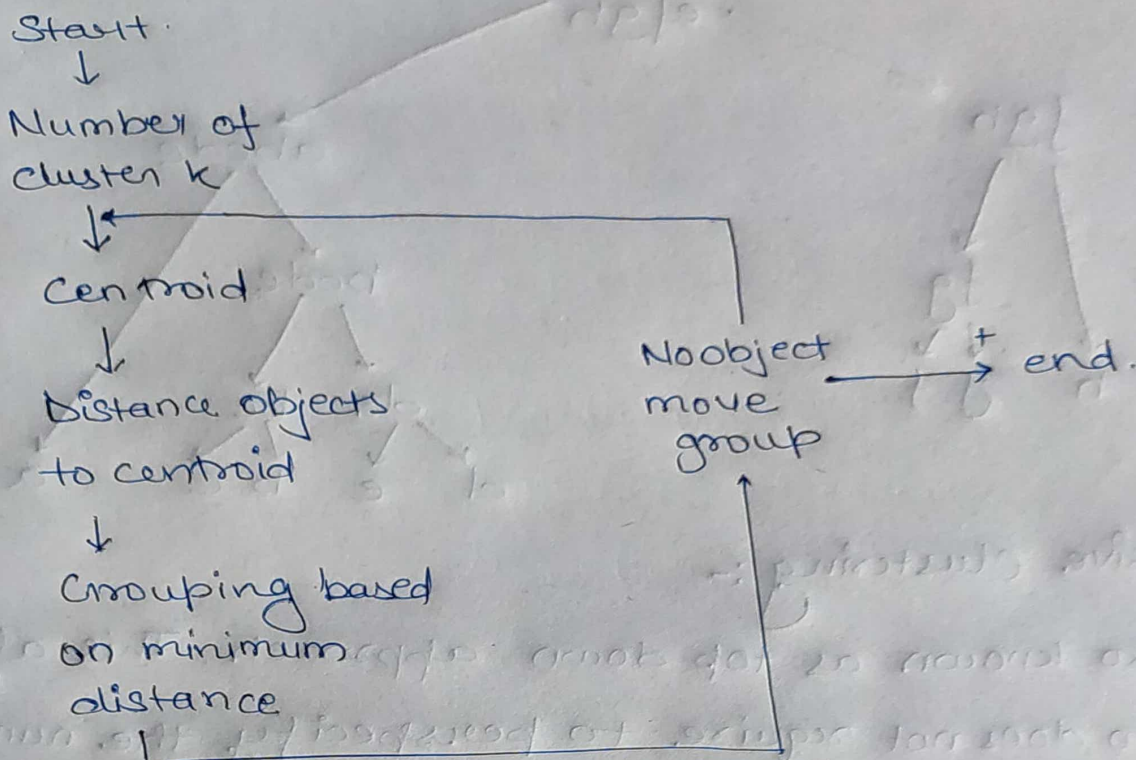
Clustering Algorithm.

Clustering is very much important as it determines the intrinsic grouping among the unlabelled data present.

k-Means Clustering :-

It is the simplest supervised learning algorithm that solves clustering problem. The algorithm will categorize the items into k groups of similarity.

flow diagram.



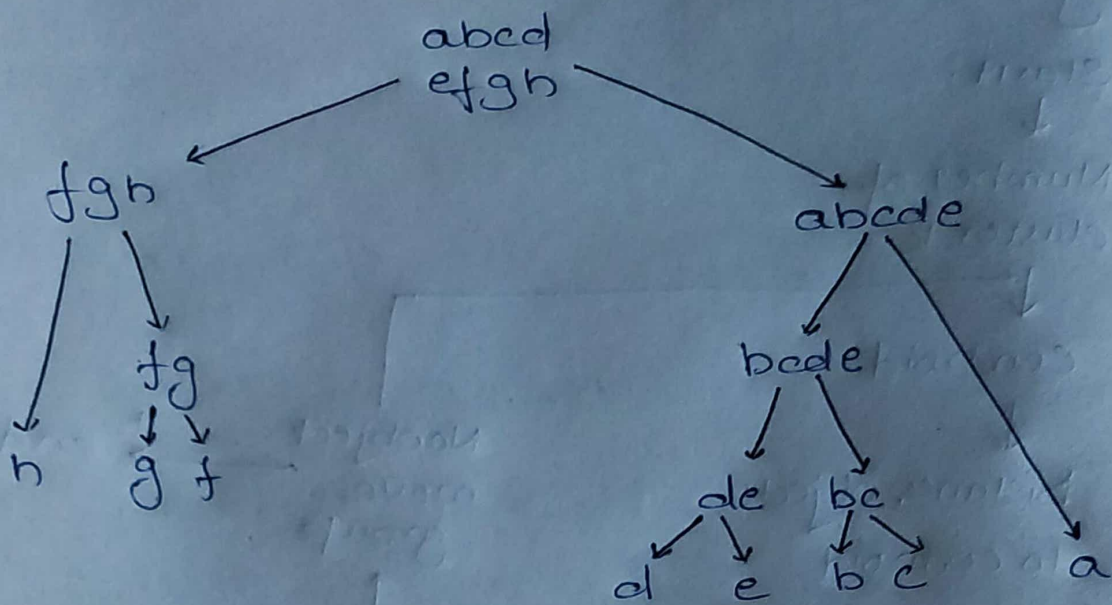
• Euclidean Distance.

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

Hierarchical clustering.

It is a method of cluster analysis which seek to build a hierarchy clusters. i.e. tree type structure based on the hierarchy.

Agglomerative Hierarchical clustering: Also known as bottom-up approach or hierarchical agglomerative clustering.



Divisive clustering :-

Also known as top down approach. This algorithm also does not require to pre-specify the number of clusters. It requires a method for splitting a cluster that contains the whole data and proceeds by splitting clusters recursively until individual data have been split into singleton cluster.

Conclusion :-

Hence, we implemented different clustering algorithms.

Implementation of Kmeans, Hierarchical, DBScan and Spectral Clustering

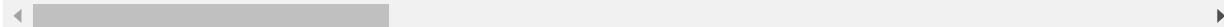
```
In [1]: import numpy as np
import pandas as pd
```

```
In [2]: data = pd.read_csv("cluster.csv")
data.head()
```

Out[2]:

| | id | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean | smoothness_me |
|---|----------|-----------|-------------|--------------|----------------|-----------|---------------|
| 0 | 842302 | M | 17.99 | 10.38 | 122.80 | 1001.0 | 0.118 |
| 1 | 842517 | M | 20.57 | 17.77 | 132.90 | 1326.0 | 0.084 |
| 2 | 84300903 | M | 19.69 | 21.25 | 130.00 | 1203.0 | 0.109 |
| 3 | 84348301 | M | 11.42 | 20.38 | 77.58 | 386.1 | 0.142 |
| 4 | 84358402 | M | 20.29 | 14.34 | 135.10 | 1297.0 | 0.100 |

5 rows × 33 columns



```
In [3]: from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
```

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 33 columns):
id                    569 non-null int64
diagnosis             569 non-null object
```

```

radius_mean          569 non-null float64
texture_mean         569 non-null float64
perimeter_mean       569 non-null float64
area_mean            569 non-null float64
smoothness_mean      569 non-null float64
compactness_mean     569 non-null float64
concavity_mean       569 non-null float64
concave points_mean  569 non-null float64
symmetry_mean        569 non-null float64
fractal_dimension_mean 569 non-null float64
radius_se            569 non-null float64
texture_se           569 non-null float64
perimeter_se         569 non-null float64
area_se              569 non-null float64
smoothness_se        569 non-null float64
compactness_se       569 non-null float64
concavity_se         569 non-null float64
concave points_se    569 non-null float64
symmetry_se          569 non-null float64
fractal_dimension_se 569 non-null float64
radius_worst         569 non-null float64
texture_worst        569 non-null float64
perimeter_worst      569 non-null float64
area_worst           569 non-null float64
smoothness_worst     569 non-null float64
compactness_worst    569 non-null float64
concavity_worst      569 non-null float64
concave points_worst 569 non-null float64
symmetry_worst       569 non-null float64
fractal_dimension_worst 569 non-null float64
Unnamed: 32          0 non-null float64
dtypes: float64(31), int64(1), object(1)
memory usage: 146.8+ KB

```

```

In [4]: cols_drop = ['id', 'Unnamed: 32']
data = data.drop(cols_drop, axis=1)
data['diagonis'] = data['diagnosis'].map({'M':1, 'B':0})
X = data.drop('diagnosis', axis=1).values
X = StandardScaler().fit_transform(X)

```

```
In [5]: #1 KMeans Clustering
from sklearn.cluster import KMeans
km = KMeans(n_clusters=2, init="k-means++", n_init=8)
km_pred = km.fit_predict(X)

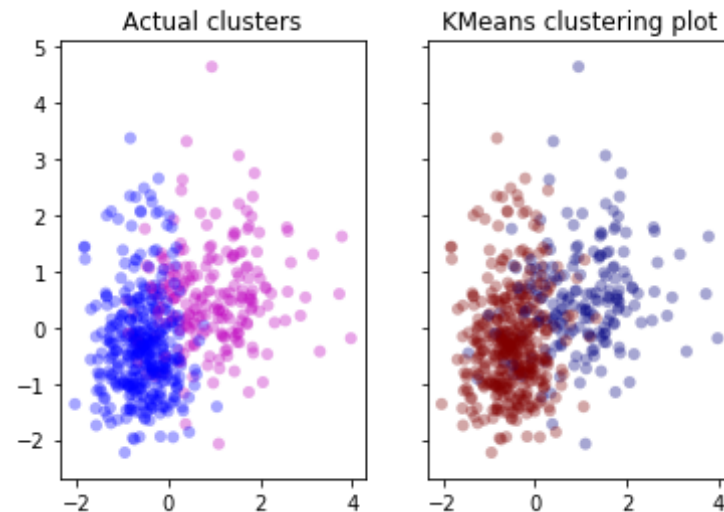
f, (ax1, ax2) = plt.subplots(1, 2, sharey=True)

ax1.scatter(X[:,0], X[:,1], c=data["diagnosis"], cmap="jet", edgecolor=
"None", alpha=0.35)
ax1.set_title("Actual clusters")

ax2.scatter(X[:,0], X[:,1], c=km_pred, cmap="jet", edgecolor="None", al
pha=0.35)
ax2.set_title("KMeans clustering plot")
```

C:\Users\admin\Anaconda3\lib\site-packages\ipykernel_launcher.py:8: MatplotlibDeprecationWarning: Support for uppercase single-letter colors is deprecated since Matplotlib 3.1 and will be removed in 3.3; please use lowercase instead.

Out[5]: Text(0.5, 1.0, 'KMeans clustering plot')



```

In [6]: #2 Hierarchical Agglomerative Clustering

from sklearn.cluster import AgglomerativeClustering
ac = AgglomerativeClustering(n_clusters=2, linkage="ward")
ac_pred = ac.fit_predict(X)

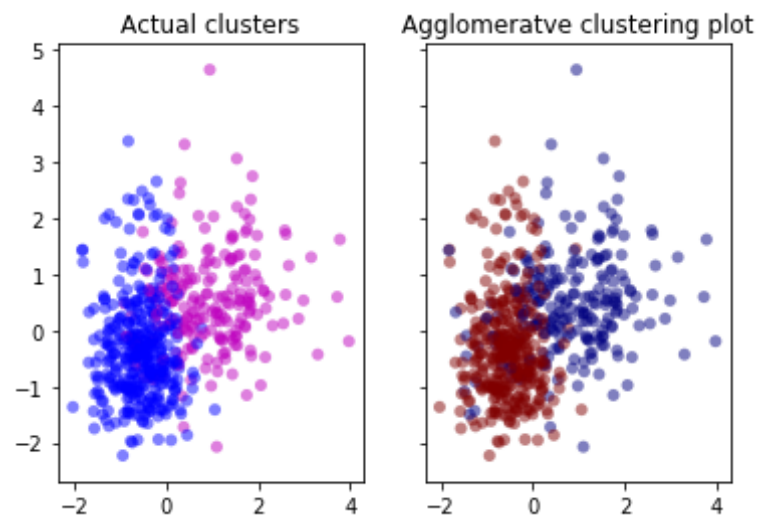
# Scatter plots
f, (ax1, ax2) = plt.subplots(1, 2, sharey=True)

ax1.scatter(X[:,0], X[:,1], c=data["diagnosis"], cmap="jet", edgecolor=
"None", alpha=0.5)
ax1.set_title("Actual clusters")

ax2.scatter(X[:,0], X[:,1], c=ac_pred, cmap="jet", edgecolor="None", al
pha=0.5)
ax2.set_title("Agglomerative clustering plot")

```

Out[6]: Text(0.5, 1.0, 'Agglomerative clustering plot')



```

In [7]: #3 DBSCAN (Density-Based Clustering of Applications with Noise)

from sklearn.cluster import DBSCAN
dbs = DBSCAN(eps=0.2, min_samples=6)

```



```

dbs_pred = dbs.fit_predict(X)

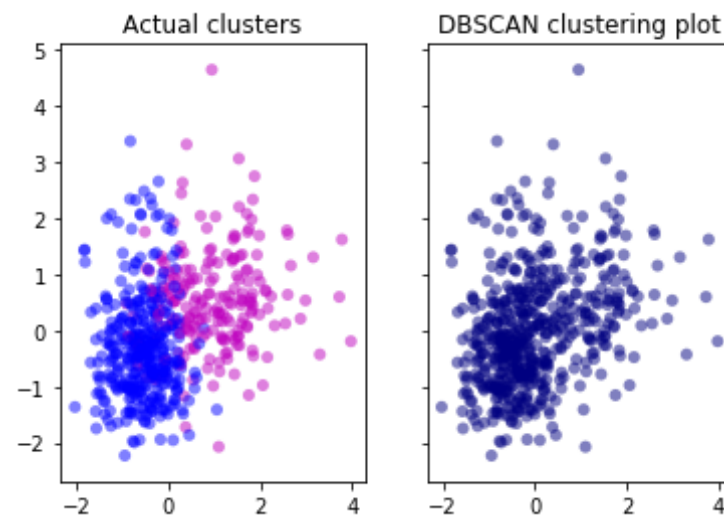
# Scatter plots
f, (ax1, ax2) = plt.subplots(1, 2, sharey=True)

ax1.scatter(X[:,0], X[:,1], c=data["diagnosis"], cmap="jet", edgecolor=
"None", alpha=0.5)
ax1.set_title("Actual clusters")

ax2.scatter(X[:,0], X[:,1], c=dbs_pred, cmap="jet", edgecolor="None", a
lpha=0.5)
ax2.set_title("DBSCAN clustering plot")

```

Out[7]: Text(0.5, 1.0, 'DBSCAN clustering plot')



```

In [8]: #4 Spectral Clustering

from sklearn.cluster import SpectralClustering
sc = SpectralClustering(n_clusters=2, gamma=0.5, affinity="rbf", assign
_labels="discretize")
sc_pred = sc.fit_predict(X)

# Scatter plots

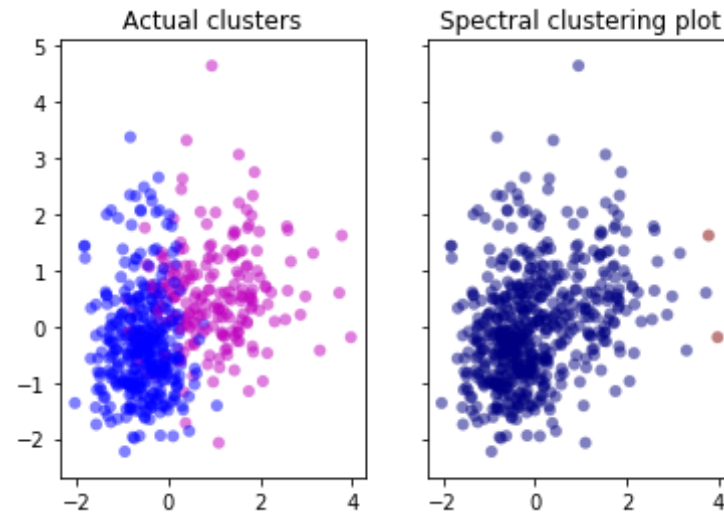
```

```
f, (ax1, ax2) = plt.subplots(1, 2, sharey=True)

ax1.scatter(X[:,0], X[:,1], c=data["diagnosis"], cmap="jet", edgecolor=
"None", alpha=0.5)
ax1.set_title("Actual clusters")

ax2.scatter(X[:,0], X[:,1], c=sc_pred, cmap="jet", edgecolor="None", al
pha=0.5)
ax2.set_title("Spectral clustering plot")
```

Out[8]: Text(0.5, 1.0, 'Spectral clustering plot')



In []: