

Assignment-05**D.Sai Kiran****21BCE9166**▼ **About**

Unsupervised ML algorithms 1) Can work with unlabelled data 2) We do not specify the output 3) Clusters and associates based on features

→ KMeans Clustering

1. It is a type of unsupervised ML algorithm
2. Clusters the data points based on the similarities in the features provided
3. k=No. of clusters

→ How does KMeans work?

1. Select the value of K (can be random / or based on the elbow method) → Elbow method

- a) plot a graph between no. of clusters and wcss
- b) Choose the point where there is an elbow transition and having lower wcss value
- c) wcss=within cluster sum of squares
summation of squares of distances between the points of each cluster to the centroid of the cluster

2. Random instantiation of centroids for each cluster/Using kMeans++ to reduce error
3. Assign each datapoint to the cluster where the distance b/w the centroid of the cluster to the data point is minimal
4. Calculate the centroid of each cluster, considering the members of the cluster
5. Reassign the datapoints as in step 3
6. If any changes made in clustering, recalculate the centroids, step 4 and move to step 5.
7. No changes, model ready

→ Working on data

1) Import the needed libraries 2) Collect the dataset and import it 3) Data preprocessing 4) Extract the needed features 5) Find k value 6) Build the model with found k 7) predict the cluster

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

```
data=pd.read_csv("Mall_Customers.csv")
```

```
data
```

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)	
0	1	Male	19	15	39	
1	2	Male	21	15	81	
2	3	Female	20	16	6	
3	4	Female	23	16	77	
4	5	Female	31	17	40	
...	
195	196	Female	35	120	79	
196	197	Female	45	126	28	
197	198	Male	32	126	74	
198	199	Male	32	137	18	
199	200	Male	30	137	83	

200 rows × 5 columns

```
data.head()
```

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)	
0	1	Male	19	15	39	
1	2	Male	21	15	81	
2	3	Female	20	16	6	
3	4	Female	23	16	77	
4	5	Female	31	17	40	

```
data.tail()
```

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)	
195	196	Female	35	120	79	
196	197	Female	45	126	28	
197	198	Male	32	126	74	
198	199	Male	32	137	18	
199	200	Male	30	137	83	

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   CustomerID            200 non-null   int64
1   Genre                 200 non-null   object
2   Age                  200 non-null   int64
3   Annual Income (k$)    200 non-null   int64
4   Spending Score (1-100) 200 non-null   int64
dtypes: int64(4), object(1)
memory usage: 7.9+ KB
```

```
data.describe()
```

	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)
count	200.000000	200.000000	200.000000	200.000000
mean	100.500000	38.850000	60.560000	50.200000
std	57.879185	13.969007	26.264721	25.823522
min	1.000000	18.000000	15.000000	1.000000
25%	50.750000	28.750000	41.500000	34.750000
50%	100.500000	36.000000	61.500000	50.000000
75%	150.250000	49.000000	78.000000	73.000000
max	200.000000	70.000000	137.000000	99.000000

```
needed_features=data.iloc[:,3:]
```

```
needed_features.head()
```

	Annual Income (k\$)	Spending Score (1-100)
0	15	39
1	15	81
2	16	6
3	16	77
4	17	40

```
x=needed_features.values
```

```
plt.scatter(x[:,0],x[:,1])
plt.title("income VS credit")
plt.xlabel("Income")
plt.ylabel("Credit score")
plt.show()
```



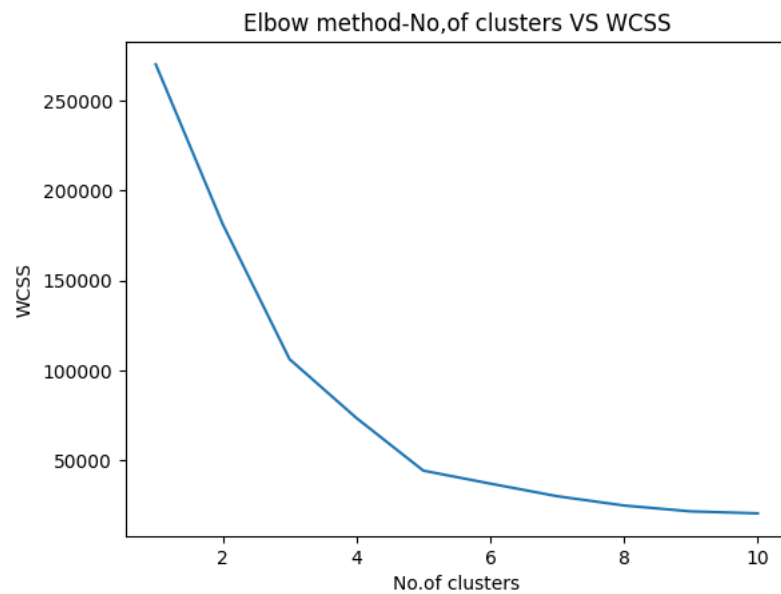
```
for i in range(1,11):
    kmeans=KMeans(n_clusters=i,init="k-means++",random_state=0)
    kmeans.fit(x)
    wcss.append(kmeans.inertia_)
```

[illegible]

4/10

```
[269981.28,
181363.59595959593,
106348.37306211122,
73679.78903948836,
44448.4554479337,
37265.86520484346,
30259.65720728547,
25095.70320999756,
21830.041978049434,
20736.679938924128]
```

```
plt.plot(range(1,11),wcss)
plt.title("Elbow method-No,of clusters VS WCSS")
plt.xlabel("No.of clusters")
plt.ylabel("WCSS")
plt.show()
```



Optimal k value=5

```
kmeans=KMeans(n_clusters=5,init="k-means++",random_state=0)
```

```
result=kmeans.fit_predict(x)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to silence this warning.
warnings.warn(
```

```
result
```

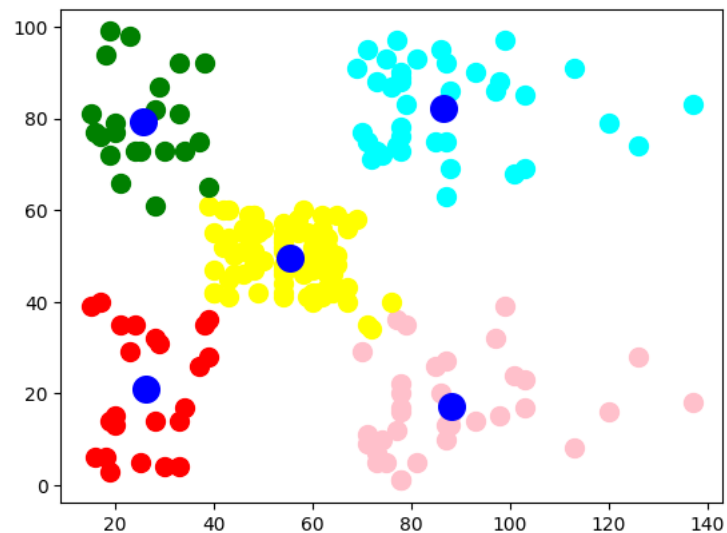
```
array([4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3,
       4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 3, 4, 1,
       4, 3, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 0, 2, 1, 2, 0, 2, 0, 2,
       1, 2, 0, 2, 0, 2, 0, 2, 0, 2, 1, 2, 0, 2, 0, 2, 0, 2, 0, 2,
       0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2,
       0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2,
       0, 2], dtype=int32)
```

```
kmeans.cluster_centers_
```

```
array([[88.2      , 17.11428571],
       [55.2962963 , 49.51851852],
       [86.53846154, 82.12820513],
       [25.72727273, 79.36363636],
       [26.30434783, 20.91304348]])
```

```
plt.scatter(x[result==0,0],x[result==0,1],s=100,color='pink',label='cluster1')
plt.scatter(x[result==1,0],x[result==1,1],s=100,color='yellow',label='cluster2')
plt.scatter(x[result==2,0],x[result==2,1],s=100,color='cyan',label='cluster3')
plt.scatter(x[result==3,0],x[result==3,1],s=100,color='green',label='cluster4')
plt.scatter(x[result==4,0],x[result==4,1],s=100,color='red',label='cluster5')
plt.scatter(kmeans.cluster_centers_[0,0],kmeans.cluster_centers_[0,1],s=200,color='blue',label="centroids")
```



```
<matplotlib.collections.PathCollection at 0x7f5d9e77f850>
```



--> Clustering Considering More than 2 features

```
x=data.iloc[:,1:]
```

```
x
```

	Genre	Age	Annual Income (k\$)	Spending Score (1-100)	
0	Male	19	15	39	
1	Male	21	15	81	
2	Female	20	16	6	
3	Female	23	16	77	
4	Female	31	17	40	
...	
195	Female	35	120	79	
196	Female	45	126	28	
197	Male	32	126	74	
198	Male	32	137	18	
199	Male	30	137	83	

200 rows × 4 columns

label Encoding Gender

```
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
val=x.iloc[:,0]
le.fit(val)
print("Classes = ",le.classes_)
print("After transformation = ",le.transform(val))
```

```
Classes = ['Female' 'Male']
After transformation = [1 1 0 0 0 0 0 0 1 0 1 0 0 0 1 1 0 1 1 0 1 1 0 1 0 1 0 1 0 0 1 0 1 0 0 0 0
0 0 0 0 1 1 0 0 0 0 0 0 0 0 1 0 1 0 1 0 1 0 1 1 1 0 0 1 1 0 0 1 0 1 0 0 0
1 1 0 1 0 0 1 1 1 0 0 1 0 0 0 0 0 1 1 0 0 1 0 0 1 1 0 0 1 1 1 0 0 1 1 1 1
0 0 1 0 0 0 0 0 0 1 0 0 1 0 0 1 1 1 1 1 0 0 1 0 0 1 1 0 0 1 0 0 1 1 1 0
0 1 1 1 0 0 0 0 1 0 1 0 0 0 1 0 1 0 1 0 0 1 1 1 1 0 0 1 1 1 1 0 0 1 0 0
1 0 1 0 0 0 0 1 0 0 0 0 1 1 1]
```

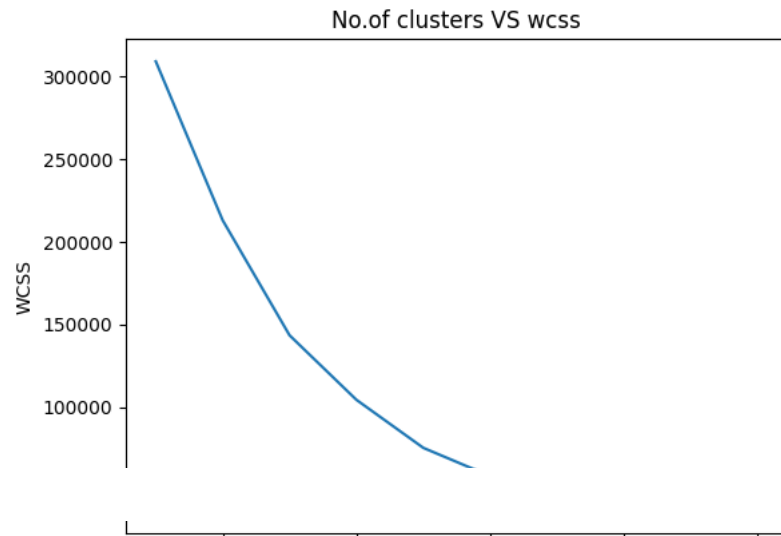
```
x.iloc[:,0]=le.transform(val)
```

```
<ipython-input-25-d0835554298a>:1: DeprecationWarning: In a future version, `df.iloc[:, i] = newvals` will attempt to set the values inplace instead of always setting a new array. To
x.iloc[:,0]=le.transform(val)
```

```
x.head()
```

[illegible]

8/10



k=6

```
kmeans=KMeans(n_clusters=6,init="k-means++",random_state=0)
ymeans=kmeans.fit_predict(x)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to silence this warning.
warnings.warn(
```

ymean

```
array([[5, 4, 5, 4, 5, 4, 5, 4, 5, 4, 5, 4, 5, 4, 5, 4, 5, 4, 5, 4,
        5, 4, 5, 4, 5, 4, 5, 4, 5, 4, 5, 4, 5, 4, 1, 4, 1, 0,
        5, 4, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0,
        1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0,
        0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1,
        1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0, 3, 0, 3, 2, 3, 2, 3, 2, 3,
        0, 3, 2, 3, 2, 3, 2, 3, 2, 3, 0, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3,
        2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3,
        2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3, 2, 3,
        2, 3], dtype=int32)
```

kmeans.cluster_centers_

```
array([[ 0.34210526, 27.          , 56.65789474, 49.13157895],
       [ 0.44444444, 56.15555556, 53.37777778, 49.08888889],
       [ 0.57142857, 41.68571429, 88.22857143, 17.28571429],
       [ 0.46153846, 32.69230769, 86.53846154, 82.12820513],
       [ 0.40909091, 25.27272727, 25.72727273, 79.36363636],
       [ 0.38095238, 44.14285714, 25.14285714, 19.52380952]])
```

