

K-Means Algorithm Assignment

```
In [6]: import pandas as pd
import numpy as np
from sklearn.cluster import KMeans

import matplotlib.pyplot as plt
import seaborn as sns

%matplotlib inline
```

```
In [14]: import warnings
warnings.filterwarnings('ignore')
```

```
In [3]: df = pd.read_csv(r"C:\Users\HP\Documents\ML-Assignment\DecisonTree\DecisonTree\Mall_Customers.csv")
df.head()
```

```
Out[3]:
```

	CustomerID	Gender	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

```
In [4]: df.shape
```

```
Out[4]: (200, 5)
```

```
In [5]: df.describe().T
```

```
Out[5]:
```

	count	mean	std	min	25%	50%	75%	max
CustomerID	200.0	100.50	57.879185	1.0	50.75	100.5	150.25	200.0
Age	200.0	38.85	13.969007	18.0	28.75	36.0	49.00	70.0
Annual Income (k\$)	200.0	60.56	26.264721	15.0	41.50	61.5	78.00	137.0
Spending Score (1-100)	200.0	50.20	25.823522	1.0	34.75	50.0	73.00	99.0

```
In [7]: df.isnull().sum()
```

```
Out[7]: CustomerID      0
Gender      0
Age      0
Annual Income (k$)      0
Spending Score (1-100)  0
dtype: int64
```

```
In [8]: df.columns
```

```
Out[8]: Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k$)',
              'Spending Score (1-100)'],
              dtype='object')
```

```
In [9]: df.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   Gender                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
```

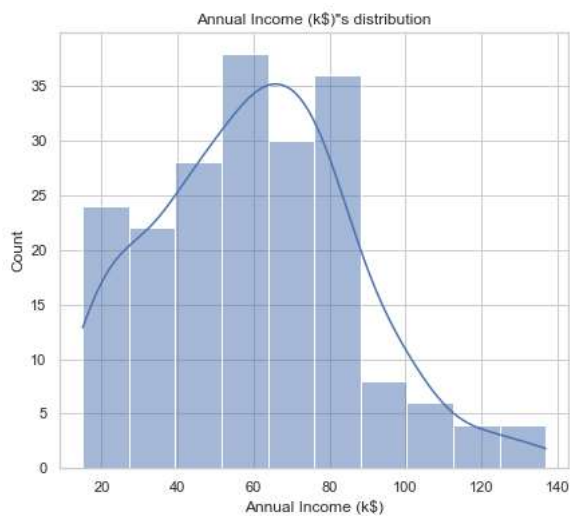
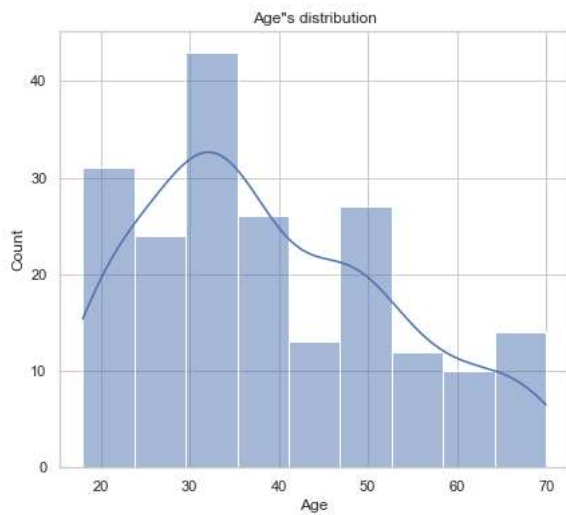
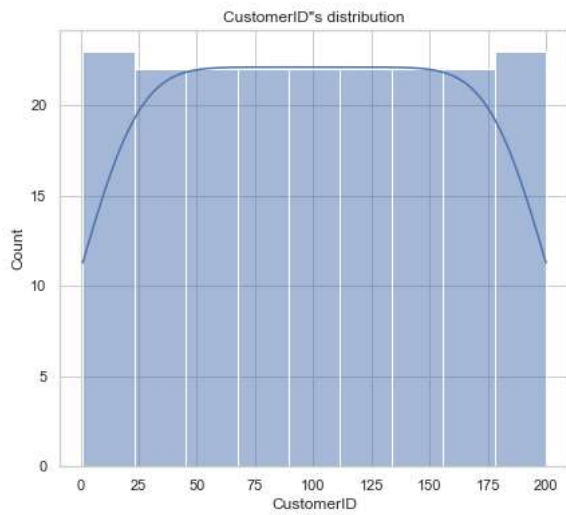
```
In [10]: numerical_col = [i for i in df.columns if df[i].dtype!='O']
numerical_col
```

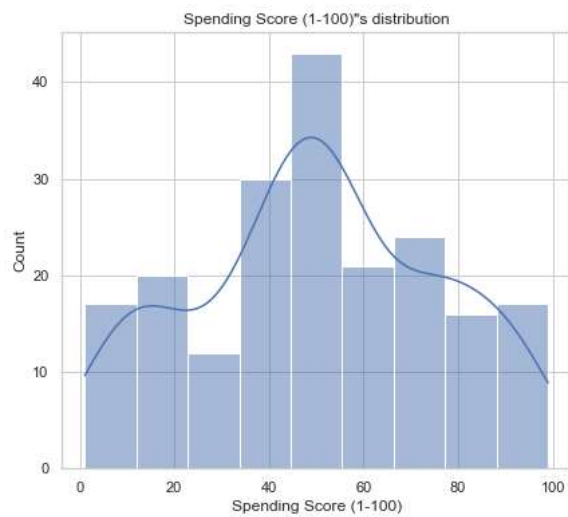
```
Out[10]: ['CustomerID', 'Age', 'Annual Income (k$)', 'Spending Score (1-100)']
```

```
In [11]: cat_col = [i for i in df.columns if df[i].dtype=='O']  
cat_col
```

```
Out[11]: ['Gender']
```

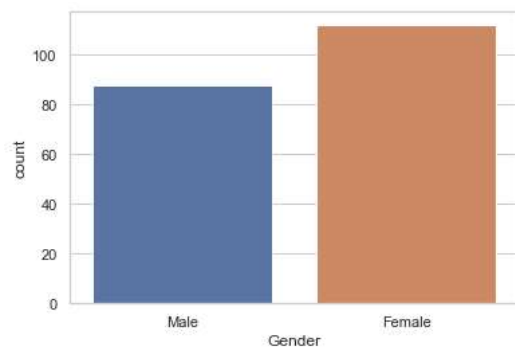
```
In [16]: for x in numerical_col:
sns.set(style = 'whitegrid')
plt.figure(figsize=(15,6))
plt.subplot(121)
sns.histplot(df,x=x,kde=True)
plt.title(f'{x}"s distribution')
plt.show()
```





```
In [17]: #for categorical plot  
sns.countplot(x='Gender',data=df)
```

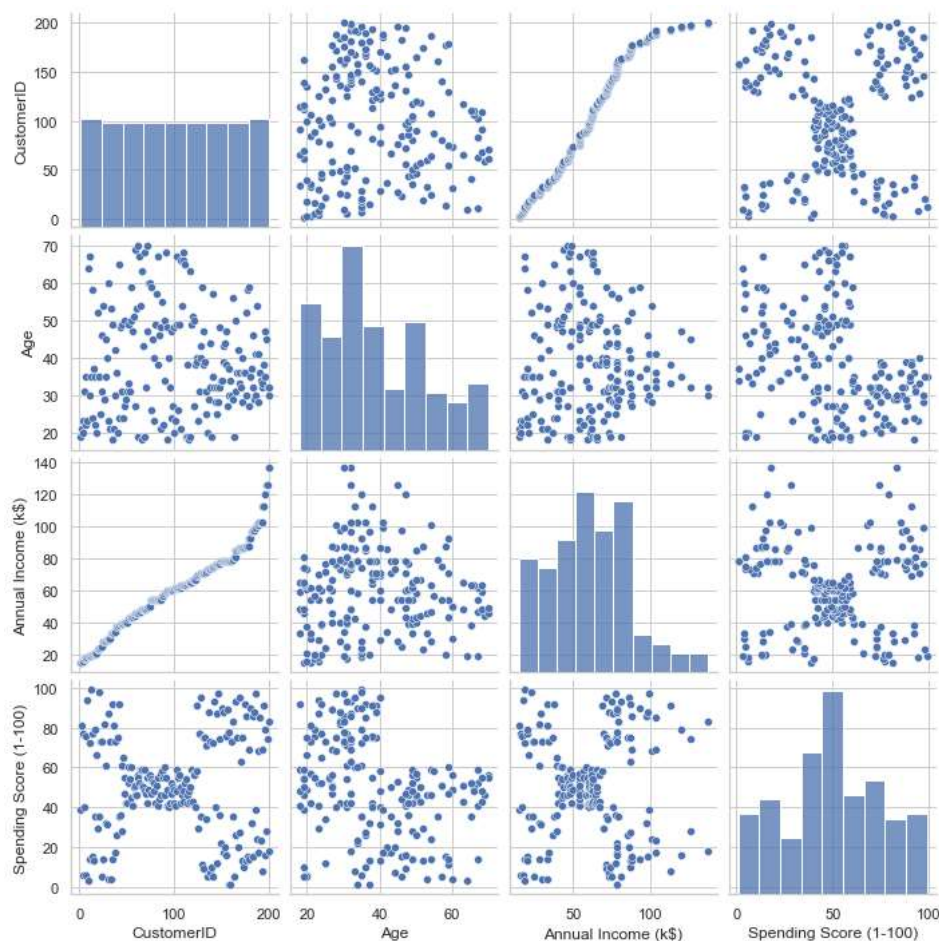
```
Out[17]: <AxesSubplot:xlabel='Gender', ylabel='count'>
```



Females are more spending

```
In [18]: sns.pairplot(df)
```

```
Out[18]: <seaborn.axisgrid.PairGrid at 0x22fc9ecbb0>
```



```
In [19]: df.columns
```

```
Out[19]: Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k$)',  
              'Spending Score (1-100)'],  
              dtype='object')
```

Clustering the data for 'Annual Income (k)', 'Spending Score (1-100)'

```
In [23]: X = df.iloc[:,[3,4]].values  
X
```

```
Out[23]: array([[ 15,  39],  
               [ 15,  81],  
               [ 16,   6],  
               [ 16,  77],  
               [ 17,  40],  
               [ 17,  76],  
               [ 18,   6],  
               [ 18,  94],  
               [ 19,   3],  
               [ 19,  72],  
               [ 19,  14],  
               [ 19,  99],  
               [ 20,  15],  
               [ 20,  77],  
               [ 20,  13],  
               [ 20,  79],  
               [ 21,  35],  
               [ 21,  66],  
               [ 23,  29],  
               ...,  
               ...,  
               ...])
```

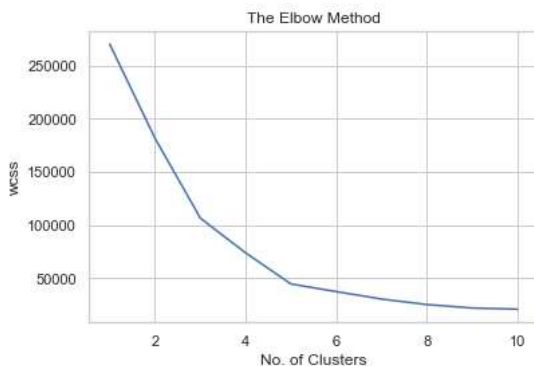
```
In [24]: wcss=[]
for i in range(1,11):
    km = KMeans(n_clusters=i,init = 'k-means++', max_iter = 300, n_init = 10, random_state = 0)
    km.fit(X)
    wcss.append(km.inertia_)
```

```
In [25]: wcss
```

```
Out[25]: [269981.28,
181363.59595959596,
106348.37306211119,
73679.78903948834,
44448.45544793371,
37265.86520484346,
30259.65720728547,
25095.703209997548,
21830.041978049438,
20736.679938924128]
```

Elbow Curve

```
In [26]: plt.plot(range(1,11),wcss)
plt.title('The Elbow Method')
plt.xlabel('No. of Clusters')
plt.ylabel('wcss')
plt.show()
```

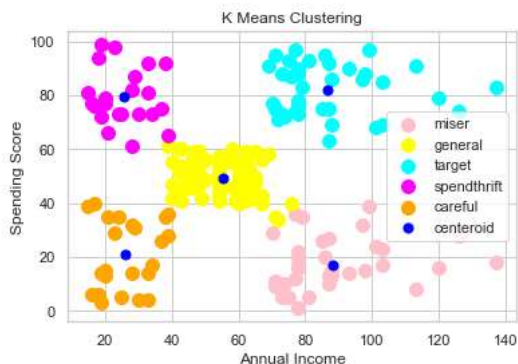


If we observe above diagram we can see the abrupt change at the 5 so we can keep the 5

```
In [27]: km = KMeans(n_clusters = 5, init = 'k-means++', max_iter = 300, n_init = 10, random_state = 0)
y_means = km.fit_predict(X)

plt.scatter(X[y_means == 0, 0], X[y_means == 0, 1], s = 100, c = 'pink', label = 'miser')
plt.scatter(X[y_means == 1, 0], X[y_means == 1, 1], s = 100, c = 'yellow', label = 'general')
plt.scatter(X[y_means == 2, 0], X[y_means == 2, 1], s = 100, c = 'cyan', label = 'target')
plt.scatter(X[y_means == 3, 0], X[y_means == 3, 1], s = 100, c = 'magenta', label = 'spendthrift')
plt.scatter(X[y_means == 4, 0], X[y_means == 4, 1], s = 100, c = 'orange', label = 'careful')
plt.scatter(km.cluster_centers[:,0], km.cluster_centers[:, 1], s = 50, c = 'blue', label = 'centroid')

plt.title('K Means Clustering')
plt.xlabel('Annual Income')
plt.ylabel('Spending Score')
plt.legend()
plt.show()
```



```
In [28]: km.labels_
```

[illegible]

```
In [29]: km.cluster_centers_
```

```
Out[29]: array([[88.2      , 17.11428571],
 [55.2962963 , 49.51851852],
 [86.53846154, 82.12820513],
 [25.72727273, 79.36363636],
 [26.30434783, 20.91304348]])
```

```
In [30]: km.inertia
```

Out[30]: 44448.45544793371

Hirachial Clustering

```
In [31]: from sklearn.cluster import AgglomerativeClustering
```

```
In [32]: ag = AgglomerativeClustering()
```

```
In [33]: ag.fit(X)
```

```
Out[33]: AgglomerativeClustering()
```

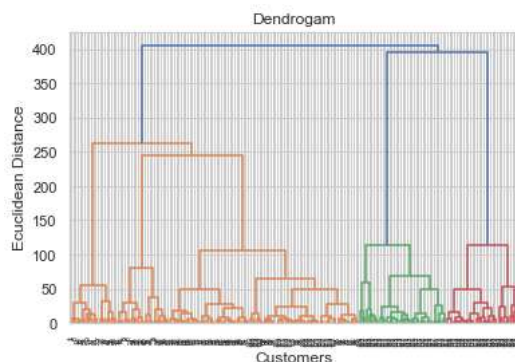
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [34]: ag.labels_
```

[illegible]

```
In [35]: import scipy.cluster.hierarchy as sch
```

```
In [37]: dendrogram = sch.dendrogram(sch.linkage(X, method = 'ward'))
plt.title('Dendrogram')
plt.xlabel('Customers')
plt.ylabel('Euclidean Distance')
plt.show()
```



In []: