

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
import matplotlib
warnings.filterwarnings('ignore')

In [2]: df=pd.read_excel('Customer Segmentation.xlsx')

In [3]: df.head()

Out[3]:
   InvoiceNo  StockCode      Description      Quantity      InvoiceDate  UnitPrice  CustomerID      Country
0   536365   85123A    WHITE HANGING HEART T-LIGHT HOLDER          6  2010-12-01 08:26:00      2.55   17850.0  United Kingdom
1   536365   71053              WHITE METAL LANTERN              6  2010-12-01 08:26:00      3.39   17850.0  United Kingdom
2   536365   844068      CREAM CUPID HEARTS COAT HANGER          8  2010-12-01 08:26:00      2.75   17850.0  United Kingdom
3   536365   84029G    KNITTED UNION FLAG HOT WATER BOTTLE          6  2010-12-01 08:26:00      3.39   17850.0  United Kingdom
4   536365   84029E      RED WOOLLY HOTTIE WHITE HEART.          6  2010-12-01 08:26:00      3.39   17850.0  United Kingdom

In [4]: df.shape

Out[4]: (541908, 8)

In [5]: df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 541908 entries, 0 to 541908
Data columns (total 8 columns):
 #   Column            Non-Null Count  Dtype
---  --
 0   InvoiceNo         541908 non-null  object
 1   StockCode        541908 non-null  object
 2   Description       540485 non-null  object
 3   UnitPrice         541908 non-null  float64
 4   InvoiceDate       541908 non-null  datetime64[ns]
 5   CustomerID       498829 non-null  float64
 6   Country          541908 non-null  object
dtypes: datetime64[ns](1), float64(2), int64(1), object(4)
memory usage: 33.1+ MB

In [6]: df.isnull().sum().sort_values(ascending=False)

Out[6]:
CustomerID    135088
Description    1454
InvoiceNo      0
StockCode      0
Description     0
Quantity       0
InvoiceDate     0
UnitPrice      0
Country        0
dtypes: int64      0

In [7]: df[df.dropna()
df.isnull().sum().sort_values(ascending=False)

Out[7]:
InvoiceNo      0
StockCode      0
Description     0
Quantity       0
InvoiceDate     0
UnitPrice      0
dtypes: int64      0
Country        0

In [8]: # we can see that quantity has negative values

Out[8]:
      Quantity  UnitPrice  CustomerID
count  40820.000000  40609.000000  40609.000000
mean    12.082323      3.460471  15294.335371
std     240.495370     49.315162  1713.600353
min    -40995.000000  0.000000    12346.000000
25%     2.000000     1.250000    13963.000000
50%     5.000000     1.950000    15152.000000
75%    12.000000     3.750000    16791.000000
max    80995.000000   8142.750000   18287.000000

In [9]: for i in df.columns:
    print(i, '-', len(df[i].unique()))

InvoiceNo = 22190
StockCode = 3884
Description = 3886
Quantity = 3896
InvoiceDate = 20460
UnitPrice = 628
CustomerID = 4372
Country = 37

In [10]: df=df.drop(df[df['Quantity']<0].index)

In [11]: df.describe()

Out[11]:
      Quantity  UnitPrice  CustomerID
count  39704.000000  39704.000000  39704.000000
mean    12.023823      3.116174   15294.335371
std     180.420210     22.096788   1713.169877
min     1.000000     0.000000    12346.000000
25%     2.000000     1.250000    13969.000000
50%     5.000000     1.950000    15159.000000
75%    12.000000     3.750000    16795.000000
max    80995.000000   8142.750000   18287.000000

In [12]: # adding a column named amount spent
df['AmountSpent']=df['Quantity']*df['UnitPrice']

In [13]: df.head()

Out[13]:
   InvoiceNo  StockCode      Description      Quantity      InvoiceDate  UnitPrice  CustomerID      Country  AmountSpent
0   536365   85123A    WHITE HANGING HEART T-LIGHT HOLDER          6  2010-12-01 08:26:00      2.55   17850.0  United Kingdom      15.30
1   536365   71053              WHITE METAL LANTERN              6  2010-12-01 08:26:00      3.39   17850.0  United Kingdom      20.34
2   536365   844068      CREAM CUPID HEARTS COAT HANGER          8  2010-12-01 08:26:00      2.75   17850.0  United Kingdom      22.00
3   536365   84029G    KNITTED UNION FLAG HOT WATER BOTTLE          6  2010-12-01 08:26:00      3.39   17850.0  United Kingdom      20.34
4   536365   84029E      RED WOOLLY HOTTIE WHITE HEART.          6  2010-12-01 08:26:00      3.39   17850.0  United Kingdom      20.34

In [14]: df['InvDay']=df['InvoiceDate'].dt.date
df['InvTime']=df['InvoiceDate'].dt.time
df['month']=df['InvoiceDate'].dt.month_name()
df['day']=df['InvoiceDate'].dt.day_name()
df.head()

Out[14]:
   InvoiceNo  StockCode      Description      Quantity      InvoiceDate  UnitPrice  CustomerID      Country  AmountSpent  InvDay  InvTime  month      day
0   536365   85123A    WHITE HANGING HEART T-LIGHT HOLDER          6  2010-12-01 08:26:00      2.55   17850.0  United Kingdom      15.30  2010-12-01 08:26:00  December  Wednesday
1   536365   71053              WHITE METAL LANTERN              6  2010-12-01 08:26:00      3.39   17850.0  United Kingdom      20.34  2010-12-01 08:26:00  December  Wednesday
2   536365   844068      CREAM CUPID HEARTS COAT HANGER          8  2010-12-01 08:26:00      2.75   17850.0  United Kingdom      22.00  2010-12-01 08:26:00  December  Wednesday
3   536365   84029G    KNITTED UNION FLAG HOT WATER BOTTLE          6  2010-12-01 08:26:00      3.39   17850.0  United Kingdom      20.34  2010-12-01 08:26:00  December  Wednesday
4   536365   84029E      RED WOOLLY HOTTIE WHITE HEART.          6  2010-12-01 08:26:00      3.39   17850.0  United Kingdom      20.34  2010-12-01 08:26:00  December  Wednesday
```

EDA

```
In [15]: # Top 5 Customers with the highest no of orders

In [16]: # Highest money spent by customers
sns.barplot(x=df['CustomerID'], y=df['AmountSpent'])
sns.sort_values(by='AmountSpent', ascending=False).iloc[:5]

Out[16]:
   CustomerID      Country  InvoiceNo
4019    17841.0  United Kingdom      7847
1888    14911.0         EIRE         5677
1298    14096.0  United Kingdom      5111
334     12748.0  United Kingdom      4096
1679    14086.0  United Kingdom      2700

In [17]: # Top 18 StockCodes by quantity
df4=df.groupby('StockCode').sum()
df4.sort_values(['Quantity'], ascending=False, inplace=True)
df4.reset_index(inplace=True)
df5=df4[['StockCode','Quantity']][:30]
df5

Out[17]:
   StockCode  Quantity
0          23013      80995
1          23166      77936
2          84077      54415
3          22197      49183
4          85096      46181
5          85123A      36782
6          84879      35362
7          21212      33693
8          2304      27254
9          22482      26076

In [18]: # Top 18 StockCodes by quantity
plt.figure(figsize=(12,6))
sns.barplot(x=df5['Quantity'], y=df5['StockCode'])
plt.title('Top 18 StockCodes by quantity')

Out[18]:
Text(0.5, 1.0, 'Top 18 StockCodes by quantity')

In [19]: TopCountries=df.groupby('Country')['AmountSpent'].sum().reset_index().sort_values('AmountSpent',ascending=False)
TopCountries

Out[19]:
   Country  AmountSpent
35  United Kingdom  7308391.554
23  Netherlands    285466.340
10         EIRE     26545.900
14         Germany  228867.140
13         France   209024.050
0         Australia  138521.310
30         Spain    61577.110
31  Switzerland    56443.950
3         Belgium   41196.340
31         Sweden   38278.330
19         Japan    37416.370
24         Norway   36165.440
26         Portugal  33439.890
12         Finland  22546.080
29         Singapore  21279.290
6         Channel Islands  20462.440
9         Denmark   18959.340
16         Italy    17483.240
7         Cyprus    13590.380
1         Austria   10196.680
25         Poland    7334.650
17         Israel    7221.690
15         Greece    4760.520
16         Iceland   4310.000
5         Canada     3466.340
33         USA       3060.390
22         Malta     2725.590
36         Unspecified  2667.070
34  United Arab Emirates  1902.280
20         Lebanon   1693.880
21         Lithuania  1661.060
11  European Community   1300.250
4         Brazil     1143.000
27         RSA        1002.310
8         Czech Repdica   828.140
2         Bahrain     546.600
28         Saudi Arabia  145.920

In [20]: # Top 5 countries where maximum sale happens.
plt.figure(figsize=(12,6))
sns.barplot(x=TopCountries['Country'], y=TopCountries['AmountSpent'].head(5))
plt.title('Top 5 Countries based on highest Amount Spent')

Out[20]:
Text(0.5, 1.0, 'Top 5 Countries based on highest Amount Spent')

In [21]: # Top 5 countries where least sell happens.
plt.figure(figsize=(25,6))
sns.barplot(x=TopCountries['Country'], y=TopCountries['AmountSpent'].tail(5))
plt.title('Top 5 Countries based on last store revenue contributors')

Out[21]:
Text(0.5, 1.0, 'Top 5 Countries based on last store revenue contributors')

In [22]: # Sales by Month
SalesByMonth=df.groupby('month')['AmountSpent'].sum().reset_index().sort_values('AmountSpent',ascending=False)
SalesByMonth

Out[22]:
   month  AmountSpent
9  November  1161817.380
2  December  1090069.680
10 October   1039318.790
11 September  952838.382
8  May       676594.560
6  June      66213.690
7  August    60543.900
5  July      60090.011
7  March     595500.700
4  January   566465.040
0  April     466200.361
3  February  447137.390

In [23]: # Sales in different months
plt.figure(figsize=(25,6))
sns.barplot(x=SalesByMonth['month'], y=SalesByMonth['AmountSpent'])
plt.title('Sales in different Months')

Out[23]:
Text(0.5, 1.0, 'Sales in different Months')

In [24]: # Sales on day
sales_on_day_basis=df.groupby('day')['AmountSpent'].sum().reset_index().sort_values('AmountSpent',ascending=False)
sales_on_day_basis

Out[24]:
   day  AmountSpent
3  Thursday  376869.070
4  Tuesday   370034.031
5  Wednesday  158836.170
0  Friday    1485917.401
1  Monday    1367146.411
2  Sunday    782514.221

In [25]: plt.figure(figsize=(10,6))
sns.barplot(x=sales_on_day_basis['day'], y=sales_on_day_basis['AmountSpent'])
plt.title('Sales on different Days')

Out[25]:
Text(0.5, 1.0, 'Sales on different Days')

In [26]: # outliers
for i in df[['Quantity','UnitPrice','AmountSpent']] :
    sns.boxplot(df[i])
    plt.show()

In [27]: cols = ['AmountSpent']
Q1 = df[cols].quantile(0.25)
Q3 = df[cols].quantile(0.75)
IQR = Q3 - Q1
df[~((df[cols] < (Q1 - 1.5 * IQR)) | (df[cols] > (Q3 + 1.5 * IQR)))].any(axis=1)

In [28]: <AxesSubplot: xlabel='AmountSpent'>

Out[28]:
In [29]: sns.distplot(df['AmountSpent'])

Out[29]:
<AxesSubplot: xlabel='AmountSpent', ylabel='Density'>

In [30]: from sklearn.preprocessing import StandardScaler
col_names = ['CustomerID','Quantity','UnitPrice']
features = df[col_names]
scaler = StandardScaler().fit(features.values)
scaled_features = pd.DataFrame(features, columns = col_names)
scaled_features.head()

Out[30]:
   CustomerID  Quantity  UnitPrice
0   1.480393  0.080345  0.080362
1   1.480393  0.080345  0.224636
2   1.480393  0.080345  0.090260
3   1.480393  0.080345  0.224636
4   1.480393  0.080345  0.224636

In [31]: # Machine Learning

In [32]: #convert amount spent to int
df['AmountSpent']=df['AmountSpent'].astype('int64')
df['UnitPrice']=df['UnitPrice'].astype(int)
df['CustomerID']=df['CustomerID'].astype('int64')

In [33]: from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()

In [34]: df['InvDescription']=le.fit_transform(df['Description'])
df['Country']=le.fit_transform(df['Country'])

In [35]: df=df.drop(columns=['StockCode','Description','InvDay','InvTime','month','day'])

In [36]: df.head()

Out[36]:
   InvoiceNo  Quantity      InvoiceDate  UnitPrice  CustomerID  Country  AmountSpent  InvDescription
0   536365          6  2010-12-01 08:26:00          2.55         17850      15      3674
1   536365          8  2010-12-01 08:26:00          2.75         17850      22      849
2   536365          6  2010-12-01 08:26:00          3.39         17850      20      1782
3   536365          6  2010-12-01 08:26:00          3.39         17850      20      2741

In [37]: df.drop_duplicates(inplace=True)

In [38]: df.shape

Out[38]: (361471, 8)

In [39]: X = df[['CustomerID','AmountSpent']]

In [40]: X

array([[17850,    15],
       [17850,    22],
       [17850,    22],
       [17850,    16],
       [17850,    16],
       [17850,    14]], dtype=int64)

In [41]: from sklearn.cluster import KMeans

In [42]: wcss = []

for i in range(1,11):
    kmeans = KMeans(n_clusters=i, init='k-means++', random_state=42)
    kmeans.fit(X)
    wcss.append(kmeans.inertia_)

print(wcss)

In [43]: fig = plt.figure(figsize=(10,5))
sns.lineplot(range(1,11), wcss, marker='o', color='red')

plt.title('The Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')

plt.show()

In [44]: kmeans = KMeans(n_clusters=4, init='k-means++', random_state=10)
y_means = kmeans.fit_predict(X)

In [45]: y_means

array([3, 3, 3, ..., 8, 8, 0])

In [46]: #visualising the clusters
plt.figure(figsize=(12,6))
sns.scatterplot(data = X, x = X[y_means==0,0], y = X[y_means==0,1], color='yellow', label='cluster1',s=5)
sns.scatterplot(data = X, x = X[y_means==1,0], y = X[y_means==1,1], color='blue', label='cluster2',s=5)
sns.scatterplot(data = X, x = X[y_means==2,0], y = X[y_means==2,1], color='green', label='cluster3',s=5)
sns.scatterplot(data = X, x = X[y_means==3,0], y = X[y_means==3,1], color='grey', label='cluster4',s=5)
sns.scatterplot(data = X, x = X[y_means==4,0], y = X[y_means==4,1], color='orange', label='cluster5',s=5)

sns.scatterplot(data=X,kmeans.cluster_centers_[0,0], y=kmeans.cluster_centers_[0,1], color='red',label='centroids', s=300, markers='x')
plt.xlabel('CustomerID')
plt.ylabel('AmountSpent')
plt.legend()
plt.show()

In [47]: from sklearn.metrics import DBSCAN
clustering_model=DBSCAN(eps=9, min_samples=4)

clustering_model.fit(X)

predict_labels = clustering_model.labels_
predict_labels

In [48]: plt.scatter(X[:,0],X[:,1],c=predict_labels)

plt.xlabel('Annual Income')
plt.ylabel('Spending Score')
plt.title('DBSCAN')

plt.show()

In [49]: from sklearn.cluster import DBSCAN
clustering_model= DBSCAN(eps=9, min_samples=4)

clustering_model.fit(X)

predict_labels = clustering_model.labels_
predict_labels

In [50]: # dendrogram = hierarchy.dendrogram(hierarchy.linkage(X,'average'))
# plt.title('Dendrogram')
# plt.xlabel('Customers')
# plt.ylabel('Linkage distance')
# plt.savefig('img,color='red')
# plt.show()

In [51]: tuned_clustering = KMeans(n_clusters=4, init='k-means++', random_state=10)
label = tuned_clustering.fit_predict(X)

In [52]: label

array([3, 3, 3, ..., 8, 8, 0])

In [53]: silhouette_score(X, tuned_clustering.labels_, metric='euclidean')

In [54]: from sklearn.cluster import AgglomerativeClustering

clustering_model = AgglomerativeClustering(n_clusters=4,affinity='euclidean')

clustering_model.fit(X)

clustering_prediction = clustering_model.fit_predict(X)
clustering_prediction

In [55]: silhouette_score(X, clustering_prediction, metric='euclidean')

In [56]: plt.figure(figsize=(16,10))
plt.scatter(X[clustering_prediction==0,0],
            X[clustering_prediction==0,1], s=30, c='green',
            label='cluster1')

plt.scatter(X[clustering_prediction==1,0],
            X[clustering_prediction==1,1], s=30,
            c='red',label='cluster2')

plt.scatter(X[clustering_prediction==2,0],
            X[clustering_prediction==2,1], s=30,
            c='blue',label='cluster3')

plt.scatter(X[clustering_prediction==3,0],
            X[clustering_prediction==3,1], s=30,
            c='orange',label='cluster4')

plt.xlabel('CustomerID')
plt.ylabel('InvDescription')
plt.legend()
plt.grid()
plt.show()

In [57]: from sklearn.cluster import DBSCAN
clustering_model= DBSCAN(eps=9, min_samples=4)

clustering_model.fit(X)

predict_labels = clustering_model.labels_
predict_labels
```