

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 import warnings #ignore the warnings
7 warnings.filterwarnings("ignore")
8 1.Dataset was downloaded
9 2.Loading the dataset
10 df = pd.read_csv('Mall_Customers.csv')
11 df
12 CustomerID→Gender→Age→Annual Income (k$)→Spending Score (1-100)
13 0→1→Male→19→15→39
14 1→2→Male→21→15→81
15 2→3→Female→20→16→6
16 3→4→Female→23→16→77
17 4→5→Female→31→17→40
18 .....
19 195→196→Female→35→120→79
20 196→197→Female→45→126→28
21 197→198→Male→32→126→74
22 198→199→Male→32→137→18
23 199→200→Male→30→137→83
24 200 rows x 5 columns
25
26 df.shape
27 (200, 5)
28 df.head()
29 CustomerID→Gender→Age→Annual Income (k$)→Spending Score (1-100)
30 0→1→Male→19→15→39
31 1→2→Male→21→15→81

```

```

31 1→2→Male→21→15→81
32 2→3→Female→20→16→6
33 3→4→Female→23→16→77
34 4→5→Female→31→17→40
35 3.Performing Visualizations
36 Univariate Analysis
37 sns.displot(df.Gender)
38
39 sns.displot(df.Age)
40
41 Bi-Variate analysis
42 sns.lineplot(df.Gender,df.Age)
43
44 sns.lineplot(df.CustomerID,df.Gender)
45
46 sns.pairplot(df)
47
48 4.Descriptive Statistics
49 df.describe()
50 CustomerID→Age→Annual Income (k$)→Spending Score (1-100)
51 count→200.000000→200.000000→200.000000→200.000000
52 mean→100.500000→38.850000→60.560000→50.200000
53 std→57.879185→13.969007→26.264721→25.823522
54 min→1.000000→18.000000→15.000000→1.000000
55 25%→50.750000→28.750000→41.500000→34.750000
56 50%→100.500000→36.000000→61.500000→50.000000
57 75%→150.250000→49.000000→78.000000→73.000000
58 max→200.000000→70.000000→137.000000→99.000000
59 5.Finding Missing Values And Replaceing it.
60 df.isnull().any()
61 CustomerID                False
62 Gender                    False

```

```

65 Spending Score (1-100)    False
66 dtype: bool
67 #ther is no any null values in the dataset
68 6.Finding Outliers And Replacing them
69 sns.boxplot(df.Age)
70
71 #ther is no outliers present
72 7.Checking Categorical columns and performing encoding.
73 df.head()
74 CustomerID→Gender→Age→Annual Income (k$)→Spending Score (1-100)
75 0→1→Male→19→15→39
76 1→2→Male→21→15→81
77 2→3→Female→20→16→6
78 3→4→Female→23→16→77
79 4→5→Female→31→17→40
80 Label Encoding
81 from sklearn.preprocessing import LabelEncoder
82 le = LabelEncoder()
83 df.Gender = le.fit_transform(df.Gender)
84 df.head()
85 CustomerID→Gender→Age→Annual Income (k$)→Spending Score (1-100)
86 0→1→1→19→15→39
87 1→2→1→21→15→81
88 2→3→0→20→16→6
89 3→4→0→23→16→77
90 4→5→0→31→17→40
91 8.Scaling the data
92 from sklearn.preprocessing import scale
93 x = df
94 x_scaled=pd.DataFrame(scale(x),columns=x.columns)
95 x_scaled.head()
96 CustomerID→Gender→Age→Annual Income (k$)→Spending Score (1-100)

```

```

92 from sklearn.preprocessing import scale
93 x = df
94 x_scaled=pd.DataFrame(scale(x),columns=x.columns)
95 x_scaled.head()
96 CustomerID→Gender→Age→Annual Income (k$)→Spending Score (1-100)
97 0→-1.723412→1.128152→-1.424569→-1.738999→-0.434801
98 1→-1.706091→1.128152→-1.281035→-1.738999→1.195704
99 2→-1.608771→-0.886405→-1.352802→-1.700830→-1.715913
100 3→-1.671450→-0.886405→-1.137502→-1.700830→1.040418
101 4→-1.654129→-0.886405→-0.563309→-1.662600→-0.395980
102 9.Performing Clustering Algorithm
103 from sklearn import cluster
104 error = []
105 for i in range(1,11):
106     kmeans=cluster.KMeans(n_clusters=i,init='k-means++',random_state=0)
107     kmeans.fit(df)
108     error.append(kmeans.inertia_)
109 error
110 [975512.06,
111 387065.7137713772,
112 271384.50878286787,
113 195401.19855991477,
114 157157.757905983,
115 122625.19813553884,
116 103233.01724386725,
117 80053.67444777439,
118 76938.97565600359,
119 69231.3360761156]
120 plt.plot(range(1,11),error)
121 plt.title('Elbow method')
122 plt.xlabel('no of clus')
123 plt.ylabel('error')

```

```

126 km_model = cluster.KMeans(n_clusters = 3, init='k-means++',random_state=0)
127 km_model.fit(df)
128 KMeans(n_clusters=3, random_state=0)
129 ykmeans = km_model.predict(df)
130 ykmeans
131 array([[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
132         0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
133         0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1,
134         1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
135         1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
136         1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
137         2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
138         2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
139         2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
140         2, 2])
141 km_model.predict([[1,1,10,15,30]])
142 array([0])
143 10.Adding cluster data with primary Dataset
144 df['kclus'] = pd.Series(ykmeans)
145 df.head()
146 CustomerID→Gender→Age→Annual Income (k$)→Spending Score (1-100)→kclus
147 0→1→1→19→15→39→0
148 1→2→1→21→15→81→0
149 2→3→0→20→16→6→0
150 3→4→0→23→16→77→0
151 4→5→0→31→17→40→0
152 Splitting data into dependent and independent variables.
153 X = df.iloc[:, :-1]
154 X
155 CustomerID→Gender→Age→Annual Income (k$)→Spending Score (1-100)
156 0→1→1→19→15→39
157 1→2→1→21→15→81

```

```

158 2→3→0→20→16→6
159 3→4→0→23→16→77
160 4→5→0→31→17→40
161 .....
162 195→196→0→35→120→79
163 196→197→0→45→126→28
164 197→198→1→32→126→74
165 198→199→1→32→137→18
166 199→200→1→30→137→83
167 200 rows x 5 columns
168
169 y=df.kclus
170 12.Splitting data into training and testing data
171 from sklearn.model_selection import train_test_split
172 X_train,X_test,y_train,y_test=train_test_split(x_scaled,y,test_size=0.3,random_state=0)
173 X_train.shape
174 (140, 5)
175 X_test.shape
176 (60, 5)
177 y_train.shape
178 (140,)
179 y_test.shape
180 (60,)
181 13.Building the model.
182 from sklearn.neighbors import KNeighborsClassifier
183 model =KNeighborsClassifier()
184 14.Training the model.
185 model.fit(X_train,y_train)
186 KNeighborsClassifier()
187 ## 15.Testing the model.
188 pred_test=model.predict(X_test)
189

```

```

182 from sklearn.neighbors import KNeighborsClassifier
183 model =KNeighborsClassifier()
184 14.Training the model.
185 model.fit(X_train,y_train)
186 KNeighborsClassifier()
187 ## 15.Testing the model.
188 pred_test=model.predict(X_test)
189 pred_train=model.predict(X_train)
190 16.Evaluating the modle using evaluation metrics
191 from sklearn.metrics import accuracy_score,classification_report,confusion_matrix
192 print('Test accuracy score: ',accuracy_score(y_test,pred_test))
193 print('Training accuracy score: ',accuracy_score(y_train,pred_train))
194 Test accuracy score: 0.8333333333333334
195 Training accuracy score: 0.9571428571428572
196 pd.crosstab(y_test,pred_test)
197 col_0→0→1→2
198 kclus→→→→
199 0→16→4→0
200 1→1→15→4
201 2→0→1→10
202 print(classification_report(y_test,pred_test))
203      precision    recall  f1-score   support
204
205     0       0.94      0.80      0.86         20
206     1       0.75      0.75      0.75         20
207     2       0.83      0.95      0.88         20
208
209    accuracy                   0.83         60
210   macro avg       0.84      0.83      0.83         60
211   weighted avg       0.84      0.83      0.83         60
212

```