| | Name = Nayan.H.Kacha Project Domain = Data Science Project Name = CUSTOMER SEGMENTATION USING K-MEANS CLUSTERING Company Name = Exposys Data Labs Importing the Dependencies |
|------------------------------|--|
| In [80]: | <pre>import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns from sklearn.cluster import KMeans Data Collection & Analysis # First 5 rows in the dataframe df = pd.read_csv("Mall_Customers.csv") df.head()</pre> |
| Out[81]: In [79]: | 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 |
| Out[79]: In [82]: | #finding the number of rows and columes in dataset df.shape (200, 5) # getting the info about the dataset df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 200 entries, 0 to 199 Data columns (total 5 columns): # Column Non-Null Count Dtype</class> |
| In [83]: Out[83]: | 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 df.describe() CustomerID Age Annual Income (k\$) Spending Score (1-100) count 200.000000 200.000000 200.000000 200.000000 |
| In [84]: | mean 100.500000 38.850000 60.560000 50.200000 std 57.879185 13.969007 26.264721 25.823522 min 1.000000 18.000000 15.00000 1.000000 25% 50.750000 28.750000 41.500000 34.750000 50% 100.500000 36.000000 50.000000 75% 150.250000 49.000000 78.000000 73.000000 max 200.000000 70.000000 137.000000 99.000000 |
| Out[84]: In [85]: Out[85]: | CustomerID int64 Gender object Age int64 Annual Income (k\$) int64 Spending Score (1-100) int64 dtype: object # checking the missing values df.isnull().sum() CustomerID 0 Gender 0 |
| In [86]: | Age 0 Annual Income (k\$) 0 Spending Score (1-100) 0 dtype: int64 plt.figure(1, figsize=(15,6)) n=0 for x in ['Age','Annual Income (k\$)', 'Spending Score (1-100)']: n+=1 plt.subplot(1, 3, n) plt.subplots_adjust(hspace =0.5, wspace = 0.5) sns.distplot(df[x], bins = 20) plt.title('Distplot of {}'.format(x)) plt.show() |
| | c:\users\hp\appdata\local\programs\python\python39\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and wil l be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-lev el function for histograms). warnings.warn(msg, FutureWarning) c:\users\hp\appdata\local\programs\python\python39\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and wil l be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-lev el function for histograms). warnings.warn(msg, FutureWarning) c:\users\hp\appdata\local\programs\python\python39\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and wil l be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-lev el function for histograms). warnings.warn(msg, FutureWarning) Distplot of Age Distplot of Annual Income (k\$) Distplot of Spending Score (1-100) |
| | Distplot of Age 0.020 0.035 0.030 0.025 0.0025 0.0025 0.0025 0.0025 0.0075 0.0075 0.0075 0.0075 |
| In [87]: | 0.0050 0.0005 0. |
| | Male Male |
| In [88]: | Plt.figure(1, figsize=(15,7)) n=0 for cols in ['Age' , 'Annual Income (k\$)' , 'Spending Score (1-100)']: n+=1 plt.subplot(1 , 3 , n) |
| | <pre>sns.set(style="whitegrid") plt.subplots_adjust(hspace = 0.5 , wspace = 0.5) sns.violinplot(x = cols , y = 'Gender' , data = df) plt.ylabel('Gender' if n == 1 else '') plt.title('Violin Plot') plt.show()</pre> Violin Plot Violin Plot Violin Plot Wale Male |
| | Female Female 0 50 100 150 0 50 100 Spending Score (1-100) |
| In [89]: | Calculating the Number of Customers and Ages in the Dataset age_18_25 = df.Age[(df.Age >=18) & (df.Age <=25)] age_26_35 = df.Age[(df.Age >=26) & (df.Age <=35)] age_36_45 = df.Age[(df.Age >=36) & (df.Age <=45)] age_46_55 = df.Age[(df.Age >=36) & (df.Age <=45)] age_46_55 = df.Age[(df.Age >=46) & (df.Age <=55)] age_55above = df.Age[(df.Age >= 56] agex = ["18-25", "26-35", "36-45", "46-55", "55+"] agey = [len(age_18_25.values), len(age_26_35.values), len(age_36_45.values), len(age_46_55.values), len(age_55above.values)] plt.figure(figsize=(15,6)) sns.barplot(x= agex, y= agex, palette="rocket") plt.title("Number of Customer and Ages") plt.ylabel("Number of Customer") plt.show() Number of Customer and Ages |
| | Number of Customer and Ages 50 50 20 20 20 20 20 20 20 20 |
| In [90]: Out[90]: | 10 18-25 26-35 36-45 46-55 55+ sns.relplot(x="Annual Income (k\$)", y="Spending Score (1-100)", data=df) <seaborn.axisgrid.facetgrid 0x2a2c41df7c0="" at=""> 100 100 100 100 100 100 100 100 100 10</seaborn.axisgrid.facetgrid> |
| | 80 (001-1) 60 20 40 60 80 100 120 140 Annual Income (k\$) |
| In [91]: | Ss_1_20 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 1) & (df["Spending Score (1-100)"] <= 20)] ss_21_40 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 21) & (df["Spending Score (1-100)"] <= 40)] ss_41_60 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 41) & (df["Spending Score (1-100)"] <= 60)] ss_61_80 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 61) & (df["Spending Score (1-100)"] <= 80)] ss_81_100 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 81) & (df["Spending Score (1-100)"] <= 100)] ssx = ["1-20", "21-40", "41-60", "61-80", "81-100"] ssy = [len(ss_1_20.values), len(ss_21_40.values), len(ss_41_60.values), len(ss_61_80.values), len(ss_81_100.values)] |
| | plt.figure(figsize=(15,6)) sns.barplot(x= ssx, y= ssy, palette="mako") plt.title("Spending Scores") plt.xlabel("Score") plt.ylabel("Number of Customer having the Scores") plt.show() Spending Scores |
| | 1-20 21-40 41-60 61-80 81-100 |
| In [92]: | 1-20 21-40 41-60 61-80 81-100 Calculating the Annual Incomes of the customers in the Dataset ai0_30 = df["Annual Income (k\$)"][(df["Annual Income (k\$)"] >= 0) & (df["Annual Income (k\$)"] <= 30)] ai31_60 = df["Annual Income (k\$)"][(df["Annual Income (k\$)"] >= 31) & (df["Annual Income (k\$)"] <= 60)] ai61_90 = df["Annual Income (k\$)"][(df["Annual Income (k\$)"] >= 61) & (df["Annual Income (k\$)"] <= 90)] ai91_120 = df["Annual Income (k\$)"][(df["Annual Income (k\$)"] >= 91) & (df["Annual Income (k\$)"] <= 120)] ai121_150 = df["Annual Income (k\$)"][(df["Annual Income (k\$)"] >= 121) & (df["Annual Income (k\$)"] <= 150)] aix = ["\$ 0 - 30,000", "\$ 30,001 - 60,000", "\$ 60,001 - 90,000", "\$ 90,001 - 120,000", "\$120,001 - 150,000"] aiy = [len(ai0_30.values), len(ai31_60.values), len(ai61_90.values), len(ai91_120.values), len(ai121_150.values)] plt.figure(figsize=(15,6)) |
| | sns.barplot(x= aix, y= aiy, palette="Spectral") plt.title("Annual Incomes") plt.xlabel("Income") plt.ylabel("Number of Customer") plt.show() Annual Incomes 60 |
| | \$0 - 30,000 \$30,001 - 60,000 \$60,001 - 90,000 hoome \$90,001 - 120,000 \$120,001 - 150,000 |
| In [93]: | <pre>X1= df.loc[:, ["Age", "Spending Score (1-100)"]].values from sklearn.cluster import KMeans wcss = [] for k in range(1,11): kmeans = KMeans(n_clusters=k, init="k-means++") kmeans.fit(X1) wcss.append(kmeans.inertia_) plt.figure(figsize=(12,6)) plt.grid() plt.plot(range(1,11),wcss, linewidth=2, color="red", marker ="8") plt.xlabel("K Value")</pre> |
| | plt.ylabel("WCSS") plt.show() 160000 140000 120000 88 80000 |
| In [94]: | 60000 20000 2 4 K Value 6 8 10 kmeans = KMeans(n_clusters=4) |
| In [95]: | label = kmeans.fit_predict(X1) print(label) [2 1 0 1 2 1 0 1 0 1 0 1 0 1 0 1 2 2 0 1 2 1 0 1 0 |
| In [96]: | [30.1754386 82.35087719] [27.61702128 49.14893617] [55.70833333 48.22916667]] Cluster plotting between AGE and SPENDING SCORES plt.scatter(X1[:,0], X1[:,1], c=kmeans.labels_, cmap="rainbow") plt.scatter(kmeans.cluster_centers_[:,0], kmeans.cluster_centers_[:,1], color="black") plt.title("Clusters of Customers") plt.ylabel("Age") plt.ylabel("Spending Score(1-100)") plt.show() |
| | 100 80 60 20 20 30 40 50 60 70 Age |
| In [121 | ELBOW METHOD between Annual Income and Spending scores X2=df.loc[:, ["Annual Income (k\$)", "Spending Score (1-100)"]].values from sklearn.cluster import KMeans wcss = [] for k in range(1,11): kmeans = KMeans(n_clusters=k, init="k-means++") kmeans.fit(X2) wcss.append(kmeans.inertia_) plt.figure(figsize=(12,6)) plt.grid() plt.plot(range(1,11),wcss, linewidth=2, color="red", marker ="8") |
| | plt.xlabel("K Value") plt.ylabel("WCSS") plt.show() 250000 |
| | 100000 50000 2 4 K Value 8 10 |
| In [108 | <pre>kmeans = KMeans(n_clusters=5) label = kmeans.fit_predict(X2) print(label) [2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3 2 3</pre> |
| In [110 | [[55.2962963 |
| | Clusters of Customers 00 00 00 00 00 00 00 00 00 |
| In [122 | 20 40 60 80 100 120 140 ELBOW METHOD FOR ALL AGE ,ANNUAL INCOME ,SPENDING SCORES. X3=df.loc[:, ["Age", "Annual Income (k\$)", "Spending Score (1-100)"]].values wcss= [] for k in range(1,11): kmeans = KMeans(n_clusters=k, init="k-means++") kmeans.fit(X3) wcss.append(kmeans.inertia_) plt.figure(figsize=(12,6)) plt.grid() |
| | plt.plot(range(1,11),wcss, linewidth=2, color="red", marker ="8") plt.xlabel("K Value") plt.ylabel("WCSS") plt.show() 250000 200000 |
| | 150000 100000 50000 2 4 K Value 6 8 10 |
| In [123 In [125 | kmeans = KMeans(n_clusters = 5) label = kmeans.fit_predict(X3) print(label) [3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 3 1 |
| In []: | [[40.66666667 87.75 |
| | <pre>from mpl_toolkits.mplot3d import Axes3D fig = plt.figure(figsize=(20,10)) ax = fig.add_subplot(111, projection="3d") ax.scatter(df.Age[df.label == 0], df["Annual Income (k\$)"][df.label == 0], df["Spending Score (1-100)"][df.label == 0], c="blue", s=60) ax.scatter(df.Age[df.label == 1], df["Annual Income (k\$)"][df.label == 1], df["Spending Score (1-100)"][df.label == 1], c="red", s=60) ax.scatter(df.Age[df.label == 2], df["Annual Income (k\$)"][df.label == 2], df["Spending Score (1-100)"][df.label == 2], c="green", s=60) ax.scatter(df.Age[df.label == 3], df["Annual Income (k\$)"][df.label == 3], df["Spending Score (1-100)"][df.label == 3], c="green", s=60) ax.scatter(df.Age[df.label == 4], df["Annual Income (k\$)"][df.label == 4], df["Spending Score (1-100)"][df.label == 4], c="purple", s=60) ax.view_init(30, 185) plt.xlabel("Age") plt.ylabel("Annual Income (k\$)") ax.set_zlabel("Spending Score (1-100)") plt.show()</pre> |
| | |
| | 100 80 140 120 100 80 60 50 40 40 Annual Income (k\$) |
| In []: | |
| | |
| | |
| | |