**Predicting Customer spending category** In this assignment we are going to classify the customer according to the amount he spends. We will study the customer spending patterns and will categorize them into categories based on within how much they spend on the website. Then we will train a classification model that will classify which customer belongs to which category. In [24]: #loading the required packages import pandas as pd import matplotlib.pyplot as plt import seaborn as sb import numpy as np from sklearn.naive bayes import MultinomialNB from sklearn.model\_selection import train\_test\_split from sklearn.metrics import classification report, accuracy score from sklearn.preprocessing import LabelEncoder from sklearn.tree import DecisionTreeClassifier In [2]: # Load the datasets online sales = pd.read csv('C:\\Users\\sujoydutta\\Desktop\\Data analysis\\Projects\\Marketing insights\\Online customer data = pd.read excel('C:\\Users\\sujoydutta\\Desktop\\Data analysis\\Projects\\Marketing insights\\Cus discount coupons = pd.read csv('C:\\Users\\sujoydutta\\Desktop\\Data analysis\\Projects\\Marketing insights\\Di marketing spend= pd.read csv('C:\\Users\\sujoydutta\\Desktop\\Data analysis\\Projects\\Marketing insights\\Mark tax amount = pd.read excel('C:\\Users\\sujoydutta\\Desktop\\Data analysis\\Projects\\Marketing insights\\Tax am In [3]: # Merging datasets using 'CustomerID' and 'Product Category as the common key merged data = pd.merge(online sales, customer data, on='CustomerID', how='left') merged data = pd.merge(merged data, discount coupons, on='Product Category', how='left') merged data = pd.merge(merged data, tax amount, on='Product Category', how='left') # Convert 'Date' column in the marketing spend DataFrame to datetime data type In [4]: marketing spend['Date'] = pd.to datetime(marketing spend['Date']) merged data['Transaction Date'] = pd.to datetime(merged data['Transaction Date'], format='%Y%m%d') # Merge datasets using 'Transaction Date' as the common key merged data = pd.merge(merged data, marketing spend, left on='Transaction Date', right on='Date', how='left') # Dropping the duplicate 'Date' column if needed merged data.drop(columns=['Date'], inplace=True) # Calculating Invoice Value for each transaction in the merged dataset In [5]: merged data['Invoice Value'] = ( (merged data['Quantity'] \* merged data['Avg Price']) \* (1 - merged data['Discount pct']) \* (1 + merged data['GST']) + merged data['Delivery Charges'] ) .abs() #examining the merged dataset merged data.head() CustomerID Transaction\_ID Transaction\_Date Out[6]: Product\_SKU Product\_name Product\_Category Quantity Avg\_Price Delivery\_Charges Cc Nest Learning Thermostat 0 17850 16679 2019-01-01 GGOENEBJ079499 Nest-USA 1 153.71 6.5 3rd Gen-USA -Stainle... **Nest Learning** Thermostat 1 17850 16679 2019-01-01 GGOENEBJ079499 Nest-USA 1 153.71 6.5 3rd Gen-USA -Stainle... Nest Learning Thermostat 17850 16679 2019-01-01 GGOENEBJ079499 Nest-USA 153.71 6.5 3rd Gen-USA -Stainle... Nest Learning Thermostat 3 17850 2019-01-01 GGOENEBJ079499 6.5 16679 Nest-USA 153.71 3rd Gen-USA -Stainle... **Nest Learning** Thermostat 4 17850 16679 2019-01-01 GGOENEBJ079499 Nest-USA 153.71 6.5 3rd Gen-USA -Stainle 5 rows × 21 columns #correcting the format merged\_data['Transaction\_Date'] = pd.to\_datetime(merged\_data['Transaction\_Date']) merged\_data['Transaction\_Date'] 2019-01-01 Out[7]: 2019-01-01 2019-01-01 3 2019-01-01 2019-01-01 2019-12-31 630683 630684 2019-12-31 630685 2019-12-31 630686 2019-12-31 630687 2019-12-31 Name: Transaction Date, Length: 630688, dtype: datetime64[ns] In [11]: # Formatting the dataset customer data = merged data.groupby('CustomerID').agg({ 'Transaction ID': 'count', 'Transaction Date': 'max', 'Quantity': 'sum', 'Avg Price': 'mean', 'Delivery Charges': 'sum', 'Gender': 'first', 'Location': 'first', 'Tenure Months': 'first', 'Offline Spend': 'sum', 'Online Spend': 'sum', 'Invoice Value': 'sum', }).reset index() customer data.round() Out[11]: CustomerID Transaction\_ID Transaction\_Date Quantity Avg\_Price Delivery\_Charges Gender Location Tenure\_Months Offline\_Spend 0 12346 24 2019-09-15 36 1800.0 F New York 31 72000 13.0 12347 709 2019-11-02 4082 7910.0 20 1418000 63.0 New York 2 12348 276 2019-10-19 2508 15.0 2366.0 California 39 894000 3 12350 204 2019-12-14 252 77.0 1535.0 California 25 816000 4 12356 432 2019-09-15 672 35.0 7650.0 Chicago 31 1296000 1463 18259 73 2019-04-05 541 11.0 3268.0 F California 5 182500 1464 18260 469 2019-10-05 1549 28.0 7387.0 New York 43 1473000 1465 18269 96 2019-06-20 120 11.0 618.0 Chicago 25 252000 M 1466 12 2019-10-23 149.0 42000 18277 24 72.0 Chicago 47 New 8547.0 1467 18283 1213 2019-10-10 1816 56.0 36 3572500 Jersey 1468 rows × 12 columns # Distributing the customers according to Gaussian norms In [13]: quantiles = [0, 0.33, 0.66, 1.0]labels=['Low-value','Mid-value','High-value'] customer\_data['Invoice\_Value\_Category'] = pd.qcut(customer\_data['Invoice\_Value'], q=quantiles, labels=labels) print(customer data[['Invoice Value', 'Invoice Value Category']]) Invoice Value Invoice Value Category 6.336910e+03 0 Low-value 3.612323e+06 High-value 3.728585e+05 Mid-value 3.423656e+05 Mid-value Mid-value 3.558345e+05 1463 1.277934e+05 Low-value 1464 5.820105e+05 Mid-value 1465 2.664770e+04 Low-value 1466 7.466640e+04 Low-value 1467 1.571264e+06 High-value [1468 rows x 2 columns] In [16]: #examining the final dataset customer data.head() Out[16]: CustomerID Transaction\_ID Transaction\_Date Quantity Avg\_Price Delivery\_Charges Gender Location Tenure\_Months Offline\_Spend 0 12346 2019-09-15 36 12.745000 1800.00 F New York 72000 31 12347 709 2019-11-02 4082 63.231072 7910.18 M New York 1418000 2508 14.631304 2 12348 276 2019-10-19 M California 39 894000 2365.80 3 12350 204 2019-12-14 252 77.200000 M California 816000 1534.56 4 12356 432 2019-09-15 672 34.578611 7649.88 Chicago 31 1296000 In [17]: # Creating an instance of LabelEncoder label encoder = LabelEncoder() # List of categorical columns to encode categorical columns = ['Gender', 'Location'] # Apply Label Encoding to each categorical column for col in categorical columns: customer data[col] = label\_encoder.fit\_transform(customer\_data[col]) In [19]: # Selecting Features and Target Variable features = ['CustomerID', 'Quantity', 'Avg Price', 'Delivery Charges', 'Gender', 'Location', 'Tenure Months', X = customer data[features] y = customer data['Invoice Value Category'] In [20]: #Splitting the Data into Training and Testing Sets X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.3, random\_state=42) In [21]: # Creating and training a Multinomial Naive Bayes Classifier clf = MultinomialNB() clf.fit(X\_train, y\_train) Out[21]: ▼ MultinomialNB MultinomialNB()In [22]: # Making predictions on the test set y\_pred = clf.predict(X\_test) y\_pred array(['Mid-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'Low-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value' 'Low-value', 'Mid-value', 'Low-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Low-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'Low-value', 'Low-value', 'High-value', 'High-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'High-value', 'Low-value', 'Low-value', 'Mid-value', 'Mid-value', 'High-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'Low-value', 'Low-value', 'Low-value', 'Mid-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'Low-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'How-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'Low-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'Mid-value' 'High-value', 'High-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'Mid-value'], dtype='<U10') In [23]: # Model Performance Evaluation accuracy = accuracy score(y test, y pred) classification report result = classification report(y test, y pred) # Printing the evaluation metrics print("Accuracy:", accuracy) print("\nClassification Report:\n", classification report result) Accuracy: 0.6213151927437641 Classification Report: recall f1-score support precision 0.74 0.69 0.61 0.67 High-value 149 Low-value 0.83 0.75 140 0.44 Mid-value 0.46 0.45 152 441 0.62 accuracy 0.64 0.62 0.63 0.62 macro avg 0.63 441 441 weighted avg 0.62 **Remark:** The accuracy of the model is 64% and does a decent job in classifying the customers into categories. In [25]: # Training the Decision Tree Classifier clf = DecisionTreeClassifier(random state=42) clf.fit(X train, y train) Out[25]: DecisionTreeClassifier DecisionTreeClassifier(random\_state=42) In [26]: #Making the Predictions y pred = clf.predict(X test) y\_pred Out[26]: array(['High-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'Low-value', 'Low-value', 'Mid-value', 'Low-value', 'Low-value', 'Low-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Low-value', 'Mid-value', 'Low-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Low-value', 'Low-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'Low-value', 'Mid-value', 'Low-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'Low-value', 'Mid-value', 'Mid-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'Low-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'Low-value', 'Low-value', 'Mid-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'Low-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'Low-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'Low-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'Low-value', 'Low-value', 'High-value', 'High-value', 'Low-value', 'Mid-value', 'Low-value', 'Low-value', 'High-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'Low-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'Mid-value', 'High-value', 'Low-value', 'Low-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'Low-value', 'Mid-value', 'High-value', 'Low-value', 'Mid-value', 'High-value', 'Mid-value', 'Low-value', 'High-value', 'Mid-value', 'Mid-value', 'Low-value', 'High-value', 'High-value', 'Mid-value', 'Mid-value', 'Mid-value', 'Low-value', 'Mid-value', 'Low-value', 'Mid-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'Mid-value', 'High-value', 'High-value', 'High-value', 'High-value', 'High-value', 'High-value', 'Low-value', 'High-value', 'Low-value', 'Low-value', 'High-value', 'Mid-value', 'High-value', 'Low-value', 'Mid-value', 'Low-value', 'Low-value', 'High-value', 'High-value', 'Low-value', 'Low-value', 'High-value', 'Low-value', 'Mid-value', 'Low-value', 'Mid-value'], dtype=object) In [27]: # Model Performance Evaluation accuracy = accuracy score(y test, y pred) classification report result = classification report(y test, y pred) # Printing the evaluation metrics print("Accuracy:", accuracy) print("\nClassification Report:\n", classification report result) Accuracy: 1.0 Classification Report: precision recall f1-score High-value 1.00 1.00 1.00 Low-value 1.00 1.00 1.00 140 Mid-value 1.00 1.00 1.00 1.00 441 accuracy macro avg 1.00 1.00 1.00 weighted avg 1.00 1.00 1.00 441 **Remark:** The Decision Tree model is 100% accurate and we can say that it can classify customers into Spending categories perfectly.