

```
In [1]: import numpy as np
import pandas as pd
df = read_csv("data.csv")
df.columns

Out[1]: Index(['Customer_ID', 'Age', 'Gender', 'Income', 'Spending_Score',
              'Credit_Score', 'Loan_Amount', 'Previous_Defaults', 'Marketing_Spend',
              'Purchase_Frequency', 'Seasonality', 'Sales', 'Customer_Churn',
              'Defaulted', 'Seasonality_Date'],
              dtype='object')

In [3]: df.describe()

Out[3]:
```

	Customer_ID	Age	Income	Spending_Score	Credit_Score	Loan_Amount	Previous_Defaults	Marketing_Spend	Purchase_Frequency	Sales	Customer_Churn	Defaulted
count	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000
mean	250.500003	44.200000	84398.055556	50.860000	573.411111	29466.928889	0.974000	10568.128000	1.000000	15360000	54378.964000	0.190000
std	144.818183	15.000000	38049.938777	20.125101	149.302942	11788.254534	0.80825	5508.219008	0.8745327	27263.106468	0.435734	0.392894
min	1.000000	18.000000	20055.000000	1.000000	300.000000	1563.000000	0.000000	1024.000000	1.000000	5200.000000	0.000000	0.000000
25%	125.750000	32.000000	51745.250000	25.750000	463.000000	19063.954000	0.000000	6041.500000	1.000000	1000000	30507.500000	0.000000
50%	250.500000	45.000000	84398.055556	51.000000	573.411111	29466.928889	1.000000	10754.000000	1.000000	16000000	54032.500000	0.000000
75%	375.250000	57.000000	117492.000000	77.000000	882.250000	37371.500000	2.000000	15099.750000	23.000000	78623.750000	1.000000	0.000000
max	500.000000	69.000000	149922.000000	90.000000	848.000000	49936.000000	2.000000	19990.000000	29.000000	99835.000000	1.000000	1.000000

```
In [21]: df.isnull().sum()

Out[21]: Customer_ID    0
Age                  0
Gender               0
Income              0
Spending_Score      0
Credit_Score        0
Loan_Amount         0
Previous_Defaults    0
Marketing_Spend      0
Purchase_Frequency   0
Seasonality          0
Sales               0
Customer_Churn       0
Defaulted            0
Seasonality_Date     0
dtype: object

In [31]: # J212: Missing value by mean and double checked
df[["Income", "Credit_Score", "Loan_Amount"]] = df[["Income", "Credit_Score", "Loan_Amount"]].fillna(df[["Income", "Credit_Score", "Loan_Amount"]].mean())
print(df.dtypes)

Out[31]: Customer_ID    0
Age                  0
Gender               0
Income              0
Spending_Score      0
Credit_Score        0
Loan_Amount         0
Previous_Defaults    0
Marketing_Spend      0
Purchase_Frequency   0
Seasonality          0
Sales               0
Customer_Churn       0
Defaulted            0
Seasonality_Date     0
dtype: object

In [41]: df.dtypes

Out[41]: Age                  int64
Income              float64
Spending_Score      float64
Credit_Score        float64
Loan_Amount         float64
Previous_Defaults    int64
Marketing_Spend      int64
Purchase_Frequency   int64
Seasonality          object
Sales               int64
Customer_Churn       int64
Defaulted            int64
Seasonality_Date     datetime64[ns]
dtype: object

In [51]: df.info()

Out[51]: <class 'pandas.core.frame.DataFrame'>
Int64Index: 500 entries, 0 to 499
Data columns (total 13 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Customer_ID           500 non-null    int64
 1   Age                   500 non-null    int64
 2   Gender                500 non-null    object
 3   Income                500 non-null    float64
 4   Spending_Score        500 non-null    float64
 5   Credit_Score          500 non-null    float64
 6   Loan_Amount           500 non-null    float64
 7   Previous_Defaults     500 non-null    int64
 8   Marketing_Spend       500 non-null    int64
 9   Purchase_Frequency    500 non-null    int64
10   Seasonality           500 non-null    object
11   Sales                 500 non-null    int64
12   Customer_Churn        500 non-null    int64
13   Defaulted             500 non-null    int64
14   Seasonality_Date      500 non-null    object
dtypes: float64(5), int64(9), object(3)
memory usage: 58.74 KB

In [71]: df["Seasonality_Date"] = pd.to_datetime(df["Seasonality_Date"], format="%d/%m/%Y")
print(df.dtypes)

Out[71]: Customer_ID           int64
Age                   int64
Gender                object
Income               float64
Spending_Score       float64
Credit_Score         float64
Loan_Amount          float64
Previous_Defaults     int64
Marketing_Spend       int64
Purchase_Frequency    int64
Seasonality           object
Sales                int64
Customer_Churn        int64
Defaulted             int64
Seasonality_Date      datetime64[ns]
dtype: object

In [111]: import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style("whitegrid")

fig, axes = plt.subplots(3, 1)
plt.figure(figsize=(15, 3))
plt.subplot(3, 1, 1)
plt.title("Boxplot of Age")
plt.xlabel("Age")
plt.ylabel("Count")

plt.figure(figsize=(15, 3))
plt.subplot(3, 1, 2)
plt.title("Boxplot of Income")
plt.xlabel("Income")
plt.ylabel("Count")

plt.figure(figsize=(15, 3))
plt.subplot(3, 1, 3)
plt.title("Boxplot of Spending_Score")
plt.xlabel("Spending_Score")
plt.ylabel("Count")

plt.figure(figsize=(15, 3))
plt.subplot(3, 1, 4)
plt.title("Boxplot of Credit_Score")
plt.xlabel("Credit_Score")
plt.ylabel("Count")

plt.figure(figsize=(15, 3))
plt.subplot(3, 1, 5)
plt.title("Boxplot of Loan_Amount")
plt.xlabel("Loan_Amount")
plt.ylabel("Count")

plt.figure(figsize=(15, 3))
plt.subplot(3, 1, 6)
plt.title("Boxplot of Marketing_Spend")
plt.xlabel("Marketing_Spend")
plt.ylabel("Count")

plt.figure(figsize=(15, 3))
plt.subplot(3, 1, 7)
plt.title("Boxplot of Sales")
plt.xlabel("Sales")
plt.ylabel("Count")

In [121]: numeric_df = df.select_dtypes(include='number')
Q1 = numeric_df.quantile(0.25)
Q3 = numeric_df.quantile(0.75)
IQR = Q3 - Q1
df = df[(numeric_df < (Q1 - 1.5 * IQR)) | (numeric_df > (Q3 + 1.5 * IQR))].any(axis=1)
print("After removing outliers, rows remaining:", df.shape[0])

After removing outliers, rows remaining: 405

In [141]: df.to_csv("cleaned_data.csv", index=False)
print("Data cleaning complete and saved as 'cleaned_data.csv'")

Data cleaning complete and saved as 'cleaned_data.csv'

In [211]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style("whitegrid")

df = pd.read_csv("cleaned_data.csv")
print("Shape of data:", df.shape)
print("Data description:")
print(df.describe())

Shape of data: (405, 15)

Basic information:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 405 entries, 0 to 404
Data columns (total 13 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Customer_ID           405 non-null    int64
 1   Age                   405 non-null    int64
 2   Gender                405 non-null    object
 3   Income                405 non-null    float64
 4   Spending_Score        405 non-null    float64
 5   Credit_Score          405 non-null    float64
 6   Loan_Amount           405 non-null    float64
 7   Previous_Defaults     405 non-null    int64
 8   Marketing_Spend       405 non-null    int64
 9   Purchase_Frequency    405 non-null    int64
10   Seasonality           405 non-null    object
11   Sales                 405 non-null    int64
12   Customer_Churn        405 non-null    int64
13   Defaulted             405 non-null    int64
14   Seasonality_Date      405 non-null    object
dtypes: float64(5), int64(9), object(3)
memory usage: 47.64 KB
None

Summary statistics:
Customer_ID    Age    Income    Spending_Score    Credit_Score \
count    405.000000    405.000000    405.000000    405.000000    405.000000
mean      249.123086    44.18148    84311.487490    50.763432    572.757331
std      141.444837    14.727148    37953.082256    19.648480    148.545319
min         1.000000    18.000000    20055.000000    1.000000    300.000000
25%      127.000000    32.000000    51745.250000    25.000000    463.000000
50%      250.000000    45.000000    84398.055556    51.000000    573.411111
75%      375.000000    57.000000    117492.000000    77.000000    882.000000
max      500.000000    69.000000    149922.000000    99.000000    848.000000

Loan_Amount    Previous_Defaults    Marketing_Spend    Purchase_Frequency \
count    405.000000    405.000000    405.000000    405.000000
mean      28937.427193    0.897654    10668.232529    15.214815
std      11849.042234    0.848818    5493.911048    8.432289
min         1563.000000    0.000000    1024.000000    1.000000
25%      13956.000000    0.000000    6187.000000    1.000000
50%      28936.000000    1.000000    10755.000000    1.000000
75%      38049.000000    2.000000    15376.000000    23.000000
max      49936.000000    2.000000    19990.000000    29.000000

Sales    Customer_Churn    Defaulted
count    405.000000    405.000000    405.00
mean      14950.598028    0.184193    0.0
std      27125.289394    0.438809    0.0
min         5234.000000    0.000000    0.0
25%      30902.000000    0.000000    0.0
50%      15274.000000    0.000000    0.0
75%      78616.000000    1.000000    0.0
max      99835.000000    1.000000    0.0

In [311]: numeric_cols = ["Age", "Income", "Spending_Score", "Credit_Score", "Loan_Amount", "Marketing_Spend", "Sales"]
plt.figure(figsize=(12, 8))
for i, col in enumerate(numeric_cols, 1):
    plt.subplot(3, 3, i)
    sns.histplot(df[col], kde=True, bins=20)
    plt.title(f"Distribution of {col}")
plt.tight_layout()
plt.show()

Distribution of Age
Distribution of Income
Distribution of Spending_Score
Distribution of Credit_Score
Distribution of Loan_Amount
Distribution of Marketing_Spend
Distribution of Sales

In [511]: # Correlation Heatmap
plt.figure(figsize=(8, 5))
corr = df[numeric_cols].corr()
sns.heatmap(corr, annot=True, cmap="magma", cbar=True)
plt.title("Correlation Heatmap between Numeric Variables")
plt.show()

Correlation Heatmap between Numeric Variables
Age
Income
Spending_Score
Credit_Score
Loan_Amount
Marketing_Spend
Sales

In [711]: # Gender and Seasonality Insights
plt.figure(figsize=(8, 5))
sns.boxplot(x="Gender", y="Income", data=df)
plt.title("Income Distribution by Gender")
plt.show()

Income Distribution by Gender
Income
Gender
Female
Male

sns.boxplot(x="Seasonality", y="Sales", data=df, estimator="mean")
plt.title("Average Sales across Seasonality Levels")
plt.show()

Average Sales across Seasonality Levels
Sales
Seasonality
Low
Medium
High

In [911]: # Trend Analysis Over Time
df["Seasonality_Date"] = pd.to_datetime(df["Seasonality_Date"], dayfirst=True)
monthly_sales = df.groupby(df["Seasonality_Date"].dt.to_period("M"))["Sales"].sum()
plt.figure(figsize=(10, 5))
plt.plot(monthly_sales, marker="o", title="Monthly Sales Trend")
plt.xlabel("Month")
plt.ylabel("Total Sales")
plt.show()

Monthly Sales Trend
Total Sales
Month
Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec Jan 2025

In [1111]: # Customer Churn Analysis
plt.figure(figsize=(8, 5))
sns.boxplot(x="Customer_Churn", y="Income", data=df)
plt.title("Income vs Customer Churn")
plt.show()

Income vs Customer Churn
Income
Customer_Churn
0
1

plt.figure(figsize=(8, 5))
sns.boxplot(x="Customer_Churn", y="Income", data=df)
plt.title("Income vs Customer Churn")
plt.show()

Income vs Customer Churn
Income
Customer_Churn
0
1

plt.figure(figsize=(8, 5))
sns.boxplot(x="Defaulted", y="Credit_Score", data=df)
plt.title("Credit Score vs Loan Default")
plt.show()

Credit Score vs Loan Default
Credit_Score
Defaulted
0
1

In [1311]: # Scatter plots for strongest correlations
plt.figure(figsize=(8, 5))
sns.scatterplot(x="Income", y="Sales", data=df)
plt.title("Income vs Sales")
plt.show()

Income vs Sales
Sales
Income

plt.figure(figsize=(8, 5))
sns.scatterplot(x="Spending_Score", y="Sales", data=df)
plt.title("Spending Score vs Sales")
plt.show()

Spending Score vs Sales
Sales
Spending_Score

In [1511]: # Gender-based Numeric Comparison
plt.figure(figsize=(8, 5))
sns.boxplot(x="Gender", y="Income", data=df)
plt.title("Income by Gender")
plt.show()

Income by Gender
Income
Gender
Female
Male

plt.figure(figsize=(8, 5))
sns.boxplot(x="Gender", y="Spending_Score", data=df)
plt.title("Spending Score by Gender")
plt.show()

Spending Score by Gender
Spending_Score
Gender
Female
Male

In [1711]: # Seasonality effects
sns.boxplot(x="Seasonality", y="Sales", data=df, estimator="mean", notch=True)
plt.title("Average Sales across Seasonality Levels")
plt.show()

Average Sales across Seasonality Levels
Sales
Seasonality
Low
Medium
High

In [1911]: # Import Libraries
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create and train the model
lin_model = LinearRegression()
lin_model.fit(X_train, y_train)

# Predict
y_pred = lin_model.predict(X_test)

# Evaluate model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print("Linear Regression - Sales Prediction")
print("Mean Squared Error: ", mse)
print("R^2 Score: ", r2)

Linear Regression - Sales Prediction
Mean Squared Error: 7984871.099792
R^2 Score: 0.93740028337701

In [2111]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create and train the model
lin_model = LinearRegression()
lin_model.fit(X_train, y_train)

# Predict
y_pred = lin_model.predict(X_test)

# Evaluate model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print("Linear Regression - Sales Prediction")
print("Mean Squared Error: ", mse)
print("R^2 Score: ", r2)

Linear Regression - Sales Prediction
Mean Squared Error: 7984871.099792
R^2 Score: 0.93740028337701

In [2311]: # Feature importance
from sklearn.feature_selection import SelectFromModel
from sklearn.linear_model import LinearRegression
model = LinearRegression()
model.fit(X_train, y_train)
sfm = SelectFromModel(model)
important_features = sfm.get_support()

print("Important Features:")
print(important_features)

Important Features:
[0 1]

In [2511]: # Decision Tree
from sklearn.tree import DecisionTreeClassifier
model = DecisionTreeClassifier(max_depth=3)
model.fit(X_train, y_train)

print("Decision Tree (Depth 3) Results:")
print(model.score(X_test, y_test))

Decision Tree (Depth 3) Results:
0.81
```


Simplified Decision Tree (Top 3 Levels)

