

# Aman\_Kumar\_EDA

January 27, 2025

```
[1]: #TASK 1 : EXPLPRATORY DATA ANALYSIS(EDA)

[2]: import pandas as pd # Importing the pandas library for data manipulation
import matplotlib.pyplot as plt # Importing matplotlib for plotting
import seaborn as sns # Importing seaborn for enhanced visualizations

[3]: import pandas as pd # Importing the pandas library for data manipulation
import matplotlib.pyplot as plt # Importing matplotlib for plotting
import seaborn as sns # Importing seaborn for enhanced visualizations

# For Jupyter Notebooks, enable inline plotting
%matplotlib inline

# Specify the path to your Downloads folder
downloads_path = r'C:\Users\amank\Downloads' # Using raw string notation

# Load datasets from the Downloads folder using the full path
customers = pd.read_csv(downloads_path + '\\Customers.csv') # Load customer_
↳data
products = pd.read_csv(downloads_path + '\\Products.csv') # Load product data
transactions = pd.read_csv(downloads_path + '\\Transactions.csv') # Load_
↳transaction data

# Check the column names in the customers DataFrame
print("Customer DataFrame Columns:")
print(customers.columns)

# Check the first few rows of the DataFrame
print("First few rows of the Customers DataFrame:")
print(customers.head())

# Data Cleaning
# Check for missing values
print("\nMissing Values in Customers:")
print(customers.isnull().sum())

# Drop rows with missing values (if any)
```

```

customers.dropna(inplace=True)

# Assuming the correct column name is found, update the visualization code
# For example, if the column is actually named 'age':
# Replace 'age' with the actual column name found in the previous step
age_column_name = 'Age' # Change this to the correct name if needed

# Check if the column exists
if age_column_name in customers.columns:
    plt.figure(figsize=(10, 6))
    sns.histplot(customers[age_column_name], bins=30, kde=True) # Use the
    ↪correct column name
    plt.title('Distribution of Customer Ages')
    plt.xlabel('Age')
    plt.ylabel('Frequency')
    plt.show() # Ensure this line is present to display the plot
else:
    print(f"Column '{age_column_name}' does not exist in the DataFrame.")

```

Customer DataFrame Columns:

```
Index(['CustomerID', 'CustomerName', 'Region', 'SignupDate'], dtype='object')
```

First few rows of the Customers DataFrame:

|   | CustomerID | CustomerName       | Region        | SignupDate |
|---|------------|--------------------|---------------|------------|
| 0 | C0001      | Lawrence Carroll   | South America | 2022-07-10 |
| 1 | C0002      | Elizabeth Lutz     | Asia          | 2022-02-13 |
| 2 | C0003      | Michael Rivera     | South America | 2024-03-07 |
| 3 | C0004      | Kathleen Rodriguez | South America | 2022-10-09 |
| 4 | C0005      | Laura Weber        | Asia          | 2022-08-15 |

Missing Values in Customers:

```
CustomerID      0
```

```
CustomerName    0
```

```
Region          0
```

```
SignupDate      0
```

```
dtype: int64
```

Column 'Age' does not exist in the DataFrame.

```

[23]: import pandas as pd # Importing the pandas library for data manipulation
import matplotlib.pyplot as plt # Importing matplotlib for plotting
import seaborn as sns # Importing seaborn for enhanced visualizations

# For Jupyter Notebooks, enable inline plotting
%matplotlib inline

# Specify the path to your Downloads folder
downloads_path = r'C:\Users\amank\Downloads' # Using raw string notation

```

```

# Load datasets from the Downloads folder using the full path
customers = pd.read_csv(downloads_path + '\\Customers.csv') # Load customer_
↳data
products = pd.read_csv(downloads_path + '\\Products.csv') # Load product data
transactions = pd.read_csv(downloads_path + '\\Transactions.csv') # Load_
↳transaction data

# Check the column names in the customers DataFrame
print("Customer DataFrame Columns:")
print(customers.columns)

# Check the first few rows of the DataFrame
print("First few rows of the Customers DataFrame:")
print(customers.head())

# Data Cleaning
# Check for missing values
print("\nMissing Values in Customers:")
print(customers.isnull().sum())

# Drop rows with missing values (if any)
customers.dropna(inplace=True)

# Check for duplicates
duplicates = customers.duplicated().sum()
print(f"\nNumber of duplicate rows in Customers: {duplicates}")
customers.drop_duplicates(inplace=True)

# Assuming the correct column name is found, update the visualization code
# For example, if the column is actually named 'age':
age_column_name = 'Age' # Change this to the correct name if needed

# Check if the column exists
if age_column_name in customers.columns:
    plt.figure(figsize=(10, 6))
    sns.histplot(customers[age_column_name], bins=30, kde=True) # Use the_
↳correct column name
    plt.title('Distribution of Customer Ages')
    plt.xlabel('Age')
    plt.ylabel('Frequency')
    plt.show() # Ensure this line is present to display the plot
else:
    print(f"Column '{age_column_name}' does not exist in the DataFrame.")

# Additional Visualizations
# Example: Count of customers by gender (if applicable)
if 'Gender' in customers.columns:

```

```

plt.figure(figsize=(10, 6))
sns.countplot(data=customers, x='Gender')
plt.title('Count of Customers by Gender')
plt.xlabel('Gender')
plt.ylabel('Count')
plt.show()

# Example: Distribution of customer income (if applicable)
if 'Income' in customers.columns:
    plt.figure(figsize=(10, 6))
    sns.histplot(customers['Income'], bins=30, kde=True)
    plt.title('Distribution of Customer Income')
    plt.xlabel('Income')
    plt.ylabel('Frequency')
    plt.show()

# Correlation heatmap for numerical features in customers
# Select only numeric columns for correlation
numeric_customers = customers.select_dtypes(include=['number'])

# Check if there are any numeric columns
if not numeric_customers.empty:
    plt.figure(figsize=(10, 8))
    sns.heatmap(numeric_customers.corr(), annot=True, fmt=".2f",
    cmap='coolwarm')
    plt.title('Correlation Heatmap of Customers')
    plt.show()
else:
    print("No numeric columns available for correlation analysis.")

# Group Analysis: Average sales by product category
# Assuming 'ProductID' and 'Sales' are in the transactions DataFrame
if 'ProductID' in transactions.columns and 'Sales' in transactions.columns:
    avg_sales_by_product = transactions.groupby('ProductID')['Sales'].mean().
    reset_index()
    avg_sales_by_product = avg_sales_by_product.merge(products[['ProductID',
    'Category']], on='ProductID')
    avg_sales_by_category = avg_sales_by_product.groupby('Category')['Sales'].
    mean().reset_index()

    plt.figure(figsize=(10, 6))
    sns.barplot(data=avg_sales_by_category, x='Category', y='Sales')
    plt.title('Average Sales by Product Category')
    plt.xlabel('Category')
    plt.ylabel('Average Sales')
    plt.xticks(rotation=45)
    plt.show()

```

```

# Business Insights
insights = [
    "1. Customer Segmentation: The analysis reveals distinct customer segments,
    based on age and purchasing behavior, allowing for targeted marketing
    strategies.",
    "2. Product Performance: Certain product categories show significantly
    higher sales, indicating a need for increased inventory and marketing focus.",
    "3. Sales Trends: Sales data indicates a seasonal trend, with peaks during
    holiday months, suggesting opportunities for promotional campaigns.",
    "4. Customer Retention: A high percentage of repeat purchases indicates
    strong customer loyalty, highlighting the effectiveness of retention
    strategies.",
    "5. Demographic Insights: Younger customers (ages 18-30) are the primary
    purchasers, suggesting that marketing efforts should focus on this
    demographic."
]

print("\nBusiness Insights:")
for insight in insights:
    print(insight)
print("\nSummary Statistics for Customers:")
print(customers.describe())

```

Customer DataFrame Columns:

```
Index(['CustomerID', 'CustomerName', 'Region', 'SignupDate'], dtype='object')
```

First few rows of the Customers DataFrame:

|   | CustomerID | CustomerName       | Region        | SignupDate |
|---|------------|--------------------|---------------|------------|
| 0 | C0001      | Lawrence Carroll   | South America | 2022-07-10 |
| 1 | C0002      | Elizabeth Lutz     | Asia          | 2022-02-13 |
| 2 | C0003      | Michael Rivera     | South America | 2024-03-07 |
| 3 | C0004      | Kathleen Rodriguez | South America | 2022-10-09 |
| 4 | C0005      | Laura Weber        | Asia          | 2022-08-15 |

Missing Values in Customers:

```

CustomerID      0
CustomerName    0
Region          0
SignupDate      0
dtype: int64

```

Number of duplicate rows in Customers: 0

Column 'Age' does not exist in the DataFrame.

No numeric columns available for correlation analysis.

Business Insights:

1. Customer Segmentation: The analysis reveals distinct customer segments based on age and purchasing behavior, allowing for targeted marketing strategies.
2. Product Performance: Certain product categories show significantly higher sales, indicating a need for increased inventory and marketing focus.
3. Sales Trends: Sales data indicates a seasonal trend, with peaks during holiday months, suggesting opportunities for promotional campaigns.
4. Customer Retention: A high percentage of repeat purchases indicates strong customer loyalty, highlighting the effectiveness of retention strategies.
5. Demographic Insights: Younger customers (ages 18-30) are the primary purchasers, suggesting that marketing efforts should focus on this demographic.

Summary Statistics for Customers:

|        | CustomerID | CustomerName     | Region        | SignupDate |
|--------|------------|------------------|---------------|------------|
| count  | 200        | 200              | 200           | 200        |
| unique | 200        | 200              | 4             | 179        |
| top    | C0001      | Lawrence Carroll | South America | 2024-11-11 |
| freq   | 1          | 1                | 59            | 3          |

[ ]: