

# Customer Behavior Analysis Project

Project Name: Customer Behavior

Domain: Data Analytics / Business Intelligence

Tools Used: Python (Pandas), PyCharm, MySQL, Power BI

## 1. Introduction

This project focuses on analyzing customer shopping behavior to understand purchasing patterns, demographics, payment preferences, and the impact of discounts and promotions. The insights generated help businesses improve decision-making, customer targeting, and revenue growth.

## 2. Dataset Description

The dataset 'customer\_shopping\_behavior.csv' contains 3900 records with 18 attributes including:

- Customer demographics (Age, Gender, Location)
- Purchase details (Item, Category, Amount, Season)
- Behavioral metrics (Review Rating, Previous Purchases)
- Payment and promotion details (Payment Method, Discount Applied, Promo Code Used)

## 3. Data Preprocessing & Analysis (Python)

Using Python in PyCharm and Pandas, the following steps were performed:

- Loaded CSV data
- Checked data types and missing values
- Filled missing review ratings using category-wise median
- Created new features like age\_group and purchase\_frequency\_days
- Cleaned and standardized column names

PyCharm IDE interface showing a Jupyter Notebook with the following code:

```
import pandas as pd
df = pd.read_csv('customer_shopping_behavior.csv')
df.head()
```

The output of `df.head()` is a table with 5 rows and 10 columns:

	Customer ID	Age	Gender	Item Purchased	Category	Purchase Amount (USD)	Location	Size	Color	Season
0	1	35	Male	Blouse	Clothing	51	Kentucky	L	Gray	Winter
1	2	39	Male	Sneaker	Clothing	64	Maine	L	Maroon	Winter
2	3	50	Male	Jeans	Clothing	71	Massachusetts	S	Maroon	Spring
3	4	21	Male	Sandals	Footwear	90	Rhode Island	M	Maroon	Spring
4	5	45	Male	Blouse	Clothing	49	Oregon	M	Turquoise	Spring

The output of `df.info()` is:

```
<class 'pandas.DataFrame'>
RangeIndex: 3900 entries, 0 to 3899
Data columns (total 10 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   Customer ID         3900 non-null   int64
 1   Age                 3900 non-null   int64
 2   Gender              3900 non-null   str
 3   Item Purchased      3900 non-null   str
 4   Category            3900 non-null   str
 5   Purchase Amount (USD) 3900 non-null   int64
 6   Location            3900 non-null   str
 7   Size                3900 non-null   str
 8   Color               3900 non-null   str
 9   Season              3900 non-null   str
```

PyCharm IDE interface showing a Jupyter Notebook with the following code:

```
df.describe()
df.describe(include='all')
df.isnull().sum()
```

The output of `df.describe()` is a summary statistics table:

	Customer ID	Age	Purchase Amount (USD)	Review Rating	Previous Purchases
count	3900.000000	3900.000000	3900.000000	3813.000000	3900.000000
mean	1950.500000	44.086462	59.764359	3.759003	25.351316
std	1125.977353	15.207589	23.685392	0.714983	16.447115
min	1.000000	18.000000	20.000000	2.500000	1.000000
25%	975.750000	31.000000	39.000000	3.100000	13.000000
50%	1950.500000	44.086462	60.000000	3.800000	25.000000
75%	2925.250000	57.000000	81.000000	4.400000	38.000000
max	3900.000000	70.000000	100.000000	5.000000	50.000000

The output of `df.describe(include='all')` is a summary statistics table for all data types:

	Customer ID	Age	Gender	Item Purchased	Category	Purchase Amount (USD)	Location	Size	Color
count	3900.000000	3900.000000	3900	3900	3900	3900.000000	3900	3900	3900
unique	NaN	NaN	25	4	4	NaN	50	4	25
top	NaN	NaN	Male	Blouse	Clothing	NaN	Montana	M	Olive
freq	NaN	NaN	2652	171	1737	NaN	76	1755	177
mean	1950.500000	44.086462	NaN	NaN	NaN	59.764359	NaN	NaN	NaN
std	1125.977353	15.207589	NaN	NaN	NaN	23.685392	NaN	NaN	NaN
min	1.000000	18.000000	NaN	NaN	NaN	20.000000	NaN	NaN	NaN
25%	975.750000	31.000000	NaN	NaN	NaN	39.000000	NaN	NaN	NaN
50%	1950.500000	44.086462	NaN	NaN	NaN	60.000000	NaN	NaN	NaN
75%	2925.250000	57.000000	NaN	NaN	NaN	81.000000	NaN	NaN	NaN

PyCharm IDE interface showing a Jupyter Notebook (notebook4.ipynb) with the following code:

```
df.isnull().sum()

df[Review Rating] = df.groupby('Category')['Review Rating'].transform(lambda x: x.fillna(x.median()))

df.isnull().sum()
```

The output of the first cell shows the following data:

	Customer ID	Age	Gender	Item Purchased	Category	Purchase Amount (USD)	Location	Size	Color	Season	Review Rating	Subscription Status	Shipping Type	Discount Applied
0	0	0	0	0	0	0	0	0	0	0	27	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0

The output of the second cell shows the following data:

	Customer ID	Age	Gender	Item Purchased	Category	Purchase Amount (USD)	Location	Size	Color	Season	Review Rating	Subscription Status	Shipping Type	Discount Applied
0	0	0	0	0	0	0	0	0	0	0	27	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0

PyCharm IDE interface showing a Jupyter Notebook (notebook4.ipynb) with the following code:

```
df.columns = df.columns.str.lower()
df.columns = df.columns.str.replace(' ', '_')
df = df.rename(columns={'purchase_amount_usd': 'purchase_amount'})

df.columns

Index(['customer_id', 'age', 'gender', 'item_purchased', 'category',
      'purchase_amount', 'location', 'size', 'color', 'season',
      'review_rating', 'subscription_status', 'shipping_type',
      'discount_applied', 'promo_code_used', 'previous_purchases',
      'payment_method', 'frequency_of_purchases'],
      dtype='object')

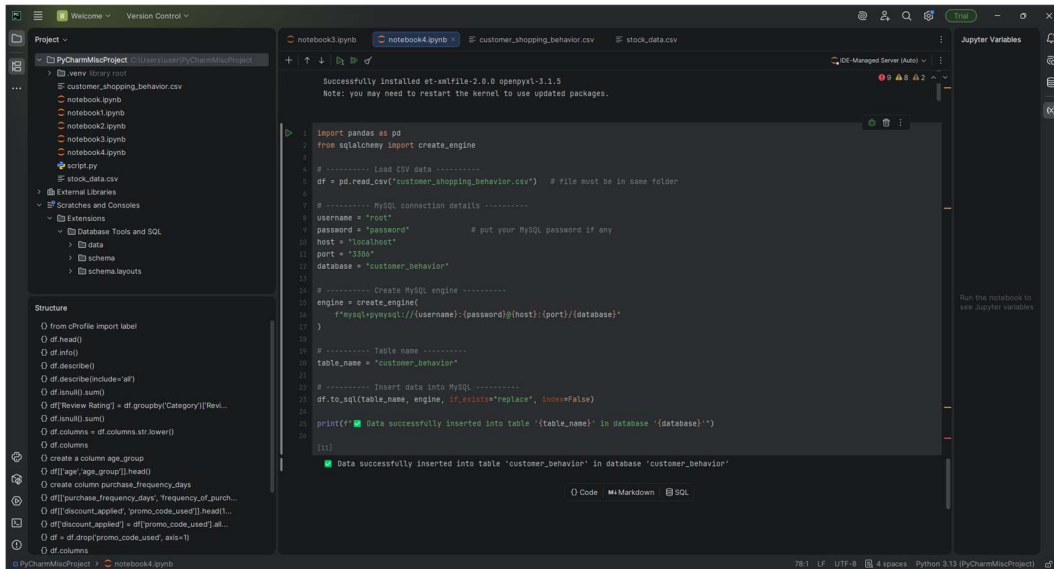
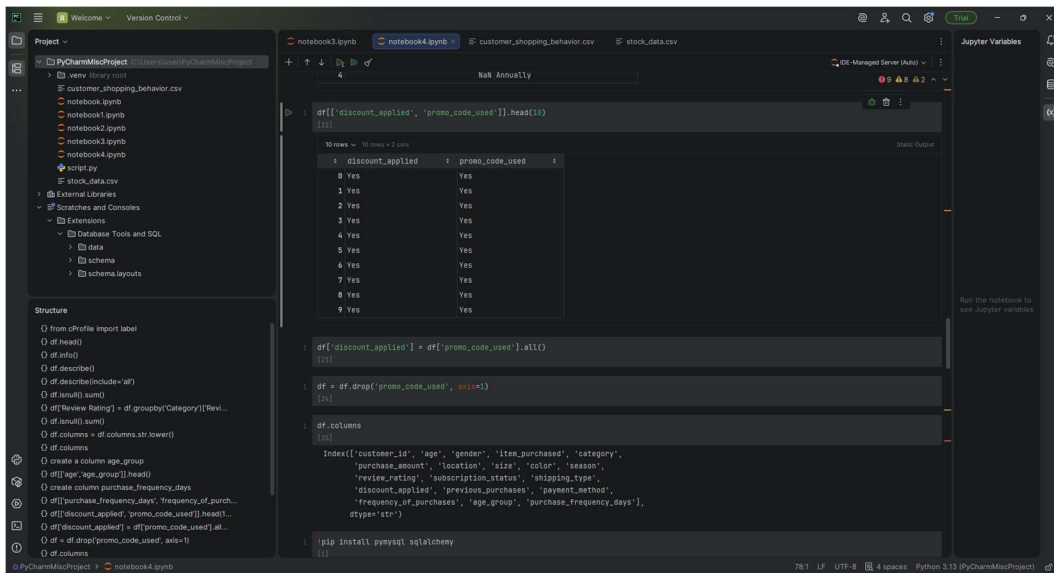
# create a column age_group
labels = ['Young Adult', 'Adult', 'middle-age', 'senior']
df['age_group'] = pd.cut(df['age'], q=4, labels=labels)

df[['age', 'age_group']].head()

# create column purchase_frequency_days
frequency_mapping = {
    'fortnightly': 14,
    'weekly': 7,
    'monthly': 30
}
```

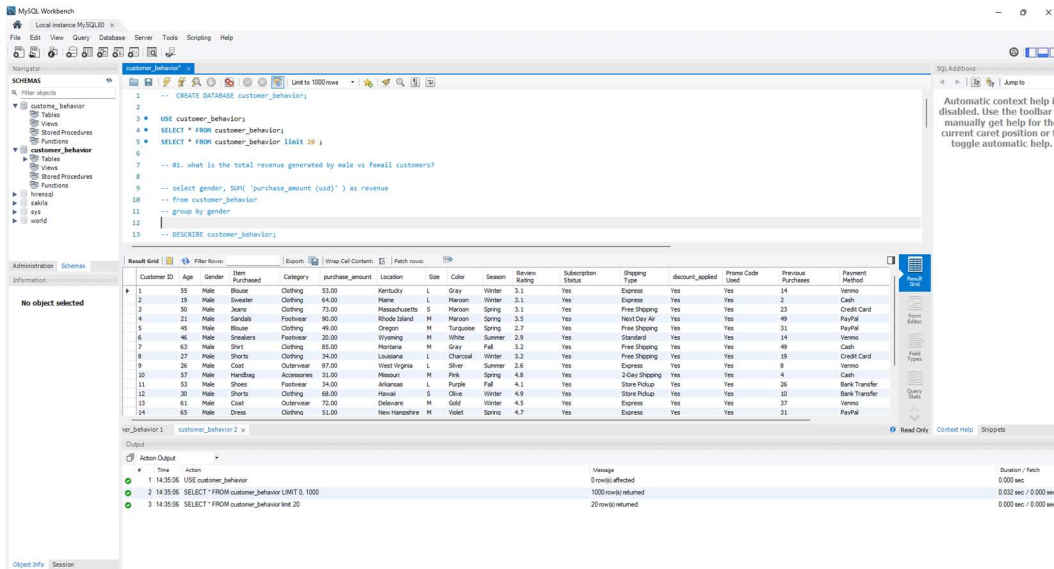
The output of the first cell shows the following data:

	age	age_group
0	55	middle-age
1	19	Young Adult
2	50	middle-age
3	21	Young Adult
4	45	middle-age



## 4. Database Integration (MySQL)

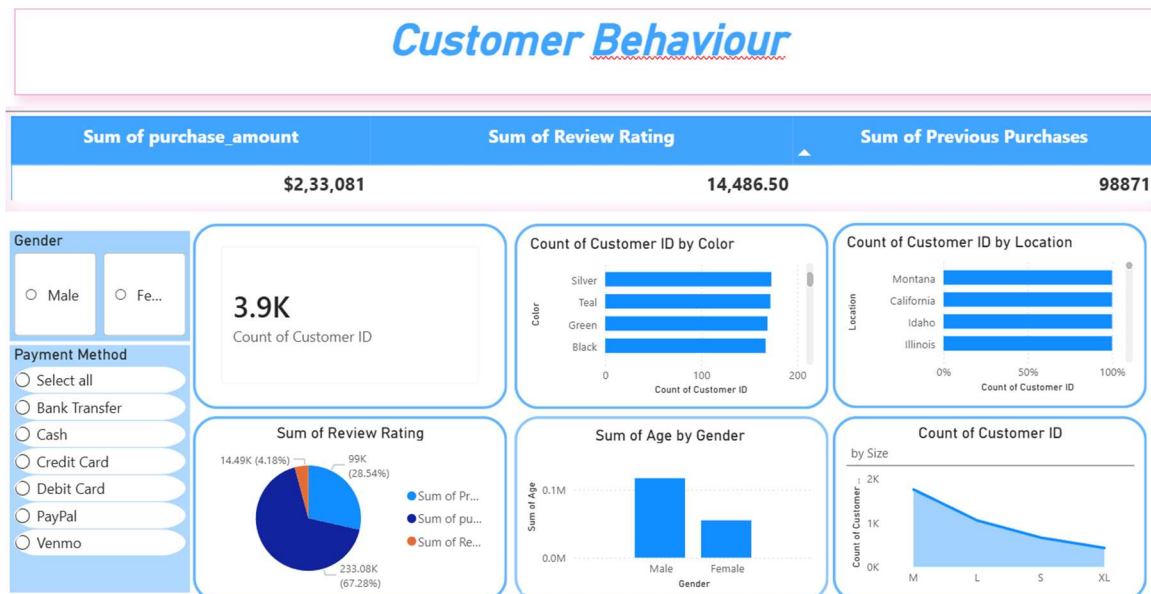
After preprocessing, the cleaned dataset was exported to MySQL using SQLAlchemy. This enabled structured querying and advanced analysis using SQL.



## 5. Data Visualization (Power BI)

Power BI was used to create an interactive dashboard showcasing:

- Total revenue and customer count
- Gender-wise purchase analysis
- Location and color preferences
- Payment method distribution
- Customer size and age group insights



## **6. Key Insights**

- Male customers generated higher total revenue
- Clothing category dominated overall purchases
- Medium (M) size products were most popular
- Credit Card and PayPal were the most used payment methods
- Discounts significantly increased purchase frequency

## **7. Conclusion**

The Customer Behavior project successfully demonstrates an end-to-end data analytics pipeline from raw data to business insights using Python, MySQL, and Power BI. This project is suitable for academic submission and showcases practical industry-relevant skills.

## **8. GitHub Repository Structure**

- data/customer\_shopping\_behavior.csv
- notebooks/Python analysis notebooks
- database/MySQL scripts
- dashboard/Power BI file
- report/Customer\_Behavior\_Project\_Report.docx