# Customer Shopping Behavior Analysis

## 1. Project Overview:

This project analyzes customer shopping behavior using transactional data from 3,900 purchases across various product categories. The goal is to uncover insights into spending patterns, customer segments, product preferences, and subscription behavior to guide strategic business decisions.

## 2. Dataset Summary:

- Rows: 3,900
- Columns: 18
- Key Features: - Customer demographics (Age, Gender, Location, Subscription Status)
- Purchase details (Item Purchased, Category, Purchase Amount, Season, Size, Color)
- Shopping behavior (Discount Applied, Promo Code Used, Previous Purchases, Frequency of Purchases, Review Rating, Shipping Type)
- Missing Data: 37 values in Review Rating column

## 3. Exploratory Data Analysis using Python:

We began with data preparation and cleaning in Python:

**Data Loading**: Imported the dataset using pandas.

**Initial Exploration**: Used df.info() to check structure and .describe() for summary statistics

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3900 entries, 0 to 3899
Data columns (total 18 columns):
 #   Column                Non-Null Count  Dtype
---  ------                --------------  -----
 0   Customer ID           3900 non-null   int64
 1   Age                   3900 non-null   int64
 2   Gender                3900 non-null   object
 3   Item Purchased        3900 non-null   object
 4   Category              3900 non-null   object
 5   Purchase Amount (USD) 3900 non-null   int64
 6   Location              3900 non-null   object
 7   Size                  3900 non-null   object
 8   Color                 3900 non-null   object
 9   Season                3900 non-null   object
 10  Review Rating         3863 non-null   float64
 11  Subscription Status   3900 non-null   object
 12  Shipping Type         3900 non-null   object
 13  Discount Applied      3900 non-null   object
 14  Promo Code Used       3900 non-null   object
 15  Previous Purchases    3900 non-null   int64
 16  Payment Method        3900 non-null   object
 17  Frequency of Purchases 3900 non-null  object
dtypes: float64(1), int64(4), object(13)
memory usage: 548.6+ KB
```

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

| | Customer ID | Age | Gender | Item Purchased | Category | Purchase Amount (USD) | Location | Size | Color | Season | Review Rating | Subscription Status | Shipping Type | Discount Applied | Promo Code Used |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 3900.000000 | 3900.000000 | 3900 | 3900 | 3900 | 3900.000000 | 3900 | 3900 | 3900 | 3900 | 3863.000000 | 3900 | 3900 | 3900 | 3900 |
| unique | NaN | NaN | 2 | 25 | 4 | NaN | 50 | 4 | 25 | 4 | NaN | 2 | 6 | 2 | 2 |
| top | NaN | NaN | Male | Blouse | Clothing | NaN | Montana | M | Olive | Spring | NaN | No | Free Shipping | No | No |
| freq | NaN | NaN | 2652 | 171 | 1737 | NaN | 96 | 1755 | 177 | 999 | NaN | 2847 | 675 | 2223 | 2223 |
| mean | 1950.500000 | 44.068462 | NaN | NaN | NaN | 59.764359 | NaN | NaN | NaN | NaN | 3.750065 | NaN | NaN | NaN | NaN |
| std | 1125.977353 | 15.207589 | NaN | NaN | NaN | 23.685392 | NaN | NaN | NaN | NaN | 0.716983 | NaN | NaN | NaN | NaN |
| min | 1.000000 | 18.000000 | NaN | NaN | NaN | 20.000000 | NaN | NaN | NaN | NaN | 2.500000 | NaN | NaN | NaN | NaN |
| 25% | 975.750000 | 31.000000 | NaN | NaN | NaN | 39.000000 | NaN | NaN | NaN | NaN | 3.100000 | NaN | NaN | NaN | NaN |
| 50% | 1950.500000 | 44.000000 | NaN | NaN | NaN | 60.000000 | NaN | NaN | NaN | NaN | 3.800000 | NaN | NaN | NaN | NaN |
| 75% | 2925.250000 | 57.000000 | NaN | NaN | NaN | 81.000000 | NaN | NaN | NaN | NaN | 4.400000 | NaN | NaN | NaN | NaN |
| max | 3900.000000 | 70.000000 | NaN | NaN | NaN | 100.000000 | NaN | NaN | NaN | NaN | 5.000000 | NaN | NaN | NaN | NaN |

**Missing Data Handling**: Checked for null values and imputed missing values in the Review Rating column using the median rating of each product category.

**Column Standardization**: Renamed columns to snake case for better readability and documentation.

**Feature Engineering**:
Created age_group column by binning customer ages.
Created purchase_frequency_days column from purchase data.

**Data Consistency Check**: Verified if discount_applied and promo_code_used were redundant; dropped promo_code_used.

**Database Integration**: Saved the final cleaned DataFrame and connected Python script to MySQL and loaded the cleaned DataFrame into the database for SQL analysis.

```python
from sqlalchemy import create_engine

# Replace with your actual MySQL Workbench details
username = 'your_username'
password = 'your_password'
host = 'localhost:XXXX'
database = 'your_database'

# Create connection string
engine = create_engine(f"mysql+mysqlconnector://{username}:{password}@{host}/{database}")

# Export the DataFrame
df.to_sql('cleaned_customer_data', con=engine, index=False, if_exists='replace')

print("DataFrame exported successfully to MySQL database!")
```

```
DataFrame exported successfully to MySQL database!
```

## 4. Data Analysis using SQL (Business Transactions):

We performed structured analysis in MySQL to answer key business questions:

**1. Revenue by Gender** – Compared total revenue generated by male vs. female customers.

| gender | revenue |
|--------|---------|
| ▶ Male | 157890 |
| Female | 75191 |

**2. High Spending Discount Users** – Identified customers who used discounts but still spent above the average purchase amount.

| customer_id | purchase_amount |
|-------------|-----------------|
| ▶ 2 | 64 |
| 3 | 73 |
| 4 | 90 |
| 7 | 85 |
| 9 | 97 |
| 12 | 68 |
| 13 | 72 |

**3. Top 5 Products by Rating** – Found products with the highest average review ratings.

| products | Average product rating |
|----------|------------------------|
| ▶ Blouse | 3.68 |
| Sweater | 3.76 |
| Jeans | 3.65 |
| Sandals | 3.84 |
| Sneakers | 3.76 |

**4. Shipping Type Comparison** – Compared average purchase amounts between Standard and Express shipping.

| shipping_type | Avg_purchase_amount |
|---------------|---------------------|
| ▶ Express | 60.48 |
| Standard | 58.46 |

**5. Subscribers vs. Non-Subscribers** – Compared average spend and total revenue across subscription status.

| subscription_status | Total_customers | Total_amount | Avg_purchase_amount |
|---|---|---|---|
| Yes | 1053 | 62645 | 59.4919 |
| No | 2847 | 170436 | 59.8651 |

**6. Discount-Dependent Products** – Identified 5 products with the highest percentage of discounted purchases.

| product | percentage |
|---|---|
| Hat | 50.00 |
| Sneakers | 49.66 |
| Coat | 49.07 |
| Sweater | 48.17 |
| Pants | 47.37 |

**7. Customer Segmentation** – Classified customers into New, Returning, and Loyal segments based on purchase history.

| customer_segment | Number of customers |
|---|---|
| Loyal | 3116 |
| Returning | 701 |
| New | 83 |

**8. Top 3 Products per Category** – Listed the most purchased products within each category.

| category | product | total_products |
|---|---|---|
| Accessories | Jewelry | 171 |
| Accessories | Sunglasses | 161 |
| Accessories | Belt | 161 |
| Clothing | Blouse | 171 |
| Clothing | Pants | 171 |
| Clothing | Shirt | 169 |
| Footwear | Sandals | 160 |
| Footwear | Shoes | 150 |
| Footwear | Sneakers | 145 |
| Outerwear | Jacket | 163 |
| Outerwear | Coat | 161 |

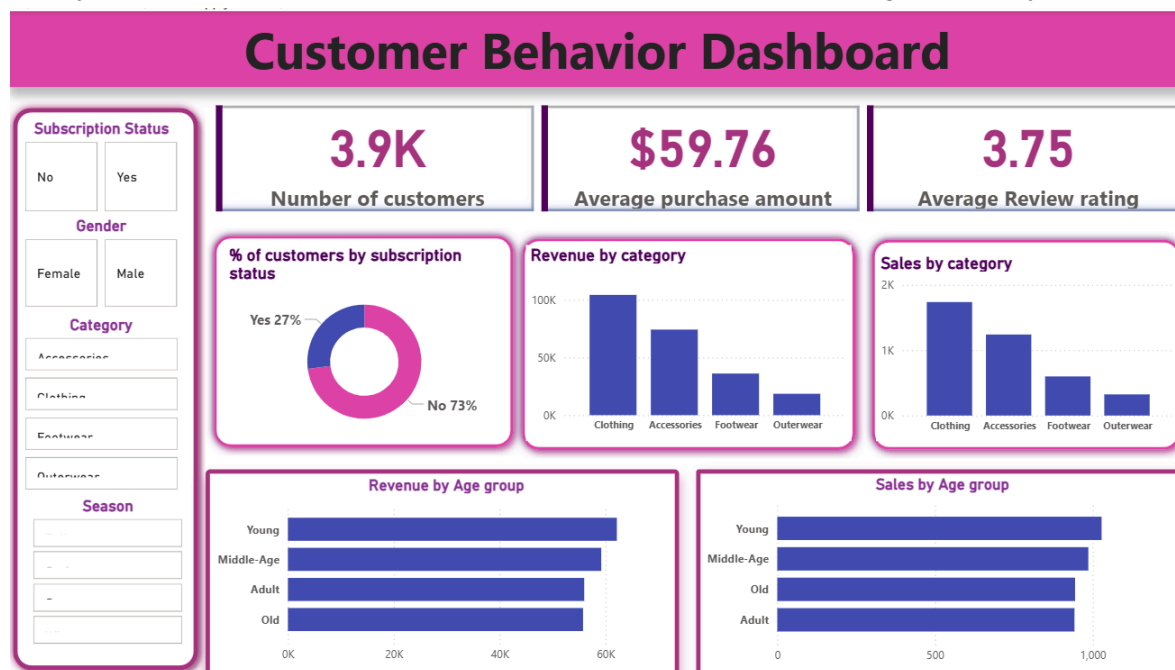**9. Repeat Buyers & Subscriptions** – Checked whether customers with >5 purchases are more likely to subscribe.

| subscription_status | repeat_buyers |
|---|---|
| ▶ Yes | 958 |
| No | 2518 |

**10. Revenue by Age Group** – Calculated total revenue contribution of each age group.

| age_group | total_revenue |
|---|---|
| ▶ Middle-Age | 59197 |
| Young | 62143 |
| Old | 55763 |
| Adult | 55978 |

**5. Dashboard in Power BI:**

Finally, we built an interactive dashboard in Power BI to present insights visually.

**6. Business Recommendations**:

 **Boost Subscriptions** – Promote exclusive benefits for subscribers.
 **Customer Loyalty Programs** – Reward repeat buyers to move them into the "Loyal" segment.
 **Review Discount Policy** – Balance sales boosts with margin control.
 **Product Positioning** – Highlight top rated and best selling products in campaigns.
 **Targeted Marketing** – Focus efforts on high revenue age groups and express-shipping users.