

Customer Shopping Behavior Analysis

1. Project Overview

This project analyses 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.

```
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)
```

- **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:** Connected Python script to MS SQL Server and loaded the cleaned DataFrame into the database for SQL analysis.

4. Data Analysis using SQL (Business Transactions)

We performed structured analysis in MS SQL to answer key business questions:

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

	gender	revenue
1	Male	157890
2	Female	75191

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

	customer_id	discount_applied	purchase_amount
1	2	Yes	64
2	3	Yes	73
3	4	Yes	90
4	7	Yes	85
5	9	Yes	97
6	12	Yes	68
7	13	Yes	72
8	16	Yes	81
9	20	Yes	90
10	22	Yes	62
11	24	Yes	88

(863 rows affected)

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

	item_purchased	Category	Average Product Rating
1	Gloves	Accessories	3.86
2	Sandals	Footwear	3.84
3	Boots	Footwear	3.82
4	Hat	Accessories	3.8
5	Skirt	Clothing	3.78

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

	shipping_type	Average Purchase Amount
1	Standard	58
2	Express	60

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

	subscription_status	Total Customers	Average Spend	Total Revenue
1	Yes	1053	59	62645
2	No	2847	59	170436

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

	item_purchased	Category	Discount_rate
1	Hat	Accessories	50%
2	Sneakers	Footwear	49%
3	Coat	Outerwear	49%
4	Sweater	Clothing	48%
5	Pants	Clothing	47%

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

	Customer_segment	Number of Customers
1	Loyal	3116
2	Returning	701
3	New	83

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

	item_rank	category	item_purchased	total_orders
1	1	Accessories	Jewelry	171
2	2	Accessories	Belt	161
3	3	Accessories	Sunglasses	161
4	1	Clothing	Blouse	171
5	2	Clothing	Pants	171
6	3	Clothing	Shirt	169
7	1	Footwear	Sandals	160
8	2	Footwear	Shoes	150
9	3	Footwear	Sneakers	145
10	1	Outerwear	Jacket	163
11	2	Outerwear	Coat	161

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

	Subscription_status	Repeat_buyers
1	Yes	958
2	No	2518

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

	age_group	Total_Revenue
1	Young Adult	62143
2	Middle-Aged	59197
3	Adult	55978
4	Senior	55763

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.