

# Customer Shopping Behavior Analysis

## 1. Project Overview

This project explores customer shopping patterns using transactional data from 3,900 purchases across multiple product categories. The objective is to identify insights on spending habits, customer segmentation, product preferences, and subscription trends to support strategic business decisions.

## 2. Dataset Summary

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

## 3. Exploratory Data Analysis using Python (Pandas)

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 `.describer()` for summary statistics.

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

  

Promo Code Used	Previous Purchases	Payment Method	Frequency of Purchases
3900	3900.000000	3900	3900
2	Nan	6	7
No	Nan	PayPal	Every 3 Months
2223	Nan	677	584
Nan	25.351538	Nan	Nan
Nan	14.447125	Nan	Nan
Nan	1.000000	Nan	Nan
Nan	13.000000	Nan	Nan
Nan	25.000000	Nan	Nan
Nan	38.000000	Nan	Nan
Nan	50.000000	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 column to **snake case** for better readability and documentation.
- **Feature Engineering:**
  - o Created **age\_group** column by binning customer ages.
  - o 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 MySQL and loaded the cleaned DataFrame into the database for SQL analysis.

## 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
	16	81
	20	90

customer 11 ×

Output

Action Output

#	Time	Action	Message
1	00:40:42	select customer_id, purchase_amount from customer where discount_applied = "Ye...	839 row(s) returned

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

	item_purchased	Average Product Rating
▶	Gloves	3.86
	Sandals	3.84
	Boots	3.82
	Hat	3.8
	Skirt	3.78

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

	shipping_type	Average 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	avg_spend	total_revenue
▶	Yes	1053	59.49	62645
	No	2847	59.87	170436

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

	item_purchased	discount_rate
▶	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.

	item_rank	category	item_purchased	total_orders
▶	1	Accessories	Jewelry	171
	2	Accessories	Sunglasses	161
	3	Accessories	Belt	161
	1	Clothing	Blouse	171
	2	Clothing	Pants	171
	3	Clothing	Shirt	169
	1	Footwear	Sandals	160
	2	Footwear	Shoes	150
	3	Footwear	Sneakers	145
	1	Outerwear	Jacket	163
	2	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
▶	Young Adult	62143
	Middle-Aged	59197
	Adult	55978
	Senior	55763

## 5. Dashboard in Power BI

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

1. **Data Loading:** Import data from a MySQL Database
2. **Using DAX to calculate KPI's:**
  - Created measure **Numbers of Customers** by calculating the total of the **customer\_id** column
  - Created measure **Average Purchase Amount** by calculating the average of the **purchase\_amount** column
  - Created measure **Average Review Rating** by calculating the average of the **review\_rating** column
3. **Create data visualization:**
  - Using **cards** to display KPIs such as **number of customers**, **average purchase amount**, and **average review rating**.
  - Displaying **% of Customers by Subscription Status** using a **donut chart**.
  - Displaying **Revenue and Sales by Category** using a **clustered column chart**.
  - Displaying **Revenue and Sales by Age Group** using a **clustered bar chart**.
4. Interactive visual filter using **slicers** to select specific categories such as **subscription status**, **gender**, **item category**, and **shipping type**



## 6. Business Recommendations

- 1. Increase Subscription Adoption** – Offer subscription-only benefits (exclusive discounts, free shipping)
- 2. Focus on High-Performing Product Categories** – Prioritize marketing campaigns for Clothing and Accessories
- 3. Optimize Marketing by Age Group** – Allocate more marketing budget to Young Adult and Middle-Aged segments
- 4. Customer Loyalty Programs** – Reward repeat buyers to move them into the “Loyal” segment.
- 5. Review Discount Policy** – Balance sales boost with margin control.