

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.

	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	1
top	NaN	NaN	Male	Blouse	Clothing	NaN	Montana	M	Olive	Spring	NaN	No	Free Shipping	1
freq	NaN	NaN	2652	171	1737	NaN	96	1755	177	999	NaN	2847	675	22
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

Discount Applied	Promo Code Used	Previous Purchases	Payment Method	Frequency of Purchases
3900	3900	3900.000000	3900	3900
2	2	NaN	6	7
No	No	NaN	PayPal	Every 3 Months
2223	2223	NaN	677	584
NaN	NaN	25.351538	NaN	NaN
NaN	NaN	14.447125	NaN	NaN
NaN	NaN	1.000000	NaN	NaN
NaN	NaN	13.000000	NaN	NaN
NaN	NaN	25.000000	NaN	NaN
NaN	NaN	38.000000	NaN	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 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 PostgreSQL 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.

Result Grid			Filter Rows:
	gender	sum(purchase_amount)	
▶	Male	111061	
	Female	53161	

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

Result Grid			Filter Rows:
	customer_id	purchase_amount	
▶	2	64	
	3	73	
	4	90	
	9	97	
	12	68	
	13	72	
	16	81	
	20	90	
	24	88	
	32	79	
	33	67	
	37	69	
	40	60	
	43	100	

updated_customer_data 3 x

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

Result Grid			Filter Rows:
	item_purchased	high_avg_review	
▶	Boots	3.95	
	Gloves	3.88	
	Sandals	3.84	
	T-shirt	3.81	
	Hat	3.81	

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

Result Grid			Filter Rows:
	shipping_type	avg_purchase	
▶	Express	60.93	
	Standard	58.23	

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

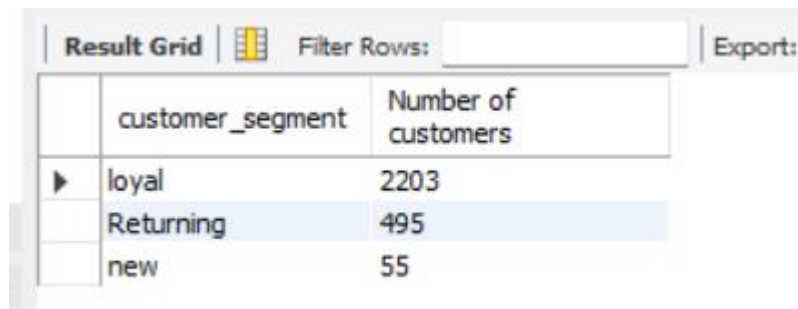
Result Grid					Filter Rows:	Export:	Wrap Cell Content:
	subscription_status	total_customers	avg_purchase	total_revenue			
▶	Yes	759	58.96	44754			
	No	1994	59.91	119468			

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



	item	discount_rate
▶	Sneakers	56.60
	Coat	51.72
	Hoodie	50.94
	Sweater	50.82
	Jewelry	47.15

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





	customer_segment	Number of customers
▶	loyal	2203
	Returning	495
	new	55

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

Result Grid		Filter Rows:	Export:	Wrap C
	item_rank	category	item_purchased	total_orders
▶	1	Accessories	Jewelry	123
	2	Accessories	Scarf	116
	3	Accessories	Belt	112
	1	Clothing	Shirt	127
	2	Clothing	Sweater	122
	3	Clothing	Dress	121
	1	Footwear	Sandals	116
	2	Footwear	Shoes	108
	3	Footwear	Sneakers	106
	1	Outerwear	Coat	116
	2	Outerwear	Jacket	109

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

Result Grid

Filter Rows:

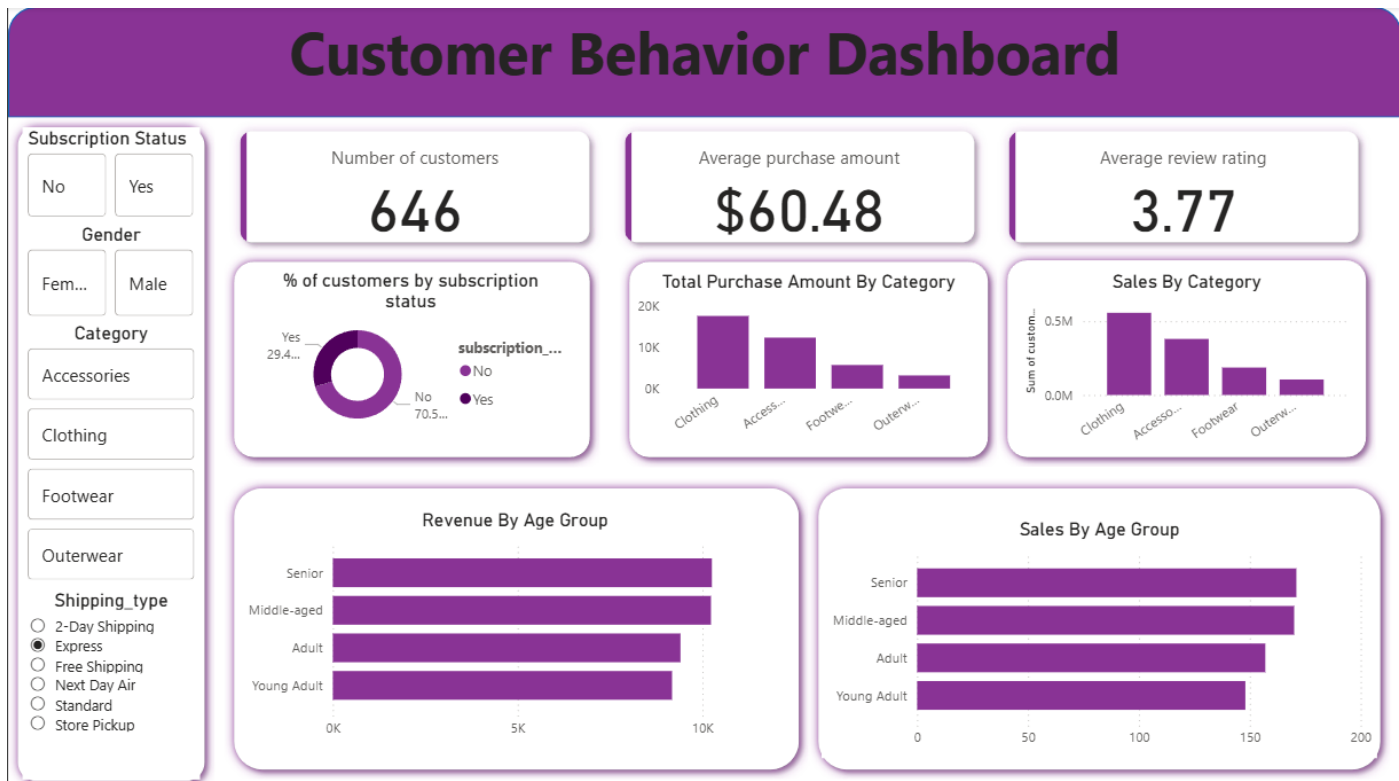
	subscription_status	repeat_buyers
▶	Yes	691
	No	1764

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

Result Grid	Filter Rows:	Ex
	total_revenue	customer_age_group
▶	94244	adult
	39964	young adult
	30014	old

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.