

# Customer Shopping Behaviour Analysis

## 1. Project Overview

This project analyses customer shopping behaviour 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 behaviour 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, Colour) Shopping behaviour (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 `df.describe()` for summary statistics.

	Customer ID	Age	Purchase Amount (USD)	Review Rating	Previous Purchases
count	3900.000000	3900.000000	3900.000000	3863.000000	3900.000000
mean	1950.500000	44.068462	59.764359	3.750065	25.351538
std	1125.977353	15.207589	23.685392	0.716983	14.447125
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	44.000000	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

```
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
```

- **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 PostgreSQL to answer key business questions:

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

Showing rows: 1 to 2  Page No

	gender text	total_revenue numeric
1	Female	75191
2	Male	157890


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

	customer_id bigint	purchase_amount bigint
1	2	64
2	3	73
3	4	90
4	7	85
5	9	97
6	12	68

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






	item_purchased text	Average Product Rating numeric
1	Gloves	3.86
2	Sandals	3.84
3	Boots	3.82
4	Hat	3.80
5	Skirt	3.78

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

Showing rows: 1 to 2  Page No:


	shipping_type text	round numeric
1	Standard	58.46
2	Express	60.48

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

Showing rows: 1 to 2  Page No:  of 1    


	subscription_status text	total_customers bigint	avg_spend numeric	total_revenue numeric
1	Yes	1053	59.49	62645.00
2	No	2847	59.87	170436.00

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

Showing rows: 1 to 5  Page No:




	item_purchased text	discount_rate numeric
1	Hat	50.00
2	Sneakers	49.66
3	Coat	49.07
4	Sweater	48.17
5	Pants	47.37

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

Showing rows: 1 to 3  Page No: 1 of 1


	customer_segment text	Number of Customers bigint
1	Loyal	3116
2	New	83
3	Returning	701

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

Showing rows: 1 to 11  Page No: 1 of 1  


	item_rank bigint	category text	item_purchased text	total_orders bigint
1	1	Accessori...	Jewelry	171
2	2	Accessori...	Sunglasses	161
3	3	Accessori...	Belt	161
4	1	Clothing	Blouse	171
5	2	Clothing	Pants	171
6	3	Clothing	Shirt	169

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

Showing rows: 1 to 2  Page No: 1

	subscription_status text	repeat_buyers bigint
1	No	2518
2	Yes	958

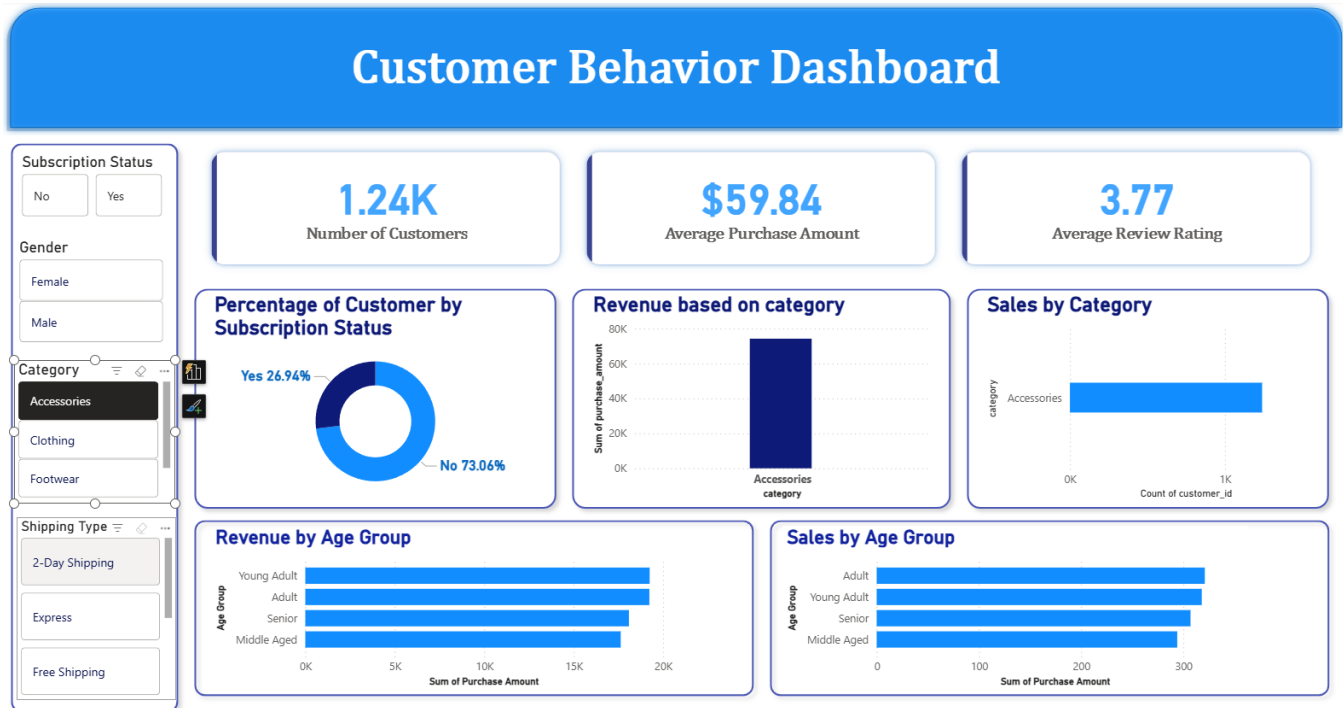
10. **Revenue by Age Group** – Calculated total revenue contribution of each age group.

Showing rows: 1 to 4  Page No: 1

	age_group text	total_revenue numeric
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 bestselling products in campaigns.
- **Targeted Marketing** – Focus efforts on high revenue age groups and express shipping users.