

# 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	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

Previous Purchases	Payment Method	Frequency of Purchases
3900.000000	3900	3900
NaN	6	7
	PayPal	Every 3 Months
NaN	677	584
25.351538	NaN	NaN
14.447125	NaN	NaN
1.000000	NaN	NaN
13.000000	NaN	NaN
25.000000	NaN	NaN
38.000000	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.

```
|: # Checking if missing data or null values are present in the dataset
df.isnull().sum()
```

```
|: Customer ID          0
Age                   0
Gender                0
Item Purchased        0
Category              0
Purchase Amount (USD) 0
Location              0
Size                  0
Color                 0
Season                0
Review Rating         37
Subscription Status   0
Shipping Type          0
Discount Applied       0
Promo Code Used        0
Previous Purchases    0
Payment Method         0
Frequency of Purchases 0
dtype: int64
```

```
# Imputing missing values in Review Rating column with the median rating of the product category

df['Review Rating'] = df.groupby('Category')['Review Rating'].transform(lambda x: x.fillna(x.median()))

df.isnull().sum()

Customer ID      0
Age              0
Gender           0
Item Purchased   0
Category          0
Purchase Amount (USD) 0
Location          0
Size              0
Color              0
Season             0
Review Rating     0
Subscription Status 0
Shipping Type     0
Discount Applied   0
Promo Code Used   0
Previous Purchases 0
Payment Method     0
Frequency of Purchases 0
dtype: int64
```

- **Column Standardization:** Renamed columns to **snake case** for better readability and documentation and replacing spaces with underscores (\_).

```
# Renaming columns according to snake casing for better readability and documentation

df.columns = df.columns.str.lower()
df.columns = df.columns.str.replace(' ', '_')
df = df.rename(columns={'purchase_amount_(usd)':'purchase_amount'})
```

```
df.columns

Index(['customer_id', 'age', 'gender', 'item_purchased', 'category',
       'purchase_amount', 'location', 'size', 'color', 'season',
       'review_rating', 'subscription_status', 'shipping_type',
       'discount_applied', 'promo_code_used', 'previous_purchases',
       'payment_method', 'frequency_of_purchases'],
      dtype='object')
```

- **Feature Engineering:**

- Created **age\_group** column by binning customer ages.

```

: # create a new column age_group
labels = ['Young Adult', 'Adult', 'Middle-aged', 'Senior']
df['age_group'] = pd.qcut(df['age'], q=4, labels = labels)

: df[['age','age_group']].head(10)

```

	age	age_group
0	55	Middle-aged
1	19	Young Adult
2	50	Middle-aged
3	21	Young Adult
4	45	Middle-aged
5	46	Middle-aged
6	63	Senior
7	27	Young Adult
8	26	Young Adult
9	57	Middle-aged

- Created **purchase\_frequency\_days** column from purchase data.

```

# create new column purchase_frequency_days

frequency_mapping = {
    'Fortnightly': 14,
    'Weekly': 7,
    'Monthly': 30,
    'Quarterly': 90,
    'Bi-Weekly': 14,
    'Annually': 365,
    'Every 3 Months': 90
}

df['purchase_frequency_days'] = df['frequency_of_purchases'].map(frequency_mapping)

df[['purchase_frequency_days','frequency_of_purchases']].head(10)

```

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

```

: (df['discount_applied'] == df['promo_code_used']).all()

```

: True

```
# Dropping promo code used column

df = df.drop('promo_code_used', axis=1)

df.columns

Index(['customer_id', 'age', 'gender', 'item_purchased', 'category',
       'purchase_amount', 'location', 'size', 'color', 'season',
       'review_rating', 'subscription_status', 'shipping_type',
       'discount_applied', 'previous_purchases', 'payment_method',
       'frequency_of_purchases', 'age_group', 'purchase_frequency_days'],
       dtype='object')
```

- **Database Integration:** Connected Python script to MySql and loaded the cleaned DataFrame into the database for SQL analysis.

```
[19]: from sqlalchemy import create_engine

# MySQL connection
username = "████████"
password = "████████"
host = "localhost"
port = "3306"

DATABASE_URL = 'mysql+mysqlconnector://{}:{}@{}:{}/customer_behavior'

engine = create_engine(DATABASE_URL)

# Write DataFrame to MySQL
table_name = "customer"
df.to_sql(table_name, engine, if_exists="replace", index=False)

# Read back sample
pd.read_sql("SELECT * FROM customer LIMIT 5;", engine)
```

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

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

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

	gender text 	revenue numeric 
1	Female	75191
2	Male	157890

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

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

Total rows: 839    Query complete 00:00:00

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

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

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

	<b>shipping_type</b> text		<b>round</b> numeric	
1	Standard		58.46	
2	Express		60.48	

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

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

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

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

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

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

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

	item_rank bigint	category text	item_purchased text	total_orders bigint
1	1	Accessories	Jewelry	171
2	2	Accessories	Sunglasses	161
3	3	Accessories	Belt	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

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

subscription_status	repeat_buyers
Yes	2049
No	1427

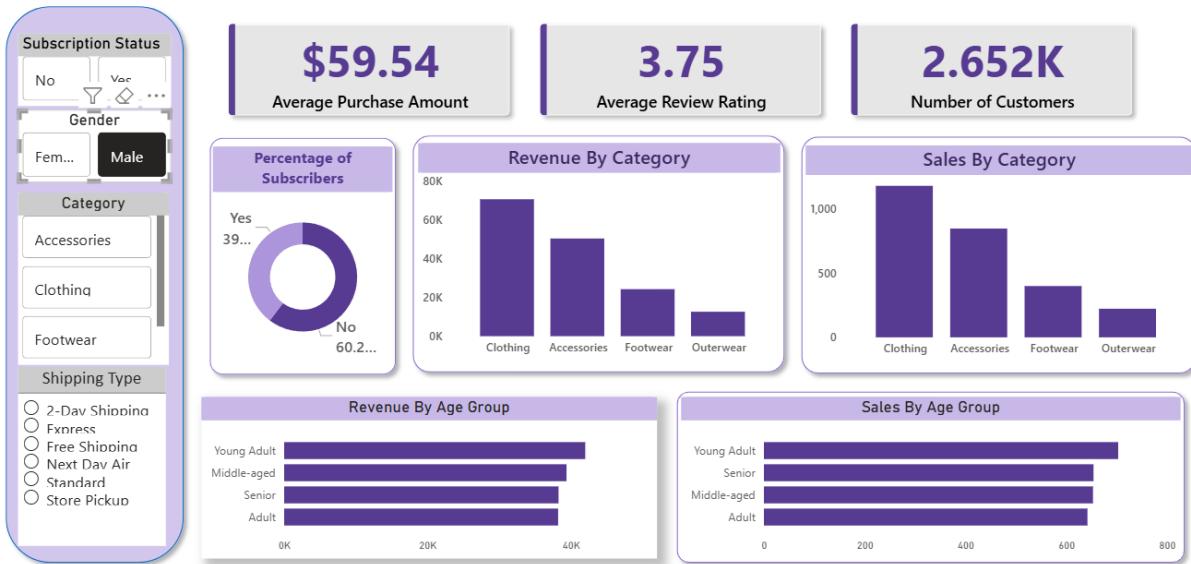
- **Repeat Buyers & Subscriptions** – Checked whether customers with >5 purchases are more likely to subscribe.

	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.

# Customer Behavior Dashboard



## 6. Business Recommendations

- Increase Subscriptions:** Offer special benefits to encourage more customers to subscribe.
- Build Customer Loyalty:** Give rewards or points to repeat buyers so they stay loyal.
- Check Discount Strategy:** Review discounts to boost sales while keeping profits stable.
- Promote Best Products:** Highlight high-rated and best-selling products in marketing campaigns.
- Target the Right Customers:** Focus marketing on high-revenue age groups and customers who prefer express delivery.