

Customer Shopping Behavior Analysis

An End-to-End Data Analysis Project Using Python, MySQL, and Power BI

Author: Layeeq Ahmed

1. Project Overview

Understanding customer shopping behavior is essential for improving sales strategies, optimizing marketing efforts, and enhancing customer experience.

This project includes:

- Data cleaning and preparation using Python
- Creating a structured relational table in MySQL
- Running analytical SQL queries for insights
- Developing an interactive dashboard in Power BI

2. Dataset Information

The dataset contains 18 columns and over 9000 customer shopping records.

3. Data Cleaning & Preparation (Python)

Key steps:

- Importing libraries
- Viewing dataset structure
- Standardizing column names
- Handling missing values
- Exporting cleaned data to MySQL

	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	39
unique	NaN	NaN	2	25	4	NaN	50	4	25	4	NaN	2	6	
top	NaN	NaN	Male	Blouse	Clothing	NaN	Montana	M	Olive	Spring	NaN	No	Free Shipping	
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										

Python SQL Connection:

```
engine =
create_engine("mysql+pymysql://root:Layeeq%401234@localhost:3306/customer_behaviour")
```

Table creation:

```
df.to_sql("customer", engine, if_exists="replace", index=False)
```

```
engine.connect()
print("Connected to customer_beaviour")
Connected to customer_beaviour

[23]: table_name = "customer"
df.to_sql(table_name, engine, if_exists="replace", index=False)

[23]: 3900

[24]: pd.read_sql("SELECT * FROM customer LIMIT 5;", engine)
[24]:   customer_id  age  gender item_purchased category purchase_amount  location  size  color  season  review_rating  subscription_status  shipping_type  disc
      0           1    55     Male       Blouse  Clothing            53  Kentucky    L    Gray  Winter      3.1        Yes    Express
      1           2    19     Male      Sweater  Clothing            64  Maine      L  Maroon  Winter      3.1        Yes    Express
      2           3    50     Male      Jeans  Clothing            73  Massachusetts  S  Maroon  Spring      3.1        Yes  Free Shipping
      3           4    21     Male     Sandals  Footwear            90  Rhode Island  M  Maroon  Spring      3.5        Yes  Next Day Air
      4           5    45     Male       Blouse  Clothing            49  Oregon     M Turquoise  Spring      2.7        Yes  Free Shipping

[25]: len(df)
[25]: 3900
```

4. Data Modeling & Analysis (MySQL)

Top 5 Products by Rating:

```
SELECT item_purchased, ROUND(AVG(review_rating),2) AS avg_rating  
FROM customer  
GROUP BY item_purchased  
ORDER BY avg_rating DESC  
LIMIT 5;
```

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

	gender text 	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
7	13	72
8	16	81
9	20	90
10	22	62
11	24	88

Total rows: 839 Query complete 00:00:00

- 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.

	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.

	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.

	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.

	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.

	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

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

	subscription_status	repeat_buyers
1	No	2518
2	Yes	958

- 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. Power BI Dashboard

Dashboard includes:

- Revenue by category
- Customer demographics
- Product ratings
- Purchase trends
- Payment method analysis



6. Key Insights

1. Winter season shows the highest order volume.
2. Clothing category contributes the most revenue.
3. Credit card users spend more on average.
4. High-rated products show better customer retention.
5. Women contribute slightly higher overall spending.

7. Business Recommendations

- Increase winter inventory and targeted promotions.
- Focus marketing on high-performing categories.
- Introduce rewards for frequent buyers.
- Highlight highly-rated products on digital platforms.
- Improve credit card payment experience.

8. Conclusion

This project demonstrates the complete workflow of a data analyst using Python, MySQL, and Power BI.