

Customer Shopping Trend Analysis

1. Project Overview

This project analyzes customer shopping trend 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)

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.

```
[5]: df.describe(include="all")
```

	Customer ID	Age	Gender	Item Purchased	Category	Purchase Amount (USD)	Location	Size	Color	Season	Review Rating	Subscription Status	Payment Method	Shipping Type	Discount Applied
count	3900.000000	3900.000000	3900	3900	3900	3900.000000	3900	3900	3900	3900	3900.000000	3900	3900	3900	3900
unique	NaN	NaN	2	25	4	NaN	50	4	25	4	NaN	2	6	6	2
top	NaN	NaN	Male	Blouse	Clothing	NaN	Montana	M	Olive	Spring	NaN	No	Credit Card	Free Shipping	No
freq	NaN	NaN	2652	171	1737	NaN	96	1755	177	999	NaN	2847	696	675	2223
mean	1950.500000	44.068462	NaN	NaN	NaN	59.764359	NaN	NaN	NaN	NaN	3.749949	NaN	NaN	NaN	NaN
std	1125.977353	15.207589	NaN	NaN	NaN	23.685392	NaN	NaN	NaN	NaN	0.716223	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	3.700000	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

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 no null found imputed few values in the `Review Rating` column where rating is < 4 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` was 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	age_group	frequency_of_purchase_days	purchase_amount
2	Young	14	64
3	Mid-Age	7	73
4	Young	7	90
7	Old	90	85
9	Young	365	97
12	Young	14	68
13	Old	14	72
16	Old	30	81
20	Old	14	90
22	Young	90	62

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

item_purchased	average_review_rating
Gloves	3.52
Hat	3.51
Sneakers	3.49
Sandals	3.48
Boots	3.48

4. **Shipping Type Comparison** – Compared average purchase amounts among shipping type.

shipping_type	average_purchase_amount
Standard	58.46
Next Day Air	58.63
Store Pickup	59.89
Free Shipping	60.41
Express	60.48
2-Day Shipping	60.73

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

no_of_customers	subscription_status	average_spend	total_revenue
1053	Yes	59.49	62645
2847	No	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	no_of_customers
Loyal	3188
Returning	629
New	83

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

product_rank	category	item_purchased
1	Accessories	Gloves
2	Accessories	Backpack
3	Accessories	Handbag
1	Clothing	Jeans
2	Clothing	T-shirt
3	Clothing	Hoodie
1	Footwear	Boots
2	Footwear	Sneakers
3	Footwear	Shoes
1	Outerwear	Coat
2	Outerwear	Jacket

9. **Repeat Buyers & Subscriptions** – Checked whether customers with > 2 purchases are more likely to subscribe.

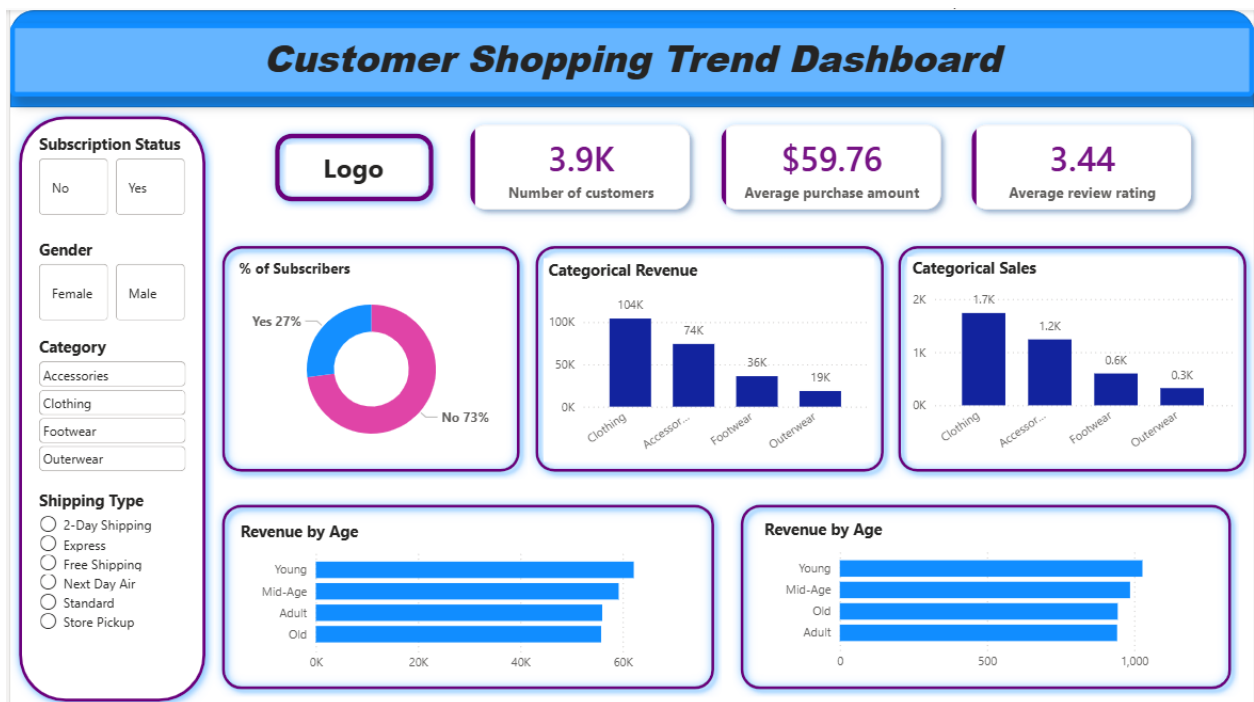
customer_count	subscription_status
1021	Yes
2724	No

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

revenue	age_group
59197	Mid-Age
62143	Young
55763	Old
55978	Adult

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.