

Customer Shopping Behavior Analysis

Created By:
Nikhil Mohan Tاتفale

February 7, 2026

1. Project Overview

This project analyzes customer shopping behavior using e-commerce transactional data from 3,900 customer purchases across multiple product categories. The primary goal is to discover meaningful patterns in customer spending habits, identify different customer segments, understand product preferences, and analyze subscription trends to support data-driven business strategies.

The analysis workflow includes data collection from Kaggle, exploratory data analysis (EDA) using Python, data cleaning and transformation, database integration with MySQL, business analytics through SQL queries, and visualization using Power BI dashboard.

2. Dataset Summary

Rows: 3,900

Columns: 18

Key Features:

- Customer Demographics: Age, Gender, Location, Subscription Status
- Purchase Information: Item Purchased, Category, Purchase Amount (USD), Season, Size, Color
- Shopping Patterns: Discount Applied, Promo Code Used, Previous Purchases, Frequency of Purchases, Review Rating, Shipping Type, Payment Method

Missing Data: 37 missing values identified in the Review Rating column

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	Previous Purchases	Payment Method	Frequency of Purchases	
0	1	55	Male	Blouse	Clothing	53	Kentucky	L	Gray	Winter	3.1	Yes	Express	Yes	Yes	14	Venmo	Fortnightly
1	2	19	Male	Sweater	Clothing	64	Maine	L	Maroon	Winter	3.1	Yes	Express	Yes	Yes	2	Cash	Fortnightly
2	3	50	Male	Jeans	Clothing	73	Massachusetts	S	Maroon	Spring	3.1	Yes	Free Shipping	Yes	Yes	23	Credit Card	Weekly
3	4	21	Male	Sandals	Footwear	90	Rhode Island	M	Maroon	Spring	3.5	Yes	Next Day Air	Yes	Yes	49	PayPal	Weekly
4	5	45	Male	Blouse	Clothing	49	Oregon	M	Turquoise	Spring	2.7	Yes	Free Shipping	Yes	Yes	31	PayPal	Annually

Figure 1: Sample view of the dataset

3. Exploratory Data Analysis using Python

The data preparation and cleaning process was conducted in Google Colab using Python and the Pandas library. Below are the key steps performed:

3.1 Data Loading and Initial Exploration

We started by importing the dataset using pandas and performing initial checks to understand the data structure.

```

1 import pandas as pd
2
3 # Load the dataset
4 df = pd.read_csv('shopping_behavior.csv')
5
6 # Check data structure
7 df.info()
8
9 # Generate summary statistics
10 df.describe()

```

Listing 1: Loading and exploring the dataset

	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	Previous Purchases	Payment Method	Frequency of Purchases
count	3900	0.000000	3900	0.000000	3900	0.000000	3900	0.000000	3900	0.000000	3900	0.000000	3900	0.000000	3900	0.000000	3900	0.000000
unique	NaN	NaN	2	25	4	NaN	50	4	25	4	NaN	2	6	2	2	NaN	6	7
top	NaN	NaN	Male	Blouse	Clothing	NaN	Montana	M	Olive	Spring	NaN	No	Free Shipping	No	No	NaN	PayPal	Every 3 Months
freq	NaN	NaN	2652	171	1737	NaN	96	1755	177	999	NaN	2847	675	2223	2223	NaN	677	584
mean	1950	0.000000	44.068462	NaN	NaN	59.754359	NaN	NaN	NaN	NaN	3.750065	NaN	NaN	NaN	NaN	25.351538	NaN	NaN
std	1125.977353	15.207589	NaN	NaN	NaN	23.685392	NaN	NaN	NaN	NaN	0.716063	NaN	NaN	NaN	NaN	14.647125	NaN	NaN
min	1.000000	18.000000	NaN	NaN	NaN	20.000000	NaN	NaN	NaN	NaN	2.500000	NaN	NaN	NaN	NaN	1.000000	NaN	NaN
25%	975.750000	31.000000	NaN	NaN	NaN	39.000000	NaN	NaN	NaN	NaN	3.100000	NaN	NaN	NaN	NaN	13.000000	NaN	NaN
50%	1950.000000	44.000000	NaN	NaN	NaN	60.000000	NaN	NaN	NaN	NaN	3.800000	NaN	NaN	NaN	NaN	25.000000	NaN	NaN
75%	2925.250000	57.000000	NaN	NaN	NaN	81.000000	NaN	NaN	NaN	NaN	4.400000	NaN	NaN	NaN	NaN	36.000000	NaN	NaN
max	3900.000000	70.000000	NaN	NaN	NaN	100.000000	NaN	NaN	NaN	NaN	5.000000	NaN	NaN	NaN	NaN	50.000000	NaN	NaN

Figure 2: Statistical summary showing distribution of numerical features

Key observations from the initial exploration:

Customer age ranges from 18 to 70 years with a mean of approximately 44 years

Purchase amounts vary between \$20 and \$100, averaging around \$59.76

Review ratings span from 2.5 to 5.0 with an average rating of 3.75

Customers have between 1 and 50 previous purchases

3.2 Missing Data Handling

We identified and addressed missing values in the dataset.

```
df.isnull().sum()
```

	0
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

Figure 3: Missing value detection - 37 null values in Review Rating

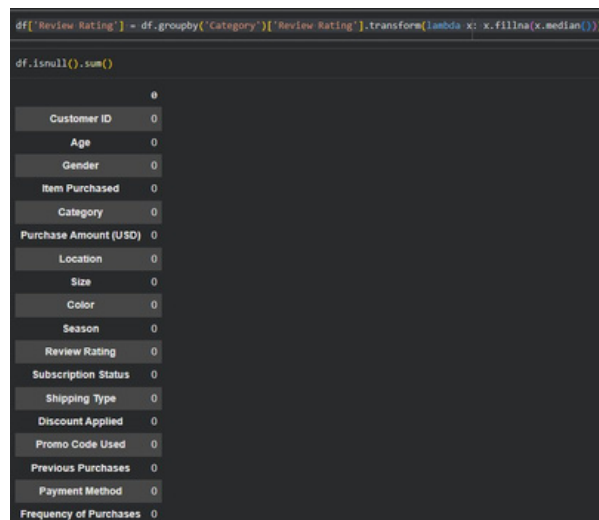
Solution: Missing values in the Review Rating column were filled using the median rating for each product category. This approach preserves category-specific rating patterns rather than using a global median.

```

1 df['Review Rating'] = df.groupby('Category')['Review Rating'].
    transform(
2     lambda x: x.fillna(x.median())
3 )
4
5 # Verify no missing values remain
6 df.isnull().sum()

```

Listing 2: Imputing missing values with category-wise median



```

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

df.isnull().sum()

```

	0
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

Figure 4: Verification - all missing values successfully handled

3.3 Column Standardization

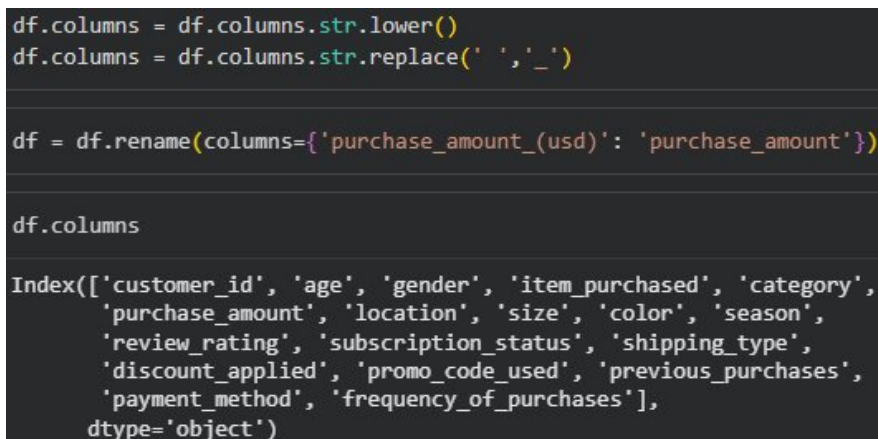
All column names were converted to snake.case format for consistency and easier database integration.

```

1 df.columns = df.columns.str.lower()
2 df.columns = df.columns.str.replace(' ', '_')

```

Listing 3: Standardizing column names



```

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')

```

Figure 5: Standardized column names in snake case

3.4 Feature Engineering

Two new features were created to enhance analytical capabilities:

1. Age Group Column: Customers were segmented into four age groups using quartiles.

```
1 labels = ['Young Adult', 'Adult', 'Middle-aged', 'Senior']
2 df['age_group'] = pd.qcut(df['age'], q=4, labels=labels)
```

Listing 4: Creating age group categories



Figure 6: Age group segmentation

2. Purchase Frequency Days: Converted categorical frequency values to numerical days for quantitative analysis.

```
1 frequency_mapping = {
2     'Fortnightly': 14,
3     'Weekly': 7,
4     'Monthly': 30,
5     'Quarterly': 90,
6     'Bi-Weekly': 14,
7     'Annually': 365,
8     'Every 3 Months': 90
9 }
10 df['purchase_frequency_days'] = df['frequency_of_purchases'].map(
    frequency_mapping)
```

Listing 5: Mapping purchase frequency to days

```
# create the column of 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)
```

	purchase_frequency_days	frequency_of_purchases
0	14	Fortnightly
1	14	Fortnightly
2	7	Weekly
3	7	Weekly
4	365	Annually
5	7	Weekly
6	90	Quarterly
7	7	Weekly
8	365	Annually
9	90	Quarterly

Figure 7: Purchase frequency converted to numerical days

3.5 Data Consistency Check

We identified redundant columns by checking if discount applied and promo code used contained identical information.

```
1 (df['discount_applied'] == df['promo_code_used']).all()
2 # Returns: True
3
4 # Drop redundant column
5 df = df.drop('promo_code_used', axis=1)
```

Listing 6: Checking for redundant columns

```
# check wheather discount_applied and promo_code_use is same data or not
(df['discount_applied'] == df['promo_code_used']).all()

np.True_

# so let drop the promo_code_used
df = df.drop('promo_code_used', axis=1)
```

Figure 8: Removing redundant promo code used column

3.6 Database Integration

After cleaning and preparing the data, the final dataset was exported and loaded into a MySQL database for structured querying and business analysis.

4. Data Analysis using SQL (Business Questions)

The cleaned dataset was uploaded to MySQL to perform structured analysis and answer critical business questions. Below are the key queries and findings:

1. Revenue by Gender

Business Question: What is the total revenue generated by male versus female customers?

```
1 SELECT gender, SUM(purchase_amount) AS revenue
2 FROM updated_customer_shopping_behavior
3 GROUP BY gender;
```

Listing 7: Analyzing revenue contribution by gender

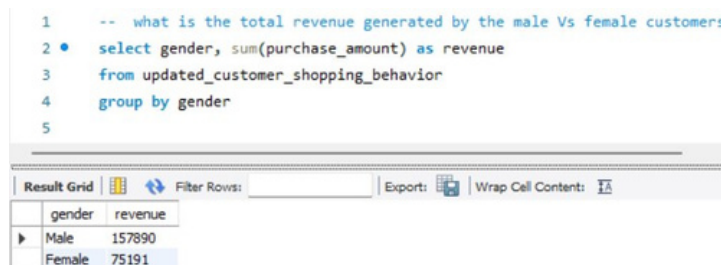


Figure 9: Total revenue comparison between genders

Key Finding: Male customers contribute \$157,690 in revenue compared to \$75,191 from female customers. This reveals a significant gender gap presenting an opportunity for targeted female customer acquisition.

2. High-Spending Discount Users

Business Question: Which customers used discounts but still spent above the average purchase amount?

```
1 SELECT customer_id, discount_applied, purchase_amount
2 FROM updated_customer_shopping_behavior
3 WHERE purchase_amount > (SELECT AVG(purchase_amount)
4                           FROM updated_customer_shopping_behavior
5                           )
6 AND discount_applied = 'yes'
7 ORDER BY purchase_amount ASC;
```

Listing 8: Identifying discount users with high spending


```

7 -- which customer use discount but still spent more than the average purchase amount
8 * select customer_id, discount_applied, purchase_amount
9 from updated_customer_shopping_behavior
10 where purchase_amount >= (select AVG(purchase_amount) from updated_customer_shopping_behavior) and discount_applied = 'yes'
11 order by purchase_amount asc

```

customer_id	discount_applied	purchase_amount
40	Yes	60
166	Yes	60
304	Yes	60
534	Yes	60
558	Yes	60
589	Yes	60
635	Yes	60
677	Yes	60
712	Yes	60
777	Yes	60
853	Yes	60
858	Yes	60
895	Yes	60
923	Yes	60
1002	Yes	60
1247	Yes	60
1296	Yes	60
1333	Yes	60
1424	Yes	60
1434	Yes	60
187	Yes	61

Figure 10: Customers using discounts while maintaining above-average spending

Key Finding: Multiple customers spend \$60 or more even with discounts applied, indicating that strategic discounting doesn't necessarily erode revenue.

3. Top 5 Products by Rating

Business Question: Which products have the highest average customer review ratings?

```

1 SELECT item_purchased, ROUND(AVG(review_rating),2) AS '
   review_rating_average'
2 FROM updated_customer_shopping_behavior
3 GROUP BY item_purchased
4 ORDER BY AVG(review_rating) DESC
5 LIMIT 5;

```

Listing 9: Finding top-rated products

```

14 -- 3. which are top 5 product with highest average review rating
15 * select item_purchased, round(Avg(review_rating),2) as 'review rating average'
16 from updated_customer_shopping_behavior
17 group by item_purchased
18 order by avg(review_rating) desc
19 limit 5;
20

```

item_purchased	review rating average
Gloves	3.86
Sandals	3.84
Boots	3.82
Hat	3.8
Skirt	3.78

Figure 11: Top 5 products by average review rating

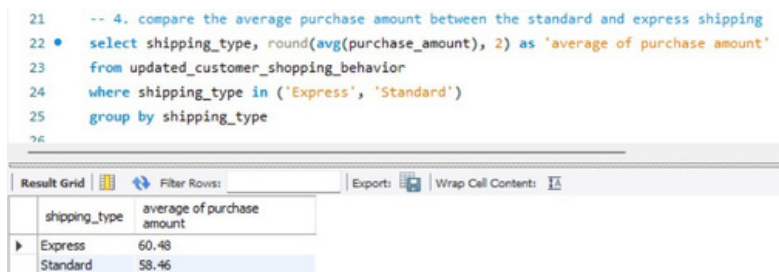
Key Finding: Gloves (3.86), Sandals (3.84), and Boots (3.82) receive the highest customer satisfaction ratings and should be prominently featured in marketing efforts.

4. Shipping Type Comparison

Business Question: Is there a difference in average purchase amount between standard and express shipping customers?

```
1 SELECT shipping_type, ROUND(AVG(purchase_amount), 2) AS '
   average_purchase_amount'
2 FROM updated_customer_shopping_behavior
3 WHERE shipping_type IN ('Express', 'Standard')
4 GROUP BY shipping_type;
```

Listing 10: Comparing purchase amounts by shipping type



```
21 -- 4. compare the average purchase amount between the standard and express shipping
22 • select shipping_type, round(avg(purchase_amount), 2) as 'average of purchase amount'
23 from updated_customer_shopping_behavior
24 where shipping_type in ('Express', 'Standard')
25 group by shipping_type
```

shipping_type	average of purchase amount
Express	60.48
Standard	58.46

Figure 12: Average purchase amount by shipping method

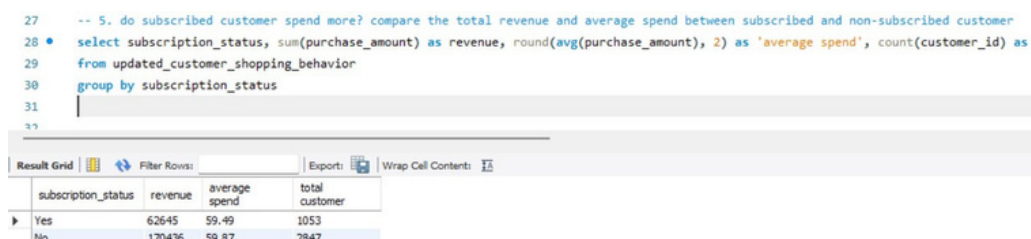
Key Finding: Express shipping customers spend slightly more on average (\$60.48 vs \$58.46), suggesting a correlation between urgency and spending.

5. Subscribers vs. Non-Subscribers

Business Question: Do subscribed customers spend more? Compare total revenue and average spend.

```
1 SELECT subscription_status,
2       SUM(purchase_amount) AS revenue,
3       ROUND(AVG(purchase_amount), 2) AS 'average_spend',
4       COUNT(customer_id) AS 'total_customer'
5 FROM updated_customer_shopping_behavior
6 GROUP BY subscription_status;
```

Listing 11: Analyzing subscription impact on spending



```
27 -- 5. do subscribed customer spend more? compare the total revenue and average spend between subscribed and non-subscribed customer
28 • select subscription_status, sum(purchase_amount) as revenue, round(avg(purchase_amount), 2) as 'average spend', count(customer_id) as
29 from updated_customer_shopping_behavior
30 group by subscription_status
```

subscription_status	revenue	average spend	total customer
Yes	62645	59.49	1053
No	170436	59.87	2847

Figure 13: Comparison of subscribed vs non-subscribed customers

Key Finding: Only 27% of customers are subscribed. While average spending per transaction is similar, the low subscription rate represents a significant growth opportunity.

6. Discount-Dependent Products

Business Question: Which products have the highest percentage of purchases with discounts applied?

```

1 SELECT item_purchased,
2        ROUND(100 * SUM(CASE WHEN discount_applied = 'yes' THEN 1
3          ELSE 0 END)
4          / COUNT(*), 2) AS discount_rate
5 FROM updated_customer_shopping_behavior
6 GROUP BY item_purchased
7 ORDER BY discount_rate DESC
8 LIMIT 5;

```

Listing 12: Calculating discount rate by product

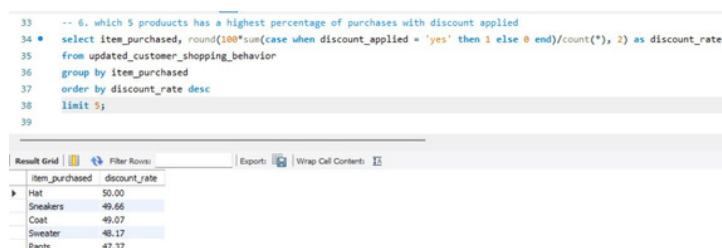


Figure 14: Products with highest discount dependency

Key Finding: Hat (50%), Sneakers (49.66%), and Coat (49.07%) show very high discount rates, suggesting potential over-reliance on promotional pricing.

7. Customer Segmentation

Business Question: How are customers distributed across New, Returning, and Loyal segments?

```

1 WITH customer_type AS
2   SELECT customer_id, previous_purchases, purchase_amount,
3          CASE
4            WHEN previous_purchases = 1 THEN 'New'
5            WHEN previous_purchases BETWEEN 2 AND 10 THEN 'Returning'
6            ELSE 'Loyal'
7          END AS customer_segmented
8   FROM updated_customer_shopping_behavior
9 )
10 SELECT customer_segmented,
11        COUNT(*) AS 'number_of_customers',
12        ROUND(SUM(purchase_amount), 2) AS 'revenue'
13 FROM customer_type
14 GROUP BY customer_segmented;

```

Listing 13: Segmenting customers by purchase history

```

41 -- 7. segment customer into new, returning and loyal based on their total number of previous purchases and show the count of each segment
42 with customer_type as (
43   select customer_id, previous_purchases, purchase_amount,
44   case when previous_purchases = 1 then 'new' when previous_purchases between 2 and 10 then 'returning' else 'loyal' end as customer_segmented
45   from updated_customer_shopping_behavior
46 )
47 select customer_segmented, count(*) as 'number of customers', round(sum(purchase_amount), 2) as 'revenue'
48 from customer_type
49 group by customer_segmented
50

```

customer_segmented	number of customers	revenue
loyal	3116	185517
returning	701	42711
new	83	4853

Figure 15: Customer distribution across segments

Key Finding: Loyal customers (80%) drive the majority of revenue, but new customer acquisition is weak at only 2% of the customer base.

8. Top 3 Products per Category

Business Question: What are the most purchased products within each category?

```

1 WITH item_count AS (
2   SELECT category, item_purchased, COUNT(customer_id) AS
3     total_orders,
4     ROW_NUMBER() OVER (PARTITION BY category
5                          ORDER BY COUNT(customer_id) DESC) AS
6     item_rank
7   FROM updated_customer_shopping_behavior
8   GROUP BY category, item_purchased
9 )
10 SELECT item_rank, category, item_purchased, total_orders
11 FROM item_count
12 WHERE item_rank <= 3;

```

Listing 14: Identifying best-sellers by category

```

52 -- 8. what are the top 3 most purchase products within each the category
53 • with item_count as (
54   select category, item_purchased, count(customer_id) as total_orders,
55   row_number() over(partition by category order by count(customer_id) desc) as item_rank
56   from updated_customer_shopping_behavior
57   group by category, item_purchased
58 )
59 select item_rank, category, item_purchased, total_orders
60 from item_count
61 where item_rank <=3
62

```

item_rank	category	item_purchased	total_orders
1	Accessories	Jewelry	171
2	Accessories	Sunglasses	161
3	Accessories	Belt	161
1	Clothing	Blouse	171
2	Clothing	Pants	171
3	Clothing	Shirt	169
1	Footwear	Sandals	160
2	Footwear	Shoes	150
3	Footwear	Sneakers	145
1	Outerwear	Jacket	163
2	Outerwear	Coat	161

Figure 16: Top 3 products in each category

Key Finding: Jewelry leads Accessories, Blouse leads Clothing, Sandals lead Footwear, and Jacket leads Outerwear. These items should be prioritized in inventory management.

9. Repeat Buyers & Subscriptions

Business Question: Are customers with more than 5 previous purchases more likely to subscribe?

```

1 SELECT subscription_status, COUNT(customer_id) AS repeat_buyers
2 FROM updated_customer_shopping_behavior
3 WHERE previous_purchases > 5
4 GROUP BY subscription_status;

```

Listing 15: Analyzing subscription behavior among repeat buyers

```

64 -- 9. are customer who are repeat buyer( more than 5 previous orders) are also likely to subscribed
65 • select subscription_status, count(customer_id) as repeat_buyers
66 from updated_customer_shopping_behavior
67 where previous_purchases > 5
68 group by subscription_status
69
70
71 -- 10. what is the revenue contribution of each age group

```

subscription_status	repeat_buyers
Yes	958
No	2518

Figure 17: Subscription status among repeat buyers

Key Finding: Even among repeat buyers, 72.5% are not subscribed, indicating subscription benefits may not be compelling enough.

10. Revenue by Age Group

Business Question: What is the revenue contribution of each age group?

```
1 SELECT age_group, ROUND(SUM(purchase_amount), 2) AS revenue
2 FROM updated_customer_shopping_behavior
3 GROUP BY age_group
4 ORDER BY revenue DESC;
```

Listing 16: Calculating revenue by age group

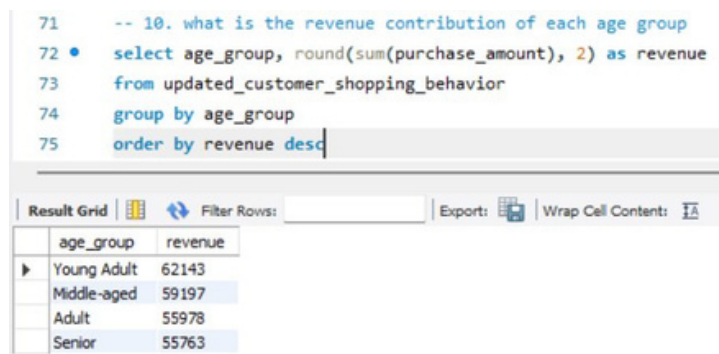


Figure 18: Revenue distribution across age groups

Key Finding: Revenue is relatively balanced across all age groups, with Young Adults contributing slightly more (\$62,143), suggesting broad market appeal.

5. Dashboard in Power BI

An interactive Power BI dashboard was created to visualize key metrics and enable dynamic exploration of the data.



Figure 19: Power BI Dashboard - Customer Shopping Behavior Overview

5.1 Dashboard Highlights

Key Performance Indicators:

Total Customers: 3,900

Average Purchase Amount: \$59.8

Average Review Rating: 3.75/5.0

Subscription Rate: 27%

Visual Components:

Revenue by Category - Clothing dominates with approximately \$100K

Revenue by Age Group - Senior customers generate the highest revenue (\$91.4K)

Sales by Category - Clothing accounts for 44.54% of all sales

Sales by Age Group - Young Adults show the highest purchase frequency

Subscription Status Distribution - Clear visualization of the subscription gap

Interactive Filters: The dashboard includes slicers for Subscription Status, Gender, Category, and Shipping Type, enabling users to drill down into specific customer segments and analyze behavior patterns dynamically.

6. Business Recommendations

Based on the comprehensive data analysis, the following strategic recommendations are proposed:

1. Expand Female Customer Base

Male customers generate 68% of total revenue. Implement targeted marketing campaigns for female demographics, expand product offerings that appeal to women, and create female-focused promotional events to bridge this revenue gap.

2. Enhance Subscription Value Proposition

Only 27% of customers subscribe. Implement exclusive subscriber-only discounts, free express shipping for members, loyalty points system, and clear communication of subscription benefits to increase penetration.

3. Optimize Discount Strategy

Products like Hats and Sneakers show 50% discount dependency. Test higher base prices with selective discounting, implement time-limited flash sales, create product bundles to maintain value, and reserve heavy discounts for customer acquisition.

4. Strengthen Customer Acquisition

Only 2% of customers are new. Accelerate growth through referral programs, first- purchase discount codes, enhanced digital marketing presence (SEO, social media, paid ads), and strategic brand partnerships.

5. Leverage High-Rated Products

Products with ratings above 3.80 (Gloves, Sandals, Boots) should be featured prominently on the homepage, highlighted in email campaigns, showcased in customer testimonials, and prioritized in inventory decisions.

6. Implement Age-Targeted Marketing

Different age groups show distinct behaviors. Target Young Adults with social media and influencer partnerships, Middle-aged customers with email campaigns and family bundles, and Seniors with easy-to-use interfaces and traditional marketing channels.