

Customer Shopping Behavior Analysis

1. Project Overview

This project explores customer shopping behavior using transactional data from 3,900 purchases across multiple product categories. The objective is to uncover insights related to spending patterns, customer segments, product preferences, and subscription behavior, enabling data-driven strategic decisions for improved customer engagement and revenue growth.

2. Dataset Summary

- **Total Rows:** 3,900
 - **Columns:** 18
 - **Key Features Include:**
 - **Customer Demographics:** Age, Gender, Location, Subscription Status
 - **Purchase Details:** Item Purchased, Category, Purchase Amount, Season, Size, Color
 - **Behavioral Indicators:** Discount Applied, Previous Purchases, Purchase Frequency, Review Ratings, Shipping Type
 - **Missing Data:**
 - 37 missing entries in the *Review Rating* column, handled via imputation
-

3. Exploratory Data Analysis (Python)

The EDA phase focuses on cleaning, transforming, and preparing the data for deeper analysis:

Data Preparation

- Loaded the dataset using **pandas**.
- Performed structural checks with `.info()` and statistical summaries using `.describe()`.

Data Cleaning

- Identified and imputed missing values in *Review Rating* using the **median rating per product category**.
- Standardized column names to **snake_case** for better readability and consistency.

Feature Engineering

- Created an **age_group** feature by binning age ranges.
- Generated a **purchase_frequency_days** metric from timestamps.
- Assessed redundancy between *discount_applied* and *promo_code_used* — removed *promo_code_used* to avoid duplication.

Database Integration

- Connected the cleaned dataset to **PostgreSQL**.
 - Loaded the transformed DataFrame into the database for SQL-based business analysis.
-

4. SQL-Based Business Analysis

Using PostgreSQL, we answered several key business questions:

1. **Revenue by Gender** – Compared total spending between male and female customers.
 2. **High-Spending Discount Users** – Identified customers who applied discounts but still spent above the overall average.
 3. **Top 5 Products by Review Rating** – Ranked products based on their average customer ratings.
 4. **Shipping Type Analysis** – Compared average purchase amounts for Standard vs. Express shipping.
 5. **Subscriber vs. Non-Subscriber Behavior** – Analyzed spend levels and revenue contribution by subscription status.
 6. **Discount-Dependent Products** – Identified products with the highest proportion of discounted purchases.
 7. **Customer Segmentation** – Categorized customers as **New**, **Returning**, or **Loyal** based on purchase frequency.
 8. **Top Products per Category** – Extracted the top 3 most purchased products within each category.
 9. **Repeat Buyers & Subscription Likelihood** – Examined whether customers with more than 5 purchases are more likely to be subscribers.
 10. **Revenue by Age Group** – Measured revenue contribution across age segments.
-

5. Power BI Dashboard

An interactive dashboard was built using **Power BI**, featuring:

- Revenue insights
- Customer segmentation visuals
- Product performance metrics
- Discount utilization patterns
- Age, gender, and geographic breakdowns

The dashboard provides a clear visual narrative of customer behavior and purchasing trends.

6. Business Recommendations

Based on our analysis, we propose the following strategies:

- **Boost Subscriptions:** Highlight exclusive benefits and run targeted subscription campaigns.
- **Strengthen Loyalty Programs:** Reward returning customers to convert them into loyal, high-value customers.
- **Optimize Discount Strategy:** Offer discounts strategically to boost conversions without eroding margins.
- **Enhance Product Positioning:** Promote top-rated and best-selling products to increase conversions.
- **Targeted Marketing Efforts:** Prioritize high-revenue age groups and express-shipping customers for targeted outreach.

Problem Statement

Problem Statement

A leading retail company wants to better understand its customers' shopping behavior in order to improve sales, customer satisfaction, and long-term loyalty. The management team has noticed changes in purchasing patterns across demographics, product categories, and sales channels (online vs. offline). They are particularly interested in uncovering which factors, such as discounts, reviews, seasons, or payment preferences, drive consumer decisions and repeat purchases.

You are tasked with analyzing the company's consumer behavior dataset to answer the following overarching business question:

"How can the company leverage consumer shopping data to identify trends, improve customer engagement, and optimize marketing and product strategies?"

A leading retail company wants to better understand its customers' shopping behavior in order to

- improve sales, customer satisfaction, and long-term loyalty. The management team has noticed
- changes in purchasing patterns across demographics, product categories, and sales channels
- (online vs. offline). They are particularly interested in uncovering factors, such as discounts,
- reviews, seasons, or payment preferences, drive consumer decisions and repeat purchases.

You are tasked with analyzing the company's consumer behavior dataset to answer the following overarching business question:

"How can the company leverage consumer shopping data to identify trends, improve customer engagement, and optimize marketing and product strategies?"

Q1. What is the total revenue generated by male vs. female customers?

SQL Query Used:

```
sql
SELECT
    gender,
    SUM(purchase_amount) AS total_revenue
FROM customer
GROUP BY gender;
```

Result Summary:

- **Female Customers:** ₹75,191
- **Male Customers:** ₹157,890

	gender	revenue
	text	numeric
1	Female	75191
2	Male	157890

Insight:

Male customers contribute **significantly higher total revenue** compared to female customers—more than 2x in this dataset. This suggests potential opportunities to improve engagement and conversion among female customers.

Q2. Which customers used a discount but still spent more than the average purchase amount?

To identify high-value customers who applied a discount yet spent above the overall average purchase amount, we use the following query:

```
sql

SELECT
    customer_id,
    purchase_amount
FROM customer
WHERE
    discount_applied = 'Yes'
    AND purchase_amount >= (
        SELECT AVG(purchase_amount)
        FROM customer
    );
```

Explanation:

- Filters only those customers who **used a discount**.
- Compares each customer's spending against the **average purchase amount** across all customers.
- Returns customers who, despite discounts, still demonstrate **above-average spending behavior**, indicating strong purchasing intent or loyalty.

	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

Q3. Which are the top 5 products with the highest average review rating?

Answer:

The following SQL query identifies the **top 5 highest-rated products** by calculating the **average review rating** for each item purchased. It first groups the data by item purchased, then computes the average rating for each group. Finally, it sorts the results in descending order to show the best-rated products at the top and limits the output to the top five:

```
sql

SELECT
    item_purchased,
    ROUND(AVG(review_rating::numeric), 2) AS "Average Product Rating"
FROM customer
GROUP BY item_purchased
ORDER BY AVG(review_rating) DESC
LIMIT 5;
```

Explanation:

- `AVG(review_rating)` calculates the average rating for each product.
- `GROUP BY item_purchased` ensures the average is computed per product.
- `ORDER BY ... DESC` sorts products from highest to lowest rating.
- `LIMIT 5` returns only the top 5 best-rated items.

	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

Q4. Compare the average Purchase Amount between Standard and Express Shipping.

The goal of this query is to compare how much customers spend on average when selecting **Standard** versus **Express** shipping options.

```
sql

SELECT shipping_type,
       ROUND(AVG(purchase_amount),2)
  FROM customer
 WHERE shipping_type IN ('Standard','Express')
 GROUP BY shipping_type;
```

Explanation:

- The query filters the dataset to only include orders where the shipping method is **Standard** or **Express**.
- It then calculates the **average purchase amount** for each shipping type using `AVG(purchase_amount)`.
- Finally, `GROUP BY shipping_type` ensures that the results are shown separately for each shipping method.

Purpose of the Analysis:

This comparison helps identify whether customers who choose **Express** shipping tend to spend more than those selecting **Standard** shipping, which is useful for understanding customer behavior, pricing impact, and shipping preference trends.

	shipping_type text	avg numeric
1	Standard	58.4602446483180428
2	Express	60.4752321981424149

Q5. Do subscribed customers spend more? Compare average spend and total revenue between subscribers and non-subscribers.

Answer:

```
SELECT
    subscription_status,
    COUNT(customer_id) AS total_customers,
    ROUND(AVG(purchase_amount), 2) AS avg_spend,
    ROUND(SUM(purchase_amount), 2) AS total_revenue
FROM customer
GROUP BY subscription_status
ORDER BY total_revenue, avg_spend DESC;
```

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

Based on the aggregated results from the query:

- **Subscribed customers have a higher average spend per customer**, indicating stronger engagement and more frequent or higher-value purchases.
- **Subscribers also contribute a significantly larger share of total revenue**, even if their overall customer count is smaller.
- **Non-subscribers show lower average spend and lower total revenue**, meaning they add less financial value despite possibly forming a larger portion of the customer base.

Conclusion:

Yes, subscribed customers spend more.

They generate **higher revenue per person** and **higher total revenue overall**, making them the most valuable customer segment for the business.

Q6. Which 5 products have the highest percentage of purchases with discounts applied?

```
SELECT
    item_purchased,
    ROUND(100* SUM(CASE WHEN discount_applied = 'Yes' THEN 1 ELSE 0 END) / COUNT(*), 2) AS discount_rate
FROM customer
GROUP BY item_purchased
ORDER BY discount_rate DESC
LIMIT 5;
```

	item_purchased	discount_rate
1	Hat	50.00
2	Sneakers	49.00
3	Coat	49.00
4	Sweater	48.00
5	Pants	47.00

Answer:

This query calculates the **discount rate** for each product by dividing the number of purchases with discounts applied by the total purchases of that product. It then returns the **top 5 products** with the highest percentage of discounted purchases.

Insights:

- These top 5 products exhibit the **highest reliance on discounts**, suggesting customers are more likely to purchase them when promotions are available.
- A high discount rate often indicates:
 - The product is **highly price-sensitive**,
 - **frequent discount campaigns** are run for that product,
 - customers prefer waiting for offers rather than buying at full price.
- This insight can be used by marketing and pricing teams
 - adjust discount frequency,
 - improve product pricing strategy,
 - or identify items where discounts are driving significant sales volume.

Q7. Segment customers into New, Returning, and Loyal based on their total number of previous purchases and show the count of each segment.

```
WITH customer_type AS (
    SELECT
        customer_id,
        previous_purchases,
        CASE
            WHEN previous_purchases = 1 THEN 'New'
            WHEN previous_purchases BETWEEN 2 AND 10 THEN 'Returning'
            ELSE 'Loyal'
        END AS customer_segment
    FROM customer
)
SELECT
    customer_segment,
    COUNT(*) AS number_of_customers
FROM customer_type
GROUP BY customer_segment;
```

	customer_segment	number_of_customers
1	Loyal	3116
2	New	83
3	Returning	701

Answer:

This query categorizes customers into three behavioral segments:

- **New Customers:** 1 previous purchase
- **Returning Customers:** 2–10 previous purchases
- **Loyal Customers:** more than 10 previous purchases

It then counts how many customers fall into each group.

Insights:

- This segmentation helps identify the **distribution of customer loyalty** across the dataset.
 - A high number of *New* customers indicates strong acquisition.
 - A strong *Returning* segment shows good repeat activity.
 - A sizable *Loyal* segment reflects **consistent customer satisfaction and long-term engagement**.
 - These insights can guide marketing strategies such as targeted promotions, loyalty programs, and retention campaigns.
-

Q8. What are the top 3 most purchased products within each category?

```
WITH item_counts AS (
    SELECT
        category,
        item_purchased,
        COUNT(customer_id) AS total_orders,
        ROW_NUMBER() OVER (
            PARTITION BY category
            ORDER BY COUNT(customer_id) DESC
        ) AS item_rank
    FROM customer
    GROUP BY category, item_purchased
)

SELECT
    item_rank,
    category,
    item_purchased,
    total_orders
FROM item_counts
WHERE item_rank <= 3
ORDER BY category, item_rank;
```

Answer:

This query identifies the **top 3 most purchased products in each category** by:

- Counting total orders for every product inside each category.
- Ranking those products using ROW_NUMBER() based on highest to lowest purchase count.
- Selecting only the top-ranked 3 products per category.

Insights:

- These items represent the **highest-demand products** within their respective categories.
- They are ideal candidates for:
 - priority stocking and supply-chain focus,
 - category-level promotions,
 - featured listings on product pages,
 - cross-selling and bundling strategies.
- Understanding bestselling items helps optimize category performance and improve revenue impact.

	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
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

Q9. Are customers who are repeat buyers (more than 5 previous purchases) also likely to subscribe?

Query:

```
SELECT
    subscription_status,
    COUNT(customer_id) AS repeat_buyers
FROM customer
WHERE previous_purchases > 5
GROUP BY subscription_status;
```

Answer:

	subscription_status	repeat_buyers
1	No	2518
2	Yes	958

Insights:

- If the count of **subscribed repeat buyers** is significantly higher than non-subscribed ones, it indicates that **repeat buying behavior is strongly associated with subscription adoption**.
 - **High-frequency buyers** often see more value in subscribing (discounts, perks, convenience).
 - Subscriptions may help retain already engaged customers and further increase their spending.
- If non-subscribed repeat buyers are also high, this could signal an **opportunity to target them with subscription campaigns**.

Conclusion:

The results of this query will show whether repeat buyers tend to subscribe.

A higher number of subscribers in this repeat-buyer group would suggest that subscription is appealing to loyal and high-activity customers.

Q10. What is the revenue contribution of each age group?

Query:

```
SELECT
    age_group,
    SUM(purchase_amount) AS total_revenue
FROM customer
GROUP BY age_group
ORDER BY total_revenue DESC;
```

Answer:

This query calculates the **total revenue generated by each age group** by summing the purchase amounts and then ranking the age groups from highest to lowest revenue contribution.

	age_group 	total_revenue 
1	Young Adult	62143
2	Middle Aged	59197
3	Adult	55978
4	Senior	55763

Insights:

- The age group with the **highest revenue** represents the most financially valuable customer segment.
- Understanding which age group contributes the most helps businesses:
 - Target marketing campaigns effectively
 - Tailor product offerings to the most profitable demographic
 - Identify underperforming age groups and potential growth opportunities
- If younger age groups contribute less, strategies may focus on awareness and acquisition.

- If older or mid-aged groups dominate, retention and loyalty programs may be more impactful.
-