# Sales Insights

E-Commerce customers
& transactions dataset

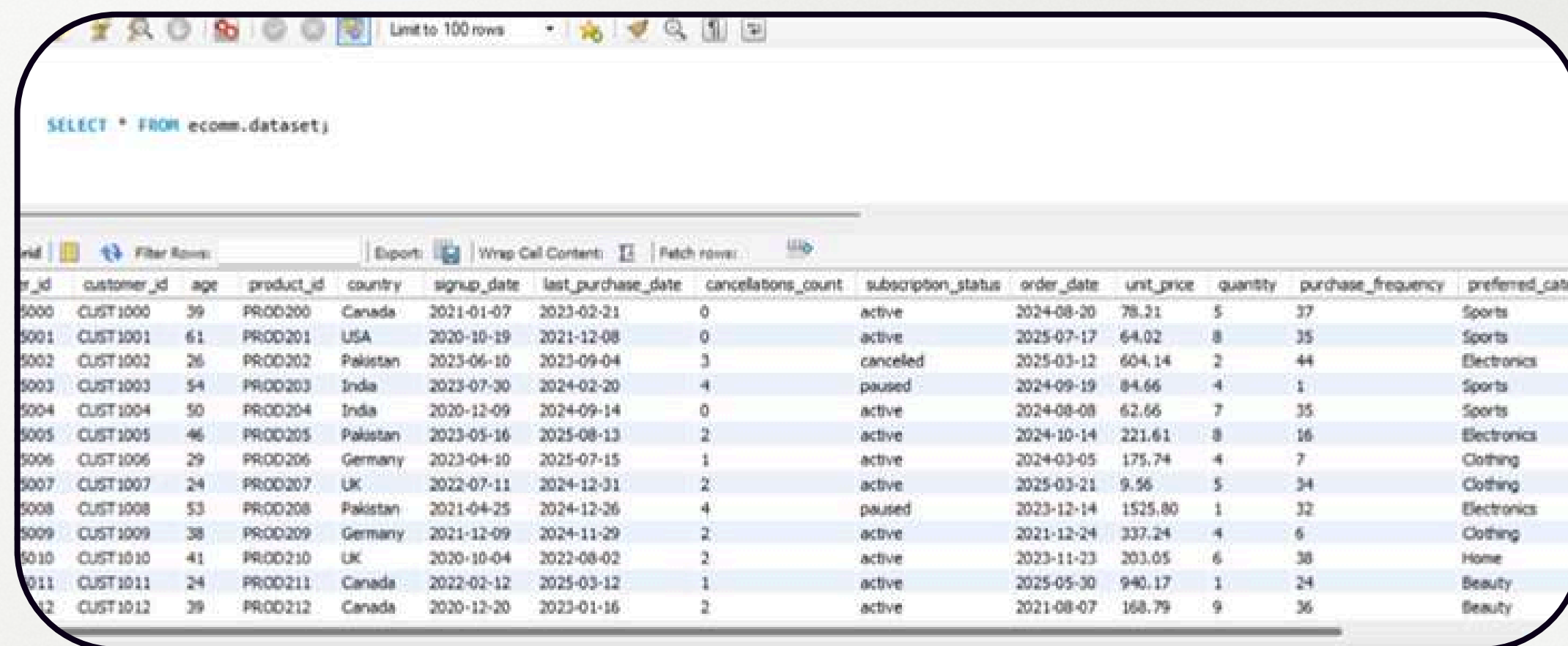Exploratory data analysis using MySQL

by ks.milda.ks@gmail.com

# Table of contents

# Dataset Introduction

This dataset contains 2,000 customer-level and order-level records from a simulated e-commerce platform (from kaggle). The customers ranging from several countries, age, gender, and order date.

Before the exploration begins, the dataset is checked and cleaned to ensure there are no missing values nor duplicate rows, as documented in this file. The process include formatting the data type (normalizing) and setting the column order_id as the primary key.

# Number of Orders

The number of order were increase every year, which may indicate growing customer engagement.

Total orders

2000

Number of Orders by Year



```
8 ●   SELECT COUNT(DISTINCT order_id) AS number_of_orders FROM ecomm.dataset;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content:

| number_of_orders |
| --- |
| 2000 |

```
8 ●   SELECT COUNT(DISTINCT order_id),YEAR(order_date) AS number_of_orders FROM ecomm.dataset GROUP BY YEAR(order_d
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content:

| COUNT(DISTINCT order_id) | number_of_orders |
| --- | --- |
| 8 | 2020 |
| 136 | 2021 |
| 304 | 2022 |
| 545 | 2023 |
| 638 | 2024 |
| 369 | 2025 |

# Number of Orders (2)

The number of orders based on each country. Picture below showed that country with the most number of orders were Germany with 360 orders.

Total orders

**2000**

Number of Orders by Country



```
12      ## Total orders per country
13  •   SELECT country, COUNT(order_id) AS total_orders FROM ecomm.dataset
14      GROUP BY country ORDER BY total_orders DESC;
15
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content:

| country | total_orders |
|---------|--------------|
| Germany | 360 |
| UK | 350 |
| Pakistan | 332 |
| India | 324 |
| USA | 319 |
| Canada | 315 |

# Average Order Value

The AOV for every country means on average each transaction contributes as number below in revenue. Countries like India and Pakistan show higher AOVs, while the UK has lower AOV.

AOV

$1025.85

AOV by Country



```
2  •  ⊝  SELECT round(sum(unit_price * quantity) /
33         COUNT(DISTINCT order_id),2) AS "AOV" FROM ecomm.dataset;
```

Result Grid | Filter Rows:

| AOV |
|---|
| 1025.85 |

```
## Average Order Value (AOV)
31  •  SELECT round(sum(unit_price * quantity) / COUNT(DISTINCT order_id),2) AS "AOV",
32         country FROM ecomm.dataset GROUP BY country ORDER BY AOV Desc;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: IA

| AOV | country |
|---|---|
| 1039.48 | India |
| 1033.62 | Pakistan |
| 1030.27 | USA |
| 1022.91 | Germany |
| 1020.69 | Canada |
| 1009.46 | UK |

# Total Revenue

The total revenue for 5 years based on country as attached below, show that the percentage revenue for each country is relatively balanced, with no single country accounting for more than one-fifth of total sales. This suggests the business has a diversified customer base across markets.

**Total Revenue**

$2.051.690,65.

**Revenue by Country**



- Canada 15.7%
- Germany 17.9%
- USA 16%
- UK 17.2%
- India 16.4%
- Pakistan 16.7%

```sql
11    ## Total revenue
12 •  SELECT round(sum(unit_price * quantity),2) AS revenue FROM ecomm.dataset;
13
14
```

Result Grid | Filter Rows: | Export: | Wrap Ce

| revenue |
|---|
| 2051690.65 |

```sql
      ## Revenue per country
18 •  SELECT country, SUM(unit_price * quantity) AS total_revenue,
19        ROUND(SUM(unit_price * quantity) * 100.0 /
20            (SELECT SUM(unit_price * quantity) FROM ecomm.dataset), 2) AS revenue_percentage
21    FROM ecomm.dataset
22    GROUP BY country
23    ORDER BY total_revenue DESC;
24
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: IA

| country | total_revenue | revenue_percentage |
|---|---|---|
| Germany | 368249.31 | 17.95 |
| UK | 353312.03 | 17.22 |
| Pakistan | 343162.68 | 16.73 |
| India | 336791.53 | 16.42 |
| USA | 328656.78 | 16.02 |
| Canada | 321518.32 | 15.67 |

# Total Revenue (2)

The revenue growth were increased every year, along with sales growth. This indicates that the growth is being driven by both higher transaction counts and consistent customer spending.

Total Revenue

$2.051.690,65.

```
     ## Revenue growth per year (by order date)
15 ●  SELECT YEAR(order_date) AS order_year, SUM(unit_price * quantity) AS total_revenue
16    FROM ecomm.dataset GROUP BY YEAR(order_date) ORDER BY order_year;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: IA

| order_year | total_revenue |
|------------|---------------|
| 2020 | 8903.09 |
| 2021 | 146137.08 |
| 2022 | 315063.40 |
| 2023 | 559328.21 |
| 2024 | 644013.74 |
| 2025 | 378245.13 |

Total Revenue by Year

# Customers Demographic

There are 2000 customers. By the time this dataset used (early second half of 2025), based on subscription status there is 1204 active customers and there number of customers with more than 1 cancellation orders are 1308 cust.

### Subscription Status



cancelled
24.7%

paused
15.2%

active
60.2%

```
## Number of customers by subscription status
27  ● SELECT subscription_status, COUNT(order_id) AS total_Cust FROM ecomm.dataset
28    WHERE subscription_status IN ('active','paused','cancelled') GROUP BY subscription_status;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: 

| subscription_status | total_Cust |
|---|---|
| active | 1204 |
| paused | 303 |
| cancelled | 493 |

```
32  ● SELECT COUNT(cancellations_count) FROM ecomm.dataset WHERE cancellations_count > 1;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: 

| COUNT(cancellations_count) |
|---|
| 1308 |

The numbers of customers with more than one cancellation are 1308 customers.

# Customers Demographic (2)

There are several high-value customers doing repeated cancellations (3–5 times each) while still generating close to $2,000 in revenue. This suggests that cancellations do not necessarily equate to lost customers — instead, they highlight friction points in the purchase journey for otherwise loyal and profitable users.

```sql
38  •  SELECT customer_id, MAX(COALESCE(cancellations_count,0)) AS total_cancellations,
39         SUM(unit_price * quantity) AS total_revenue FROM ecomm.dataset
40         GROUP BY customer_id HAVING total_cancellations > 1 ORDER BY total_revenue DESC;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: | Fetch rows:

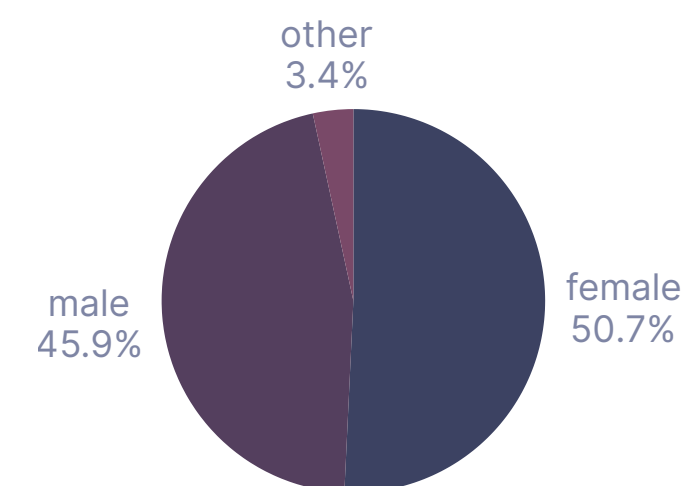| customer_id | total_cancellations | total_revenue |
|---|---|---|
| CUST1323 | 3 | 1998.08 |
| CUST2934 | 5 | 1994.70 |
| CUST1334 | 4 | 1993.11 |
| CUST2273 | 5 | 1991.63 |
| CUST2152 | 3 | 1991.34 |
| CUST1363 | 4 | 1988.67 |
| CUST2268 | 3 | 1988.48 |
| CUST1178 | 3 | 1986.96 |
| CUST1144 | 3 | 1985.20 |
| CUST2243 | 3 | 1984.80 |

Result 16 ✕

# Customers Demographic (3)

The average age of customers is 44 years, indicating that most of them are middle-aged individuals with an assumably stable income so they may be more capable of purchasing mid- to high-priced products. The gender distribution is nearly equal. Since both genders are equally represented, segmentation by gender may not be necessary — other factors (like age or income) might provide better targeting opportunities.

**Customer's Gender**

other 3.4%

female 50.7%

male 45.9%

```
33      ## Average age of customers
34  •   SELECT ROUND(AVG(age),1) AS avg_customer_age FROM ecomm.dataset;
35
```

Result Grid | Filter Rows:

| avg_customer_age |
|------------------|
| 44.1 |

```
45  •   SELECT gender, COUNT(DISTINCT customer_id) AS total_customers,
46      ROUND(COUNT(DISTINCT customer_id) * 100.0 / (SELECT COUNT(DISTINCT customer_id) FROM ecomm.dataset), 2)
47      AS percentage FROM ecomm.dataset GROUP BY gender;
```
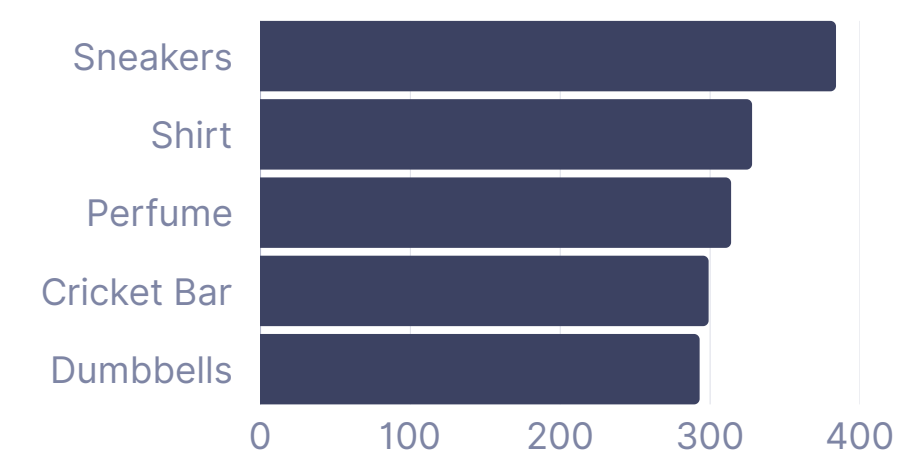
Result Grid | Filter Rows: | Export: | Wrap Cell Content: 

| gender | total_customers | percentage |
|--------|-----------------|------------|
| Female | 1015 | 50.75 |
| Male | 917 | 45.85 |
| Other | 68 | 3.40 |

# Product & Category

The top 5 products based on revenue show that "Sneakers" earned the highest revenue of $75,622.14 and were the most sold product with 384 units sold.

## Most Purchased Product



```
18    ## Top 5 products based on revenue
19  ● SELECT product_name, SUM(unit_price * quantity) AS revenue
20    FROM ecomm.dataset GROUP BY product_name ORDER BY revenue DESC LIMIT 5;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: | Fetch rows:

| product_name | revenue |
|---|---|
| Sneakers | 75622.14 |
| Skirt | 69781.61 |
| Dumbbells | 66956.52 |
| Smartphone | 58772.23 |
| Lipstick | 58218.68 |

```
      ## top 5 most purchased product
66  ● SELECT product_name, SUM(quantity) AS total_quantity_sold FROM ecomm.dataset
67    GROUP BY product_name ORDER BY total_quantity_sold DESC LIMIT 5;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: | Fetch rows:

| product_name | total_quantity_sold |
|---|---|
| Sneakers | 384 |
| Shirt | 328 |
| Perfume | 314 |
| Cricket Bat | 299 |
| Dumbbells | 293 |

# Product & Category (2)

The purchase category by gender distribution showed the most purchased were "electronics" and "clothing" by male, and "home" and "clothing" by female. The "clothing" category being one of the most purchase category were align with product "shirt" as in top 5 most purchased product.

The most average purchase frequency by category were "clothing" and "sports" as both of them leading in top 10 most purchase category.

```
51  ●   SELECT category, gender, COUNT(order_id) AS total_orders FROM ecomm.datase
52      GROUP BY category, gender ORDER BY category, total_orders DESC;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: IA

| category | gender | total_orders |
|---|---|---|
| Sports | Female | 204 |
| Sports | Male | 173 |
| Sports | Other | 13 |
| Electronics | Male | 211 |
| Electronics | Female | 193 |
| Electronics | Other | 10 |
| Clothing | Male | 207 |
| Clothing | Female | 205 |
| Clothing | Other | 14 |

```
3       ## Average purchase frequency by category
74  ●   SELECT category, AVG(purchase_frequency) AS avg_purchase_freq
```
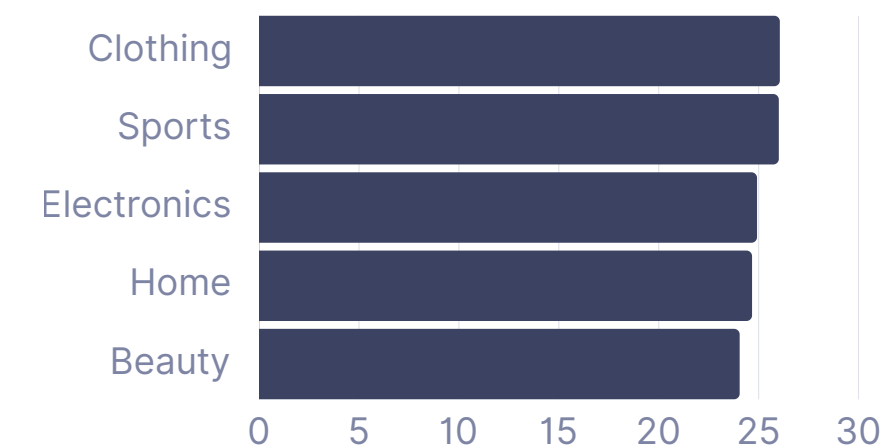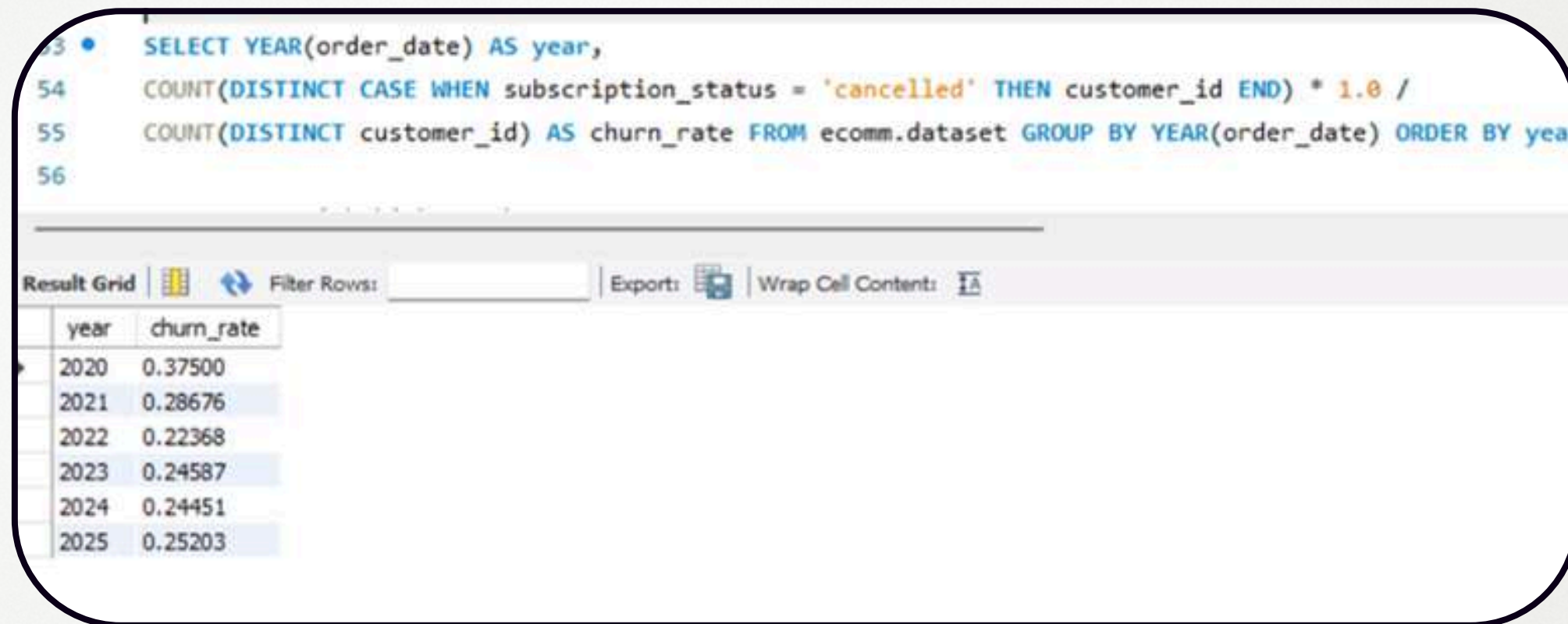
Result Grid | Filter Rows: | Export: | Wrap Cell Content: IA

| category | avg_purchase_freq |
|---|---|
| Clothing | 26.0610 |
| Sports | 26.0103 |
| Electronics | 24.9227 |
| Home | 24.6738 |
| Beauty | 24.0505 |

**Avg Purchase Frequency**

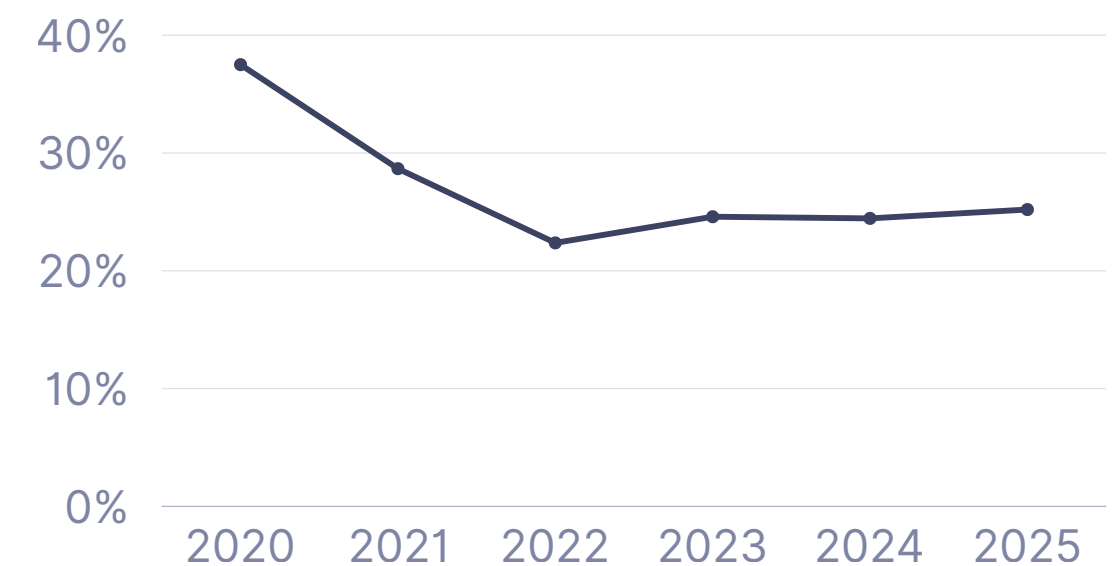| Category | 0 | 5 | 10 | 15 | 20 | 25 | 30 |
|---|---|---|---|---|---|---|---|
| Clothing | | | | | | | |
| Sports | | | | | | | |
| Electronics | | | | | | | |
| Home | | | | | | | |
| Beauty | | | | | | | |

# Churn Rate

The churn rate incredibly declined from 37% in 2020 to 22% in 2022 indicating some improvement in customer experience. However, the churn rate start to increased in 2023, suggesting that retention efforts weakened as the customer base grew.

```
53 •  SELECT YEAR(order_date) AS year,
54      COUNT(DISTINCT CASE WHEN subscription_status = 'cancelled' THEN customer_id END) * 1.0 /
55      COUNT(DISTINCT customer_id) AS churn_rate FROM ecomm.dataset GROUP BY YEAR(order_date) ORDER BY year
56
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: 𝐈𝐀

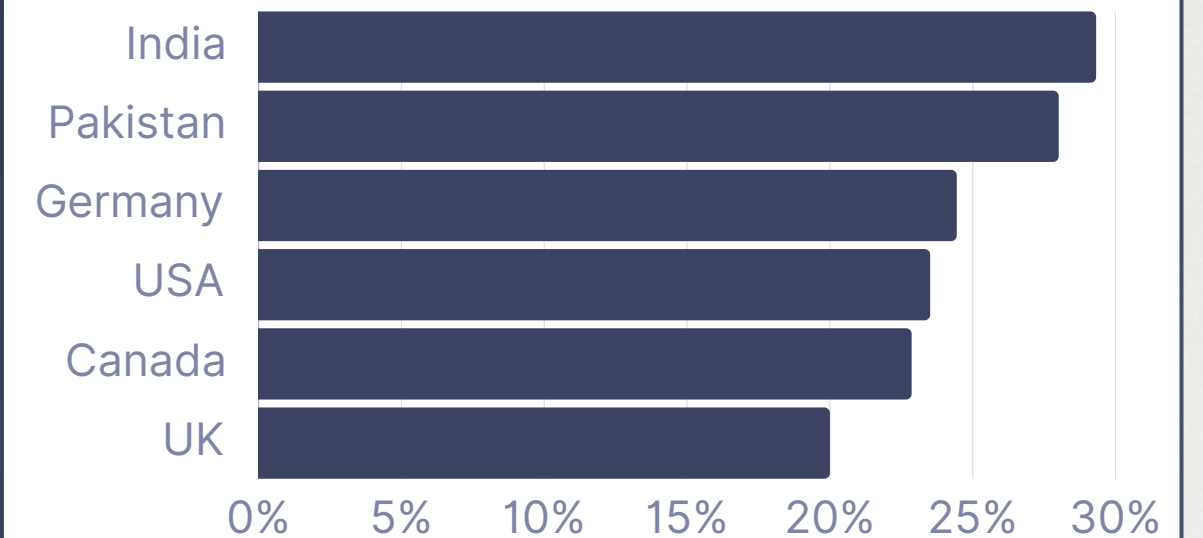| year | churn_rate |
|------|-----------|
| 2020 | 0.37500 |
| 2021 | 0.28676 |
| 2022 | 0.22368 |
| 2023 | 0.24587 |
| 2024 | 0.24451 |
| 2025 | 0.25203 |

**Churn Rate**

# Churn Rate (2)

India has the highest churn rate (0.293), indicating that nearly 29% of customers in India canceled their subscriptions. This suggests a possible issue with customer satisfaction, pricing, or local engagement strategies compared to other countries.

```
39   SELECT country,
40       COUNT(DISTINCT CASE WHEN subscription_status='cancelled' THEN customer_id END) * 1.0 /
41       COUNT(DISTINCT customer_id) AS churn_rate FROM ecomm.dataset
42   GROUP BY country ORDER BY churn_rate DESC;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content:

| country | churn_rate |
| --- | --- |
| India | 0.29321 |
| Pakistan | 0.28012 |
| Germany | 0.24444 |
| USA | 0.23511 |
| Canada | 0.22857 |
| UK | 0.20000 |

### Churn Rate by Country

| Country | |
| --- | --- |
| India | |
| Pakistan | |
| Germany | |
| USA | |
| Canada | |
| UK | |

0%  5%  10%  15%  20%  25%  30%

Total Orders
2000

Total Revenue
$2.051.690,65.

Average Order Value
$1025.85

## Churn Rate
**+0.7%** From last year • • • •



## Customer Retention Rate • • •



**74.8%**

**-0.75%** From last year

### Customer Subscription Status

- Active
- Paused
- Cancelled

15.2%
24.7%
60.2%

### Customer Gender

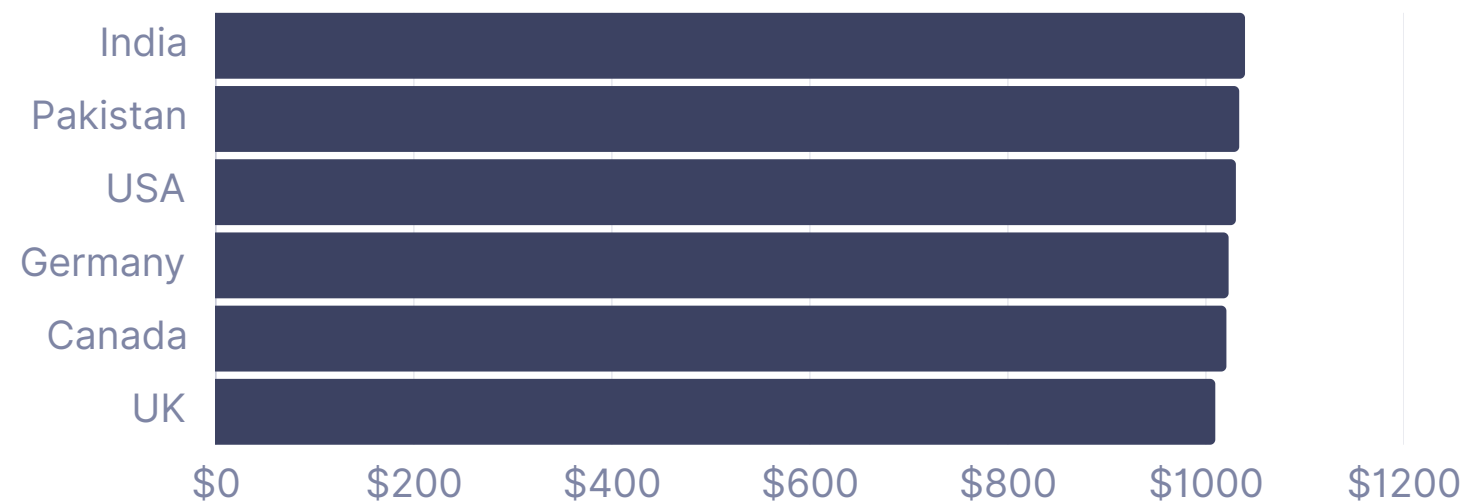- Female
- Male
- Others

3.4%
45.9%
50.7%

# Sales Performance

## AOV by Country



## Total Revenue by Year (H1 2025)



## Revenue by Country



Canada 15.7%
Germany 17.9%
USA 16%
UK 17.2%
India 16.4%
Pakistan 16.7%

## Top Country by Number of Order



Germany
UK
Pakistan
India
USA
Canada

**Insights**

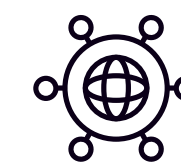Germany contributes the most orders; India and Pakistan show high AOV.

# Conclusion

## Exceeding Expectations

The business demonstrates strong overall performance with consistent growth in both orders and total revenue each year, reflecting increasing customer engagement. Revenue is well distributed across multiple countries, reducing dependency on a single market, while India and Pakistan stand out for their high Average Order Values (AOV). The customer base is diverse and stable, with an average age of 44 years and nearly equal gender distribution. Product performance is also promising, with "Sneakers" leading in both sales and revenue, supported by the popularity of categories such as Clothing and Sports.

## Areas Needing Attention

Despite earlier improvements, the churn rate began to rise again after 2023, suggesting that retention strategies have weakened as the customer base expanded. Some high-value customers were found to cancel multiple orders despite contributing significant revenue, indicating potential friction in the purchase or delivery process. Additionally, churn varies by country, highlighting inconsistencies in customer satisfaction and retention efforts across regions. These areas require closer examination to prevent future customer loss.

## Focus Areas

To sustain growth, the company should prioritize strengthening customer retention and loyalty programs while addressing issues behind repeated cancellations. Improving user experience, especially during checkout and after purchase, will be essential to maintain long-term satisfaction. Marketing strategies can be refined through demographic segmentation, leveraging insights from age, gender, and product preferences. Furthermore, focusing on high-performing product categories like Clothing and Sports while exploring underperforming ones will help balance growth. Continuous monitoring of churn and retention by country and customer group will ensure data-driven decision-making and consistent improvement.

# Thank You

___

📞 +62-851-7990-6131

🌐 Linkedin milda-s

✉️ ks.milda.ks@gmail.com

📍 Yogyakarta, Indonesia