



OLIST SQL PROJECT: E-COMMERCE DATA ANALYSIS

INDEX

1. INTRODUCTION	3
2. ENTITY RELATIONSHIP DIAGRAM (ERD)	4
3. DATA CLEANING & TRANSFORMATION	5
4. ANALYSIS 1: Order Insights	8
4.1 List of Available Payment Methods	
4.2 Order Status Distribution	
4.3 Total Number of Orders	
4.4 Total Number of Delivered Orders	
5. ANALYSIS 2: Payment and Delivery Analysis	11
5.1 Most Frequently Used Payment Method	
5.2 Average Delivery Time	
5.3 Top Product Categories by Revenue	
6. ANALYSIS 3: Seller & Customer Patterns	13
6.1 Top Revenue-Generating Sellers	
6.2 Customer Distribution by State	

1.INTRODUCTION

This project focuses on performing SQL-based analysis on an e-commerce platform's historical data. The primary goal is to gain deep insights into customer behavior, product trends, seller performance, and delivery efficiency using relational data analysis.

The dataset includes multiple interconnected tables, such as orders, products, customers, sellers, and payments, which represent different aspects of the business operations. Using SQL, we perform structured analysis across these datasets to solve real business problems.

Purpose of the Project

The objective of this project is to:

- **Understand customer behavior** and purchasing patterns across different regions.
- **Analyze order fulfilments** by evaluating delivery rates and average delivery times.
- **Identify high-performing sellers** and top-selling product categories.
- **Evaluate payment trends** to understand customer preferences in transaction methods.
- **Clean and prepare** raw datasets to ensure data integrity for accurate analysis.

The outcome of this analysis will support data-informed strategies for marketing, operations, inventory management, logistics, and customer service optimization.

Scope of Work

This project is divided into the following major phases:

1. **Exploratory Data Review:** Understanding the structure, content, and relationships of multiple datasets (orders, payments, products, customers, etc.).
2. **Data Cleaning & Transformation:** Creating filtered and standardized versions of original tables to remove anomalies and ensure consistency.
3. **SQL Analysis:** Writing and executing queries to answer real business questions, such as:
 - Which products generate the most revenue?
 - Which sellers drive the highest sales?
 - What is the most used payment method?
 - How efficient is the delivery process?
4. **Insight Documentation:** Presenting results visually and descriptively with interpretations and recommendations.

Technology Used

- **SQL (PostgreSQL flavor):** For data manipulation and querying.
- **DB Management Tool (e.g., pgAdmin):** For executing and visualizing queries.
- **ERD Tool / Diagram Editor:** For designing the entity relationship model.

Challenges Faced

Working with real-world e-commerce data presented several challenges:

1. **Data Inconsistencies:**
 - Variations in text formatting (e.g., state abbreviations in mixed case).
 - Presence of null values in critical fields like `order_id`, `product_id`.
2. **Missing or Redundant Records:**
 - Duplicates and empty records required careful filtering and deduplication using `DISTINCT` and `WHERE` clauses.
3. **Multi-table Joins and Relationships:**
 - Understanding the complex relationships between customers, orders, products, and sellers required a well-planned ERD.
4. **Volume and Complexity:**
 - The size and interconnectedness of tables demanded efficient querying strategies to ensure performance and accuracy.
5. **Time-Based Calculations:**
 - Calculating delivery time accurately involved timestamp conversions and interval calculations.

Expected Outcomes

By the end of this project, we expect to have:

- A clear overview of sales trends and customer distribution.
- Key metrics on order processing and delivery efficiency.
- Identification of top-performing sellers and categories.
- Recommendations for business improvements based on data-backed insights.

2. ENTITY RELATIONSHIP DIAGRAM (ERD)

The ERD (Entity Relationship Diagram) visually represents how different tables in the database are connected through primary and foreign keys. It helps in understanding the data structure and planning queries effectively.

Entities & Relationships:

- **Customers** place **Orders**
- Each **Order** has **Order Items**
- **Order Items** reference **Products**
- **Orders** are linked to **Payments**
- **Products** belong to **Categories**
- **Sellers** fulfil **Order Items**
- **Geolocation** links to zip codes of both customers and sellers

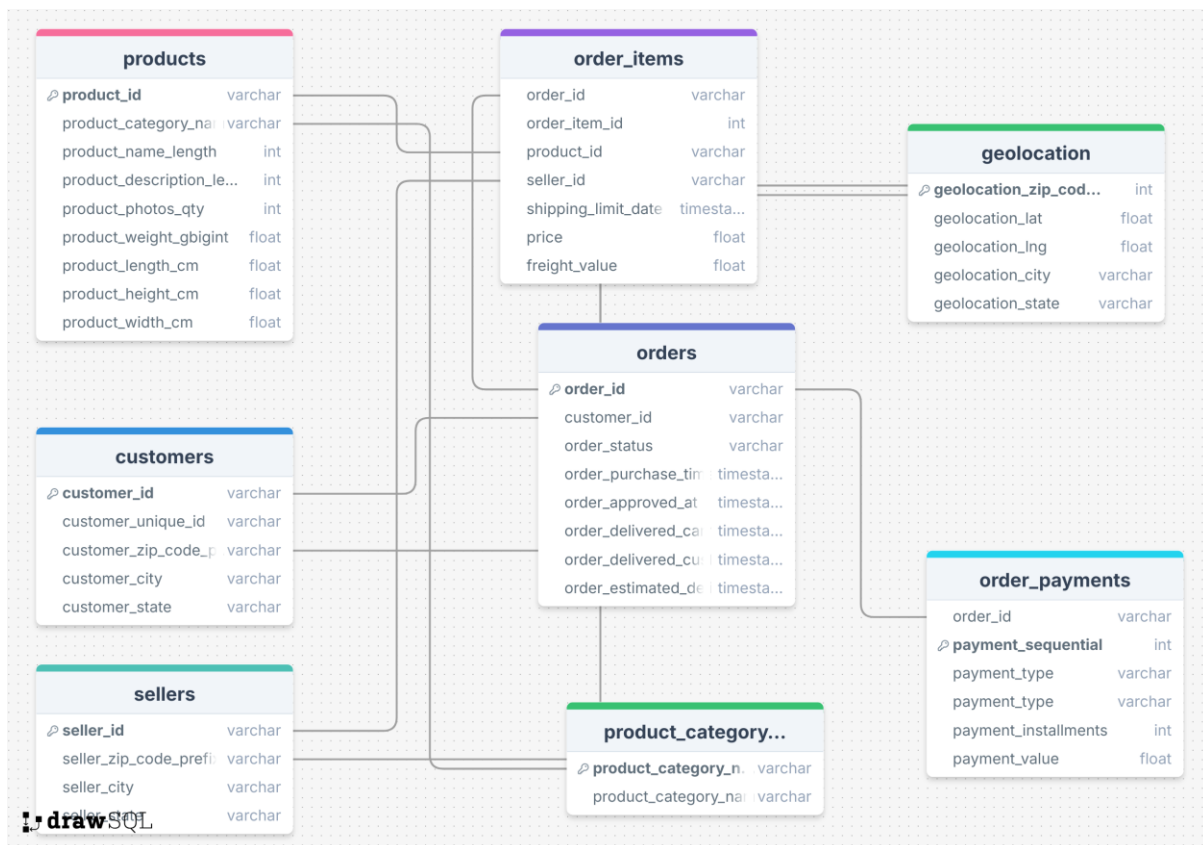


Fig 2.1 Entity Relationship Diagram

3. DATA CLEANING & TRANSFORMATION

Raw data often contains inconsistencies such as:

- Null values
- Duplicates
- Irregular formatting (e.g., mixed-case states or cities)
- Invalid entries (like zero prices or weights)

Before diving into analysis, it's essential to create **clean versions** of the tables by applying filters, removing nulls, standardizing formats, and ensuring logical integrity. This ensures all insights generated are **reliable** and **accurate**.

Cleaning Steps:

1. **clean_customers**
 - Removed nulls
 - Standardized city to lowercase and state to uppercase

```
CREATE TABLE clean_customers AS
SELECT DISTINCT
    customer_id,
    customer_unique_id,
    TRIM(LOWER(customer_city)) AS customer_city,
    TRIM(UPPER(customer_state)) AS customer_state,
    customer_zip_code_prefix
FROM customers
WHERE customer_id IS NOT NULL;
```

2. **clean_orders**

- Filtered out orders with missing IDs, customer IDs, or status

```
CREATE TABLE clean_orders AS
SELECT DISTINCT *
FROM orders
WHERE customer_id IS NOT NULL
    AND order_status IS NOT NULL
    AND order_id IS NOT NULL;
```

3. **clean_order_items**

- Ensured non-null order/product IDs
- Kept only items with positive prices

```
CREATE TABLE clean_order_items AS
SELECT DISTINCT *
FROM order_items
WHERE order_id IS NOT NULL
    AND product_id IS NOT NULL
    AND price > 0;
```

4. **clean_order_payments**

- Removed nulls
- Excluded zero-value transactions

```
CREATE TABLE clean_order_payments AS
SELECT DISTINCT *
FROM order_payments
WHERE payment_value > 0
    AND payment_type IS NOT NULL
    AND order_id IS NOT NULL;
```

5. **clean_products**

- Removed products with missing IDs or weights
- Ensured valid product category names

```
CREATE TABLE clean_products AS
SELECT DISTINCT *
FROM products
WHERE product_id IS NOT NULL
    AND product_weight_g > 0
    AND product_category_name IS NOT NULL;
```

6. **clean_product_category_translation**

- Ensured valid translations (both original and English names)

```
CREATE TABLE clean_product_category_translation AS
SELECT DISTINCT *
FROM product_category_translation
WHERE product_category_name IS NOT NULL
    AND product_category_name_english IS NOT NULL;
```

7. **clean_sellers**

- Cleaned and standardized seller city/state like customers

```
CREATE TABLE clean_sellers AS
SELECT DISTINCT
    seller_id,
    TRIM(LOWER(seller_city)) AS seller_city,
    TRIM(UPPER(seller_state)) AS seller_state,
    seller_zip_code_prefix
FROM sellers
WHERE seller_id IS NOT NULL;
```

8. **clean_geolocation**

- Filtered for valid lat/long and zip code prefixes

```
CREATE TABLE clean_geolocation AS
SELECT DISTINCT *
FROM geolocation
WHERE geolocation_zip_code_prefix IS NOT NULL
AND geolocation_lat IS NOT NULL
AND geolocation_lng IS NOT NULL;
```

ANALYSIS 4: ORDER INSIGHTS

4.1 List of Available Payment Methods

ProblemStatement:

Identify all distinct payment methods available in the dataset to understand the range of customer payment preferences.

Approach:

- Query the payment table.
- Use DISTINCT to eliminate duplicate records and list unique values in the payment_type column.

Analysis:

The query helps in understanding the types of transactions customers are most comfortable using. This insight can guide future integrations or promotions.

Recommendation:

Ensure all popular payment methods remain supported, and consider introducing newer digital wallets to enhance flexibility.

Query		Query History
1	---Which Payment method is most used---	
2	SELECT	payment_type , COUNT(*) AS usage_count
3	FROM	clean_order_payments
4	GROUP BY	payment_type
5	ORDER BY	usage_count desc
Data Output		Messages Explain X Notifications
<div> <div> <div>SQL</div> </div> </div>		
Showing rows: 1 to 4 Page No: 1 of 1		
	payment_type character varying	usage_count bigint
1	credit_card	76795
2	boleto	19784
3	voucher	5769
4	debit_card	1529

4.2 Order Status Distribution

ProblemStatement:

Determine how many orders fall into each status (e.g., delivered, shipped, canceled).

Approach:

- Group records from the orders table using order_status.
- Use COUNT() to get the frequency of each status.

Analysis:

This identifies operational efficiency and highlights problem areas like cancellations or processing delays.

Recommendation:

Track status changes in real-time and implement interventions for delayed or canceled orders to improve fulfilments rates.

Query

Query History

1

---Total number of orders per order status---

2

SELECT order_status, count(*)

3

FROM clean_orders

4

GROUP BY order_status

Data Output

Messages

Explain

×

Notifications

Showing rows: 1 to 8

Page No: 1 of 1

	order_status character varying	count bigint
1	shipped	1107
2	unavailable	609
3	invoiced	314
4	created	5
5	approved	2
6	processing	301
7	delivered	96478
8	canceled	625

4.3 Total Number of Orders

ProblemStatement:

Find the total number of orders to understand the dataset's transaction volume.

Approach:

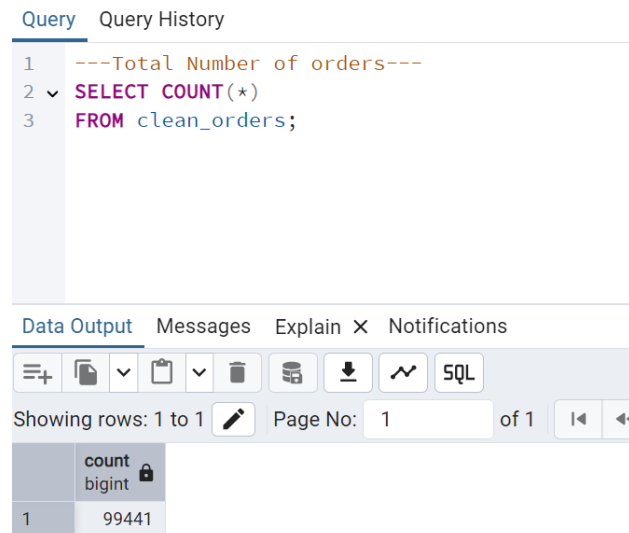
- Count the number of records in the orders table using COUNT(*).

Analysis:

Serves as a benchmark for further calculations such as delivery rate and customer engagement.

Recommendation:

Use total orders as a base metric for performance KPIs and planning for scalability.



The screenshot shows a SQL query editor with a query history tab. The query is: `---Total Number of orders---`
`SELECT COUNT(*)`
`FROM clean_orders;`

Below the query editor, the 'Data Output' tab is active, showing the results of the query. The results are displayed in a table with one row and one column.

	count bigint
1	99441

4.4 Total Number of Delivered Orders

ProblemStatement:

Count how many orders have been successfully delivered to measure order fulfilment success.

Approach:

- Filter the orders table for order_status = 'delivered'.
- Use COUNT() to calculate the total.

Analysis:

Provides insights into delivery performance. A high ratio of delivered orders indicates a smooth logistics process.

Recommendation:

Maintain or improve delivery standards by monitoring courier performance and warehouse processes.

5. ANALYSIS 2: PAYMENT AND DELIVERY ANALYSIS

5.1 Most Frequently Used Payment Method

ProblemStatement:

Identify the most commonly used payment method by customers.

Approach:

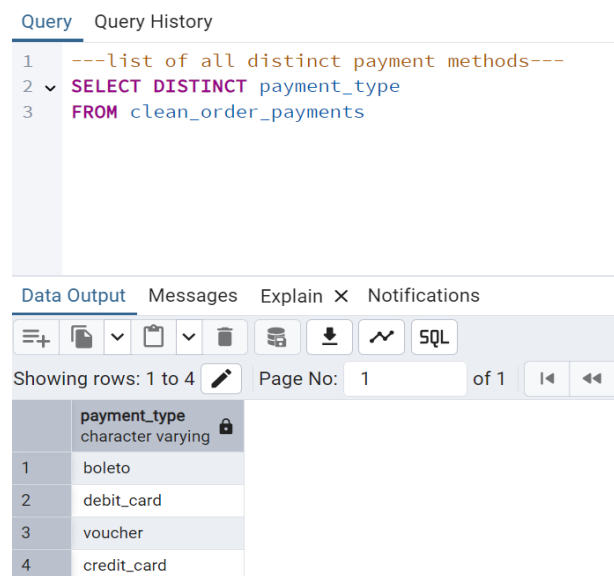
- Group the order_payments data by payment_type.
- Use COUNT() to find frequency.
- Order by count descending.

Analysis:

Shows customer preference trends and trust in specific payment types (e.g., credit card dominance).

Recommendation:

Promote frequently used payment methods with cashback or loyalty offers and ensure a smooth checkout process for them.



The screenshot shows a SQL query editor with a query history tab. The query is: `SELECT DISTINCT payment_type FROM clean_order_payments`. Below the query, there is a 'Data Output' tab showing the results of the query. The output is a table with one column, 'payment_type', and four rows of data: 'boleto', 'debit_card', 'voucher', and 'credit_card'.

	payment_type character varying
1	boleto
2	debit_card
3	voucher
4	credit_card

5.2 Average Delivery Time

ProblemStatement:

Calculate the average time it takes from order placement to delivery.

Approach:

- Use DATEDIFF or time subtraction between order_purchase_timestamp and order_delivered_customer_date.
- Apply AVG() to get the mean duration.

Analysis:

This metric reflects logistics performance and affects customer satisfaction.

Recommendation:

If average delivery time exceeds industry standards, invest in faster shipping partners or regional warehouses.

The screenshot shows a SQL query editor with the following query:

```
1 ---What is the average delivery---
2 ---time for delivered orders--
3 SELECT AVG(order_delivered_customer_date - order_purchase_timestamp) as order_status
4 FROM clean_orders
5 WHERE order_status = 'delivered'
```

Below the query editor, the 'Data Output' tab is active, showing the results of the query. The results are displayed in a table with the following structure:

order_status	interval
1	12 days 13:23:49.957272

5.3 Top Product Categories by Revenue

ProblemStatement:

Identify which product categories bring in the most revenue.

Approach:

- Join order_items with products and product_category_translation.
- Use SUM(price * quantity) and group by category.

Analysis:

Reveals best-selling product segments. Useful for inventory and promotion planning.

Recommendation:

Allocate more marketing budget to high-performing categories and optimize stock levels accordingly.

Query	Query History
1	---Which Product categories generate
2	---the highest revenue
3	SELECT pct.product_category_name_english,
4	SUM(oi.price) AS total_revenue
5	FROM clean_order_items oi
6	JOIN clean_products p ON oi.product_id = p.product_id
7	JOIN clean_product_category_translation pct ON p.product_category_name = pct.product_category_name
8	GROUP BY pct.product_category_name_english
9	ORDER BY total_revenue DESC
10	LIMIT 5;
11	
Data Output	Messages Explain X Notifications
Showing rows: 1 to 5	Page No: 1 of 1
product_category_name_english character varying	total_revenue double precision
1 health_beauty	1258681.340000017
2 watches_gifts	1205005.680000013
3 bed_bath_table	1036039.1800000409
4 sports_leisure	988048.9700000199
5 computers_accessories	911954.3200000154

6. ANALYSIS 3: SELLER & CUSTOMER PATTERNS

6.1 Top Revenue-Generating Sellers

ProblemStatement:

Determine which sellers have generated the most revenue.

Approach:

- Group by seller_id in the order_items table.
- Sum total revenue using price * quantity.

Analysis:

Highlights seller performance and can help identify reliable vendors for partnerships or promotions.

Recommendation:

Provide high-performing sellers with visibility perks and consider onboarding similar profiles

Query	Query History
1	---Who are the top 5 sellers by revenue
2	SELECT seller_id,sum(price) as total_sales
3	FROM clean_order_items
4	GROUP BY seller_id
5	ORDER BY total_sales desc
6	limit 5
Data Output	Messages Explain X Notifications
Showing rows: 1 to 5	Page No: 1 of 1
seller_id character varying	total_sales double precision
1 4869f7a5dfa277a7dca6462dcf3b52b2	229472.6299999981
2 53243585a1d6dc2643021fd1853d89...	222776.04999999952
3 4a3ca9315b744ce9f8e9374361493884	200472.91999999949
4 fa1c13f2614d7b5c4749cbc52fecda94	194042.02999999846
5 7c67e1448b00f6e969d365cea6b010ab	187923.8899999995

6.2 Customer Distribution by State

ProblemStatement:

Understand the geographic distribution of customers across states.

Approach:

- Group customers by customer_state.
- Count the number of customers in each group.

Analysis:

This shows where most of the customer base is concentrated, aiding in targeted marketing and regional logistics planning.

Recommendation:

Focus outreach in states with high customer density and explore growth opportunities in underrepresented areas.

Query

Query History

1

2

3

4

5

6

Which states have the most customers---

SELECT customer_state, COUNT(*) AS customer_count

FROM clean_customers

GROUP BY customer_state

ORDER BY customer_count DESC

LIMIT 5

Data Output

Messages

Explain X

Notifications

<