# Olist Ecommerce- Data Management & Analysis

- by Srinivas Abhilash Chintaluru

# Table of Contents

# Data Source

**Dataset:** Brazilian e-commerce dataset available in Kaggle

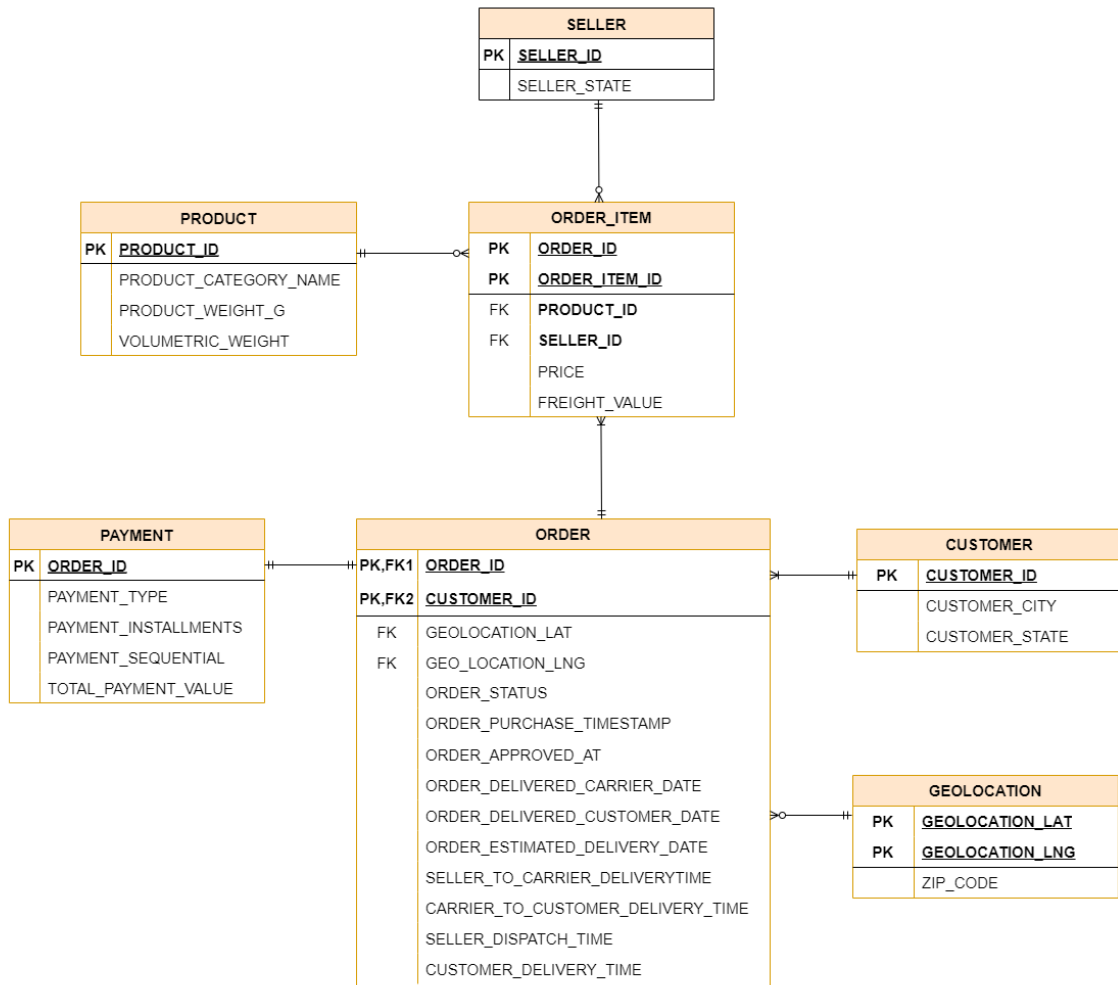A single dataset is available is normalized into seven datasets/tables as shown below

| Table name | Table Description |
|---|---|
| Orders | Information about the orders like order status, order purchase time, order approval time, etc. |
| Customer | Information about the customers like customer id, city, state, etc.. |
| Product | Information about the product like product id, product category, product weight, etc. |
| Seller | Information about the seller like seller id, seller state, etc. |
| Payment | Information about payments like payment type, payment installments, etc. |
| Delivery | Information related to delivery like the seller to carrier delivery time, carrier to customer delivery time, dispatch time etc. |
| Geolocation | Information related to latitude and longitude of the orders placed |

# Entity Relationship Diagram

**Relationships used in the ERD:**

Definition: Order item – An order can have multiple items in it, each item is assigned one order item id

1) One seller might have sold zero or more order items, and one order item has only one seller
2) One order item has only one product, while one product can be in multiple order items
3) One order can have multiple order items, while one order item belongs to only one order
4) One order has one payment, and one payment has one order
5) One order has only one customer, and one customer can place one more orders
6) One order can be placed from one location, and one location might come up have multiple orders

**SELLER**

| | |
|---|---|
| PK | SELLER_ID |
| | SELLER_STATE |

**PRODUCT**

| | |
|---|---|
| PK | PRODUCT_ID |
| | PRODUCT_CATEGORY_NAME |
| | PRODUCT_WEIGHT_G |
| | VOLUMETRIC_WEIGHT |

**ORDER_ITEM**

| | |
|---|---|
| PK | ORDER_ID |
| PK | ORDER_ITEM_ID |
| FK | PRODUCT_ID |
| FK | SELLER_ID |
| | PRICE |
| | FREIGHT_VALUE |

**PAYMENT**

| | |
|---|---|
| PK | ORDER_ID |
| | PAYMENT_TYPE |
| | PAYMENT_INSTALLMENTS |
| | PAYMENT_SEQUENTIAL |
| | TOTAL_PAYMENT_VALUE |

**ORDER**

| | |
|---|---|
| PK,FK1 | ORDER_ID |
| PK,FK2 | CUSTOMER_ID |
| FK | GEOLOCATION_LAT |
| FK | GEO_LOCATION_LNG |
| | ORDER_STATUS |
| | ORDER_PURCHASE_TIMESTAMP |
| | ORDER_APPROVED_AT |
| | ORDER_DELIVERED_CARRIER_DATE |
| | ORDER_DELIVERED_CUSTOMER_DATE |
| | ORDER_ESTIMATED_DELIVERY_DATE |
| | SELLER_TO_CARRIER_DELIVERYTIME |
| | CARRIER_TO_CUSTOMER_DELIVERY_TIME |
| | SELLER_DISPATCH_TIME |
| | CUSTOMER_DELIVERY_TIME |

**CUSTOMER**

| | |
|---|---|
| PK | CUSTOMER_ID |
| | CUSTOMER_CITY |
| | CUSTOMER_STATE |

**GEOLOCATION**

| | |
|---|---|
| PK | GEOLOCATION_LAT |
| PK | GEOLOCATION_LNG |
| | ZIP_CODE |

# Database creation, table population, and business questions

My dataset consists of a single CSV file (brazil_ecom_cleaned.csv)



**A snippet of the dataset**

**Approach:**

1) Import the dataset into the MySQL database ensuring all the data points have the correct data type.
2) Generate 6 tables using SQL queries
3) Frame business questions, and use SQL queries on the 6 tables to answer the business questions

# Database creation and table population

As my dataset is in denormalized form, I created a table for it and loaded the dataset(csv file) into this table.

**Step 1:** Created a database: ecommerce (Refer to the SQL query in DB_Table_creation.sql)

**Step 2:** Created a table: ecomm_denom_table in the ecommerce database (Refer to the SQL query in DB_Table_creation.sql)

**Step 3:** Loading data from brazil_ecom_cleaned.csv file into the ecommerce.ecomm_denom_table, below screenshots illustrate the steps of data loading

Provided the csv file path in the table data import tab



Selected the destination table, in our case it is ecommerce.ecomm_denom_table

Mapped all the columns in the csv file to the columns created in the table



Started the import process

**Data import Error 1:**

Received an error as the column appeared twice



As the column 'ORDER_DELIVERED_CARRIER_DATE' was repeating twice, I went to the previous step of mapping the columns and fixed it.

**Data import Error 2:**

After fixing this, I received another error due to the issue in the date format of date-based columns in the CSV file. Changed the date format from 'dd-mm-yyyy HH:MM:SS' to 'yyyy-mm-dd HH:MM:SS' and fixed the issue.

**Step 4:** Created 7 tables: seller, product, customer, payment, geolocation, orders, order_item (Refer to the SQL query in DB_Table_creation.sql)

**Step 5:** Load data from the denormalized table into each of the 7 tables created

# Data dictionary

**SELLER table**

| Field | Type | Null | Key | Default | Extra |
|-------|------|------|-----|---------|-------|
| SELLER_ID | varchar(100) | NO | PRI | NULL | |
| SELLER_STATE | char(2) | YES | | NULL | |

| Field | Data Description |
|-------|------------------|
| SELLER_ID | The unique ID for each seller |
| SELLER_STATE | The state that the seller belongs to |

**PRODUCT table**

| Field | Type | Null | Key | Default | Extra |
|-------|------|------|-----|---------|-------|
| PRODUCT_ID | varchar(100) | NO | PRI | NULL | |
| PRODUCT_CATEGORY_NAME | varchar(200) | YES | | NULL | |
| PRODUCT_WEIGHT_G | decimal(20,3) | YES | | NULL | |
| VOLUMETRIC_WEIGHT | decimal(20,3) | YES | | NULL | |

| Field | Data Description |
|-------|------------------|
| PRODUCT_ID | The unique ID for each product |
| PRODUCT_CATEGORY_NAME | Category of the product |
| PRODUCT_WEIGHT_G | Weight of the product |
| VOLUMETRIC_WEIGHT | Volume of the product ordered |

**PAYMENT table**

| Field | Type | Null | Key | Default | Extra |
|-------|------|------|-----|---------|-------|
| ORDER_ID | varchar(100) | NO | PRI | NULL | |
| PAYMENT_TYPE | varchar(200) | YES | | NULL | |
| PAYMENT_INSTALLMENTS | int | YES | | NULL | |
| PAYMENT_SEQUENTIAL | int | YES | | NULL | |
| TOTAL_PAYMENT_VALUE | decimal(20,3) | YES | | NULL | |

| Field | Data Description |
|-------|------------------|
| PRODUCT_ID | The unique ID for each product |
| PRODUCT_CATEGORY_NAME | Category of the product |
| PRODUCT_WEIGHT_G | Weight of the product |
| VOLUMETRIC_WEIGHT | Volume of the product ordered |

**CUSTOMER table**

| Field | Type | Null | Key | Default | Extra |
|-------|------|------|-----|---------|-------|
| CUSTOMER_ID | varchar(100) | NO | PRI | NULL | |
| CUSTOMER_CITY | varchar(100) | YES | | NULL | |
| CUSTOMER_STATE | char(2) | YES | | NULL | |

| Field | Data Description |
|-------|------------------|
| CUSTOMER_ID | The unique ID for each customer |
| CUSTOMER_CITY | City of the customer |
| CUSTOMER_STATE | State of the customer |

**GEOLOCATION table**

| Field | Type | Null | Key | Default | Extra |
|---|---|---|---|---|---|
| CUSTOMER_ZIP_CODE_PREFIX | varchar(10) | YES | | NULL | |
| GEOLOCATION_LAT | decimal(20,10) | NO | PRI | NULL | |
| GEOLOCATION_LNG | decimal(20,10) | NO | PRI | NULL | |

| Field | Data Description |
|---|---|
| CUSTOMER_ZIP_CODE_PREFIX | The unique ID for each customer |
| GEOLOCATION_LAT | City of the customer |
| CUSTOMER_STATE | State of the customer |

**ORDER_ITEM table**

| Field | Type | Null | Key | Default | Extra |
|---|---|---|---|---|---|
| ORDER_ID | varchar(100) | NO | PRI | NULL | |
| ORDER_ITEM_ID | varchar(100) | NO | PRI | NULL | |
| PRODUCT_ID | varchar(100) | YES | | NULL | |
| SELLER_ID | varchar(100) | YES | | NULL | |
| PRICE | decimal(20,3) | YES | | NULL | |
| FREIGHT_VALUE | decimal(20,3) | YES | | NULL | |

| Field | Data Description |
|---|---|
| ORDER_ID | The unique ID for each order |
| ORDER_ITEM_ID | The unique ID for each item in the order |
| PRODUCT_ID | The unique ID for the product |
| SELLER_ID | The unique ID of the seller |
| PRICE | Price of the product |
| FREIGHT_VALUE | Delivery value of the order |

**ORDERS table**

| Field | Type | Null | Key | Default | Extra |
|---|---|---|---|---|---|
| ORDER_ID | varchar(100) | NO | PRI | NULL | |
| CUSTOMER_ID | varchar(100) | NO | PRI | NULL | |
| GEOLOCATION_LAT | decimal(20,10) | YES | | NULL | |
| GEOLOCATION_LNG | decimal(20,10) | YES | | NULL | |
| ORDER_STATUS | varchar(20) | YES | | NULL | |
| ORDER_PURCHASE_TIMESTAMP | datetime | YES | | NULL | |
| ORDER_APPROVED_AT | datetime | YES | | NULL | |
| ORDER_DELIVERED_CARRIER_DATE | datetime | YES | | NULL | |
| ORDER_DELIVERED_CUSTOMER_DATE | datetime | YES | | NULL | |
| ORDER_ESTIMATED_DELIVERY_DATE | datetime | YES | | NULL | |
| SELLER_TO_CARRIER_DELIVERYTIME | decimal(20,3) | YES | | NULL | |
| CARRIER_TO_CUSTOMER_DELIVERY_... | decimal(20,3) | YES | | NULL | |
| SELLER_DISPATCH_TIME | varchar(10) | YES | | NULL | |
| CUSTOMER_DELIVERY_TIME | varchar(10) | YES | | NULL | |

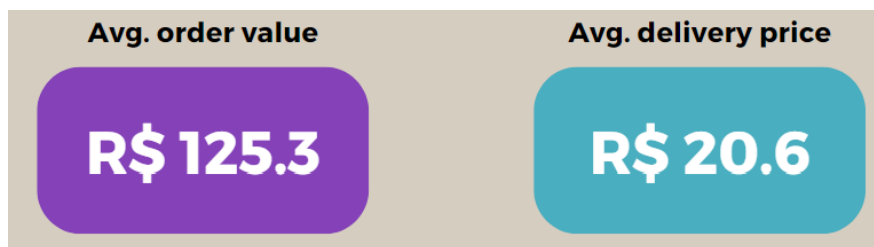| Field | Data Description |
|---|---|
| ORDER_ID | The unique ID for each order |
| CUSTOMER_ID | The unique ID for each customer |
| GEOLOCATION_LAT | Latitude of the order location |
| GEOLOCATION_LNG | Longitude of the order location |
| ORDER_STATUS | Status of the order (delivered or cancelled) |
| ORDER_PURCHASE_TIMESTAMP | The timestamp of the order purchase |
| ORDER_APPROVED_AT | The timestamp of the order approval |
| ORDER_DELIVERED_CARRIER_DATE | The timestamp of order delivery to the carrier |
| ORDER_DELIVERED_CUSTOMER_DATE | The timestamp of order delivery to the customer |
| ORDER_ESTIMATED_DELIVERY_DATE | The estimated delivery date to the customer |
| SELLER_TO_CARRIER_DELIVERYTIME | Time taken to reach from seller to carrier |
| CARRIER_TO_CUSTOMER_DELIVERYTIME | Time taken to reach from carrier to customer |
| SELLER_DISPATCH_TIME | Time taken by the seller to dispatch (binary: delay or fast) |
| CUSTOMER_DELIVERY_TIME | Time taken for the delivery (binary: delay or fast) |

# Business Hypothesis & Analysis

1. The e-commerce company started receiving complaints from the city 'santa rita'. The management wanted to have a look at the details of the orders from 'santa rita' and understand if there are any issues after looking at the data.

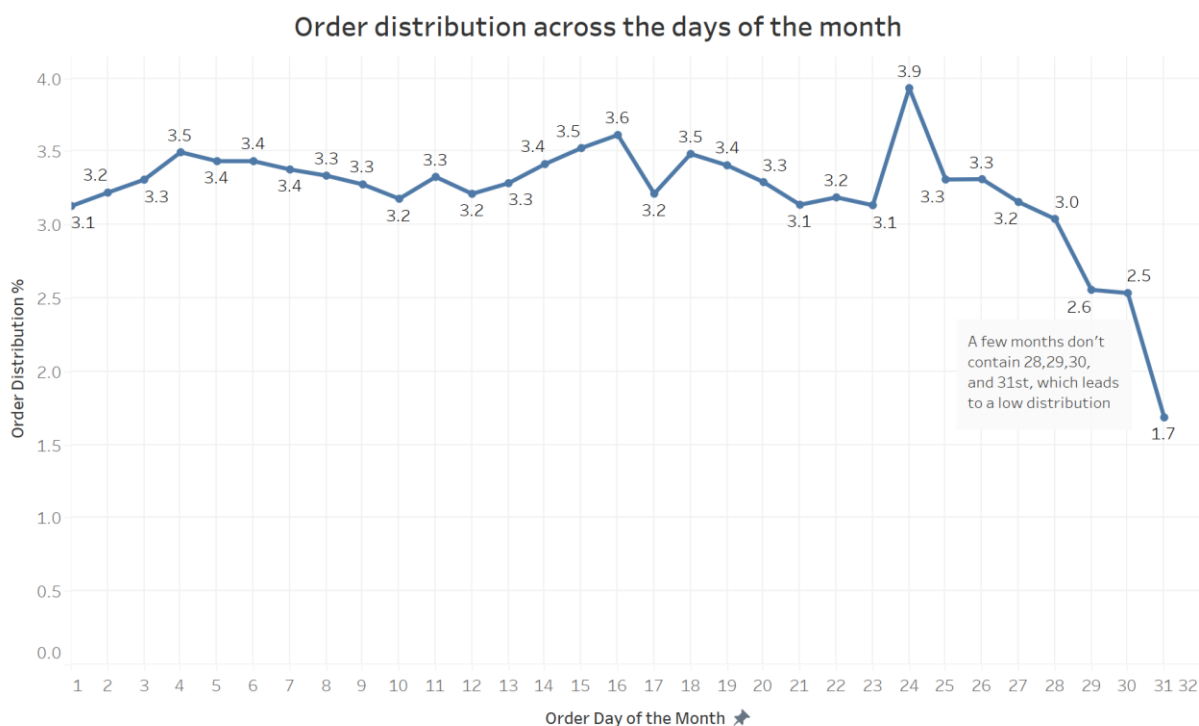| ER_DELIVERYTIME | CARRIER_TO_CUSTOMER_DELIVERY_TIME | SELLER_DISPATCH_TIME | CUSTOMER_DELIVERY_TIME | CUSTOMER_ID | CUSTOMER_CITY | CUSTOMER |
|---|---|---|---|---|---|---|
| 12.974 | | Fast | Fast | 3dcd5a57a32... | santa rita | MA |
| 17.127 | | Fast | Fast | d28e2706649... | santa rita | PB |
| 12.051 | | Fast | Fast | 4586e941d9fd... | santa rita | PB |
| 12.858 | | Fast | Fast | 67cacf586b43... | santa rita | PB |
| 9.982 | | Fast | Fast | ba1bedb9aa9... | santa rita | PB |
| 10.655 | | Fast | Fast | ff365cf639c0f... | santa rita | PB |
| 14.816 | | Fast | Fast | 54044f49e3d5... | santa rita | PB |

**There were a 7 orders from santa rita, the details of which were presented to the management.**

2. What is the average order value and average delivery cost?

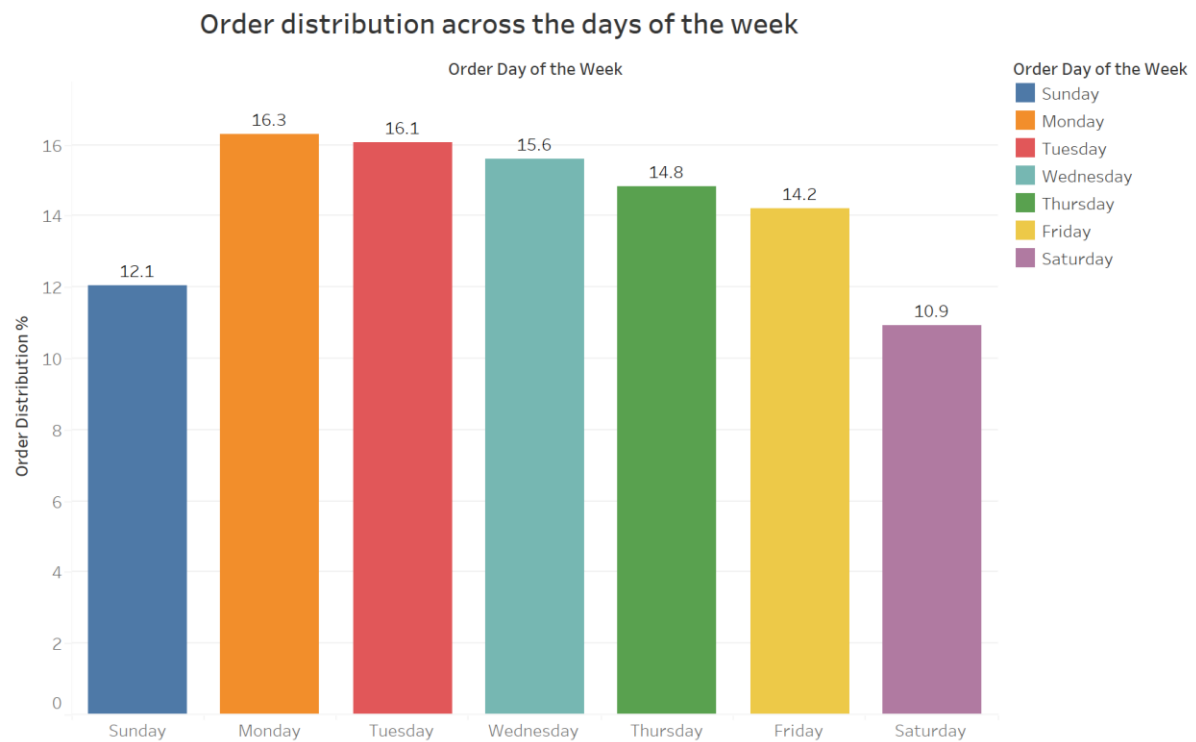| Avg. order value | Avg. delivery price |
|---|---|
| **R$ 125.3** | **R$ 20.6** |

**Avg. delivery price is approximately one-sixth of the avg. order value**

3. How are the orders distributed over the days of the month? and days of the week?

## Order distribution across the days of the month



The trend of sum of Order Distribution for Order Day Month. The marks are labelled by sum of Order Distribution.

**The number of orders peak on 16th and 24th dates of the month, while 1st, 21st, 23rd and 28th record low orders.**
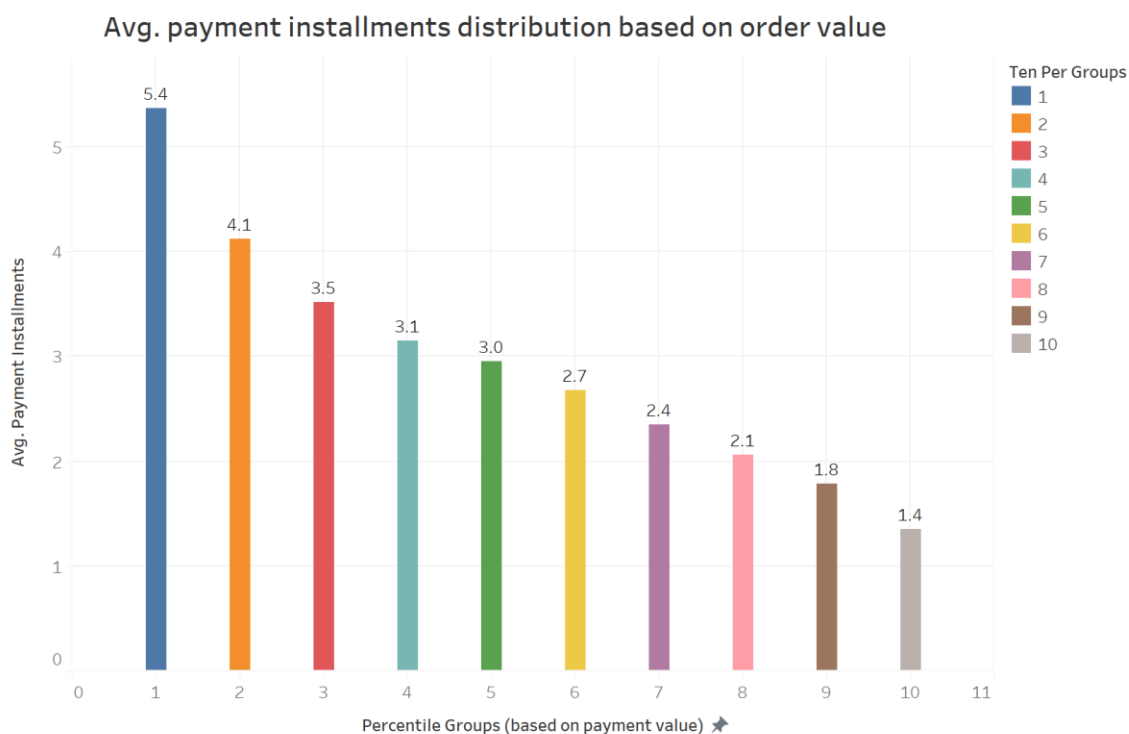
Order distribution across the days of the week



Sum of Order Distribution for each Order Day of the Week. Colour shows details about Order Day of the Week. The marks are labelled by sum of Order Distribution.

**Mondays and Tuesdays have higher % of orders, while Saturdays and Sundays have a significantly lower % of orders**

4. What are the highest and lowest ordered product categories?

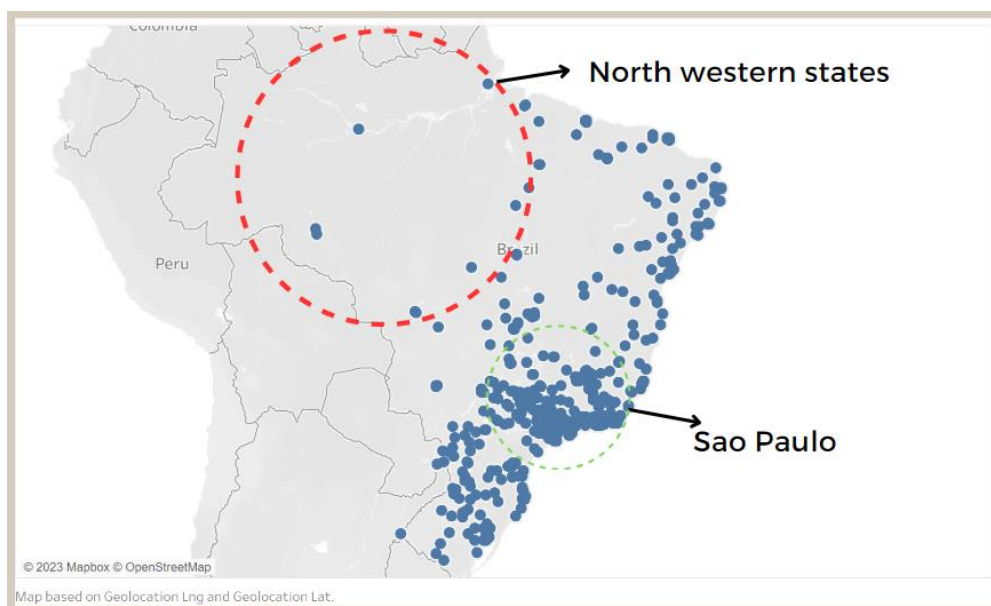| Product Categories | |
|---|---|
| **Highest order** | **Lowest ordered** |
| Bed, Table and Bath products | Kitchen |
| Beauty and Health care | Gaming PCs |
| Sport | Children clothing |
| Computer accessories | Insurance and Services |

5. Do customers prefer more installments for higher-value products?

## Avg. payment installments distribution based on order value



The plot of average of Avg Payment Installments for Ten Per Groups. Colour shows details about Ten Per Groups. The marks are labelled by sum of Avg Payment Installments.

**The average payment installments decreases as the order value increases**

6. The management is interested in knowing the states with the highest and the lowest orders as it would help allocate the workforce accordingly. The distribution of orders across different states in the country. A detailed heatmap would help the organization allocate the workforce more efficiently.



Map based on Geolocation Lng and Geolocation Lat.

**Sao Paulo state has the highest frequency of orders, while most of the north western regions have no orders**

7. What is the percentage of orders that are dispatched late by the seller (seller_dispatch_time = Delay) but delivered before the promised time to the customer (customer_delivery_time = Fast)?

**6.8% of the orders are dispatched late by the sellers but received on time by the customers**

8. As high-value and low-volume products generate a lot of revenue for an e-commerce company, What is the distribution of (high-value, low-volume), (low-value, high-volume), (low-value, low-volume), (high-value, high-volume) products are delivered before the promised time to the customer?

| Value and Volume bands | # of Products | % Distribution |
|---|---|---|
| High value and High volume | 30307 | 43.4 |
| Low value and High volume | 17243 | 16.6 |
| Low value and Low volume | 30307 | 23.4 |
| High value and Low volume | 17243 | 16.6 |

**The distribution of High value and Low volume is quite less as compared to the counterparts**

9. How many orders in total are using more than 5 vouchers to pay the order amount?

**124**

**The number of orders with more than 5 vouchers are quite less**

10. As it would decrease the delivery costs, the management wants to run marketing campaigns and generate more orders from the states where we have sellers but very few customers.

Seller/Customer ratio for state X = $\dfrac{\text{\# of orders delivered by the sellers from state X}}{\text{\# of orders delivered to customers in state X}}$

**Para**   **Sao paulo**

11. The management would like to run a credit card campaign if there are a significant number of customers who are not using credit cards to pay the orders, what is the percentage of customers who are using a credit card for a full payment/ partial payment?

**77%**

**The % of orders using credit cards for partial/complete payment are quite high**