**ANALYTIC TOOLS AND DECISION MAKING**

**PROJECT SUBMISSION**

**OLIST STORE INSIGHTS**



**Group Members**

**Manoj Kumar Ravindranath**

**Sai Sowmith Komakula**

**Jerin Thomas**

**Elda Maria Thomas**

**Siva Krishna Unnam**

**Srilekhya Byreddi**

# CONTENT

**1) PROPOSAL**

Introduction

Dataset Description

Problem Statement

Background

Insight

Motivation

Target Audience

Proposal

Approach

**2) EXPLORATORY DATA ANALYSIS**

Dataset Description

Schema

Data Segmentation

Dataset Collection

Data Imbalance

Inclusion Criteria

Preliminary Visualizations

Correlations

**3) DATA TRANSFORMATION**

**4) DATA CLEANING USING PYTHON**

**5) ANALYSIS AND VISUALIZATION USING TABLEAU**

**6) CONCLUSION**

**7) REFERENCES**

# PROJECT PROPOSAL

**Introduction:**

Olist is a Brazilian Startup that operates in the e-commerce segment, mainly through the marketplace. Olist connects small businesses to larger product marketplaces to help entrepreneurs sell their products to a larger customer base. On the one hand, the olist concentrates sellers who want to advertise on marketplaces such as Mercado Livre , B2W , Via Varejo and Amazon . On the other hand, it concentrates the products of all sellers in a single store that is visible to the final consumer. The company is headquartered in Curitiba, Paraná, and has an office in São Paulo. Currently, the business has 300 employees and more than 9,000 shopkeepers, in addition to 2 million unique consumers.

**Dataset Description:**

Brazilian E-Commerce Public Dataset by Olist was obtained from Kaggle that consist of 100k orders from 2016 to 2018 made at multiple marketplaces in Brazil. This also includes multiple Csv files based on region, customer information, seller information.
https://www.kaggle.com/olistbr/brazilian-ecommerce?select=olist_customers_dataset.csv

**Problem Statement:**

The Dataset contains the order and customer information of 98K orders and reviews of those orders. Out of these orders, around 22000 orders have least review scores of 1 and 2 and this will be the problem to be solved. 22k is almost 21 percent of the orders received in the given time period and this could have a toll on the reputation this E-commerce business. This data also has information of orders that had got good reviews and with this support what worked well among customers can be found.

These are some of the questions to be analyzed:

What is the average review score across states based on delivery delay?
Does the payment value influence review score?
What is the review score based on Price and freight value?
Does the product from a particular seller state affect the review score?
Does the price range of a particular product have an impact on customer behavior?

**Background:**

Online shopping refers to buying goods and services from merchants on the Internet. Online consumers are equally distributed between men and women and are more educated, younger, and affluent than those living in general. E-Commerce or Electronic commerce is a process of buying, selling, transferring, or exchanging products, services, and/or information via electronic networks and computers. It uses the internet and other online services to be engaged in buying and selling of digital and non digital products and services which require digital transportation or physical transportation.

The three technologies used in E-commerce are: Electronic Markets, Electronic Data Interchange Internet Commerce. E-commerce market is growing at a high rate. The online market is expected to grow by 56% in the next two years. retail e-commerce sales worldwide amounted to 2.3 trillion US dollars and e-retail revenues are projected to grow to 4.891 trillion US dollars while the growth of traditional market is gradually declining. Many larger retailers can maintain a presence offline and online by linking physical and online offerings to stay in the competition. There are different strategies for conducting business in online and traditional markets.

**Insight:**

In the contemporary world number of customers for any business is proportional to the number of good reviews provided. This data has good proportion of good and bad review score its well balanced.

It is very apparent that the E-commerce will provide next generation shopping experience and customers are leaning towards this internet experience and studies show that the customers are content with online shopping and most likely to prefer only this in future at least with regards to some products. The more people prefer online shopping the more the expectations, Hence E-commerce business must be on point when it comes to providing a smooth shopping experience.

**Motivation:**

As we all know data is one of the biggest drivers of making an analysis and predictions. With this amount of data, and the kind of data, it is feasible to deliver valuable insights just by utilizing the open resources.

Earlier it was an expensive task to analyze a data because that mean using a lot of software resources and manpower. Today with the growth in software and other valuable resources. With regards to this data, all the necessary information that can be used to perform different analysis and with the time and date data it is also possible to predict the trend.

**Target Audience:**

This analysis will be focused on the needs of customers and small business owners. Since the problem to solved is based on review score, customers will be someone getting benefited. Hence, they are the target audience.

**Proposal:**

As a part of this project, a descriptive analysis will be conducted of the dataset and correlation will be found between the Review scores and other variables. The past trend will be thoroughly examined and will be used to understand the customer behavior. In the end these insights could be used improve the business performance.

**Approach:**

- ➢ Finding the relationships among the files will be the first step towards this project.
- ➢ Each and every file will be investigated.
- ➢ All the missing values, outliers and any other discrepancy will be addressed.
- ➢ Merging tables will be performed based on the needs.
- ➢ Excel and Python will be used for plotting the graphs.
- ➢ Tableau will be used to create dashboards and relationships.

## EXPLORATORY DATA ANALYSIS

**Dataset Description:**

This is a Brazilian ecommerce public dataset of orders made at Olist Store. The dataset has information of 100k orders from 2016 to 2018 made at multiple marketplaces in Brazil. Its features allow viewing an order from multiple dimensions: from order status, price, payment and freight performance to customer location, product attributes and finally reviews scores given by customers.

This dataset has 9 csv files.

**olist_customers_dataset.csv**

This dataset has information about the customer and its location. Use it to identify unique customers in the orders dataset and to find the orders delivery location.

| Column Name | Data Type | Description |
|---|---|---|
| customer_id | Varchar | Each order has a unique customer_id |
| customer_unique_id | Varchar | Unique identifier of a customer |
| customer_zip_code_prefix | Categorical | First five digits of customer zip code |

| Customer_city | Categorical | Customer city name |
|---|---|---|
| customer_state | Categorical | Customer state name |

- There are 99441 unique customer ids.
- There are 4119 unique customer cities and most of the customers is from Sao Paulo.
- There are 27 states and most common state is SP.

**olist_geolocation_dataset.csv**

This dataset has information Brazilian zip codes and its lat/lng coordinates. Use it to plot maps and find distances between sellers and customers.

| Column Name | Data Type | Description |
|---|---|---|
| geolocation_zip_code_prefix | Categorical | First 5 digits of zip code |
| Geolocation_lat | Numerical | Latitude |
| geolocation_lng | Numerical | Longitude |
| geolocation_city | Categorical | City name |
| geolocation_state | Categorical | State |

**Olist_order_items_dataset.csv**

This dataset includes data about the items purchased within each order.

| Column Name | Data Type | Description |
|---|---|---|
| Order_id | Varchar | Order unique identifier |
| Order_item_id | Numerical | Number identifying number of items included in the same order. |
| product_id | Varchar | Product unique identifier |
| seller_id | Varchar | Seller unique identifier |
| Shipping_limit_date | Date | Shows the seller shipping limit date for handling the order over to the logistic partner. |
| price | Numerical | Item price |
| Freight_value | Numerical | Item freight value item |

**Olist_order_payments_dataset.csv**

This dataset includes data about the orders payment options.

| Column Name | Data Type | Description |
|---|---|---|
| order_id | Varchar | Unique identifier of an order. |
| payment_sequential | Numerical | Customer may pay an order with more than one payment method. If he does so, a sequence will be created to accommodate all payments. |
| payment_type | Categorical | Method of payment chosen by the customer. |
| payment_installments | Numerical | Number of installments chosen by the customer. |
| payment_value | Numercial | Transaction value. |

- Mean and standard deviation of payment_value is 154 and 217 respectively.
- There are some outliers found in payment_installments and Payment_value.

### Olist_order_reviews_dataset.csv

This dataset includes data about the reviews made by the customers.

| Column Name | Data Type | Description |
|---|---|---|
| review_id | Varchar | Unique review identifier |
| order_id | Varchar | Unique order identifier |
| review_score | Categorical | Note ranging from 1 to 5 given by the customer on a satisfaction survey. |
| review_creation_date | Date | Shows the date in which the satisfaction survey was sent to the customer. |
| review_answer_timestamp | Date | Shows satisfaction survey answer timestamp. |

### Olist_orders_dataset.csv

This is the core dataset. From each order you might find all other information.

| Column Name | Data Type | Description |
|---|---|---|
| order_id | Varchar | Unique identifier of the order. |
| customer_id | Varchar | Each order has a unique customer_id. |
| order_status | Categorical | Reference to the order status (delivered, shipped, etc.) |
| order_purchase_timestamp | Date | Shows the purchase timestamp. |
| order_approved_at | Date | Shows the payment approval timestamp. |
| order_delivered_carrier_date | Date | Shows the order posting timestamp. When it was handled to the logistic partner. |
| order_delivered_customer_date | Date | Shows the actual order delivery date to the customer. |
| order_estimated_delivery_date | Date | Shows the estimated delivery date that was informed to customer at the purchase moment. |

- There are missing values found in order_approved_at (160 values), order_delivered_carrier_date (1783 values) and order_delivered_customer_date (2965 values).
- No outliers found in this dataset.

### Olist_products_dataset.csv

This dataset includes data about the products sold by Olist.

| Column Name | Data Type | Description |
|---|---|---|
| product_id | Varchar | Unique product identifier |
| product_category_name | Categorical | Root category of product, in Portuguese. |
| product_name_length | Numerical | Number of characters extracted from the product name. |
| product_description_length | Numerical | Number of characters extracted from the product description. |
| product_photos_qty | Numerical | Number of product published photos |
| product_weight_g | Numerical | Product weight measured in grams. |
| product_length_cm | Numerical | Product length measured in centimeters. |
| product_height_cm | Numerical | Product height measured in centimeters. |

| product_width_cm | Numerical | Product width measured in centimeters |
|---|---|---|

- Except for product_id, 610 missing values are found in rest of each columns.
- There are outliers found in product dimensions.

**Olist_seller_dataset.csv:**

This dataset includes data about the sellers that fulfilled orders made at Olist. Use it to find the seller location and to identify which seller fulfilled each product.
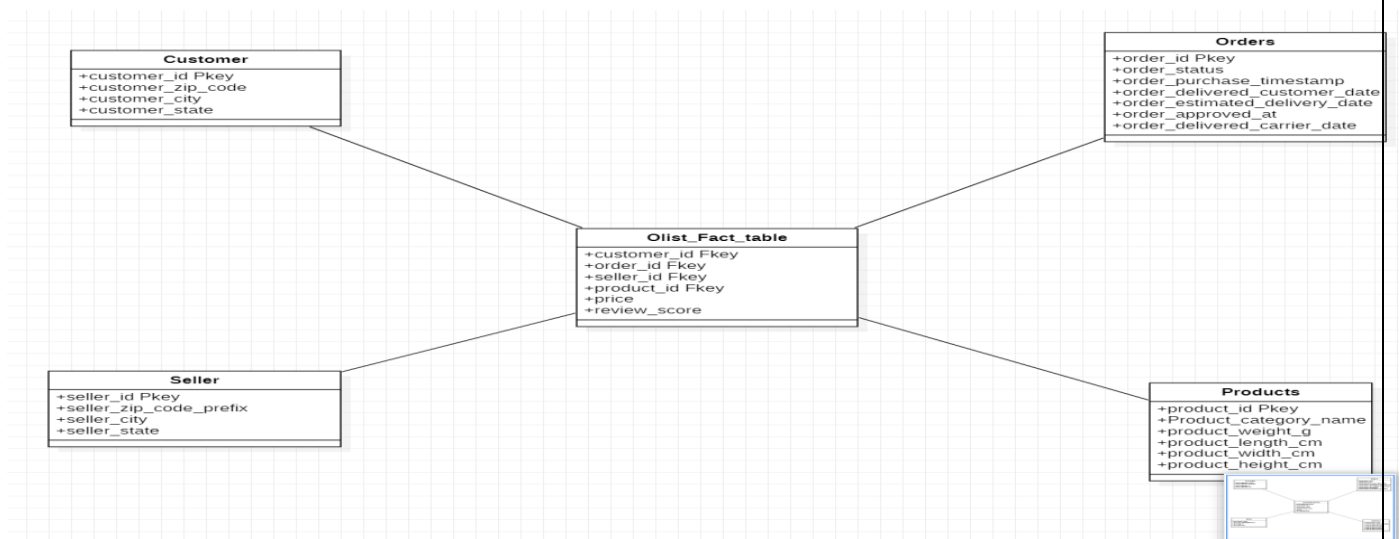
| Column Name | Data Type | Description |
|---|---|---|
| seller_id | Varchar | Seller unique identifier |
| seller_zip_code_prefix | Categorical | First 5 digits of seller zip code |
| seller_city | Categorical | Seller city name |
| seller_state | Categorical | Seller state |

**product_category_name_translation.csv**

Translates the product *category* name to English.

| Column Name | Data Type | Description |
|---|---|---|
| product_category_name | Categorical | Category name in Portuguese |
| product_category_name_english | Categorical | Category name in English |

**Schema:** The original dataset is modified by replacing the columns using the power query. Below is the star schema of the dataset. We have most of the ids in fact table as foreign keys and dimension table contains data of customers, orders, products, and sellers.



**Data Segmentation:**

Since the dataset has large verities of information it is necessary to subset the data into groups of product, customers, and sellers. Each segment holds appropriate information with respect to the data of various subsets.

**Data Imbalance:**

Review score column in our fact table looks imbalanced with 57% of the score being 5 and remaining 43% attributes to other score, but this shouldn't be a problem as the data aligns with real world data.

Mostly customers preferred payment method is credit cards compare to other methods. More than 75% customer orders from olist are paid through credit cards.
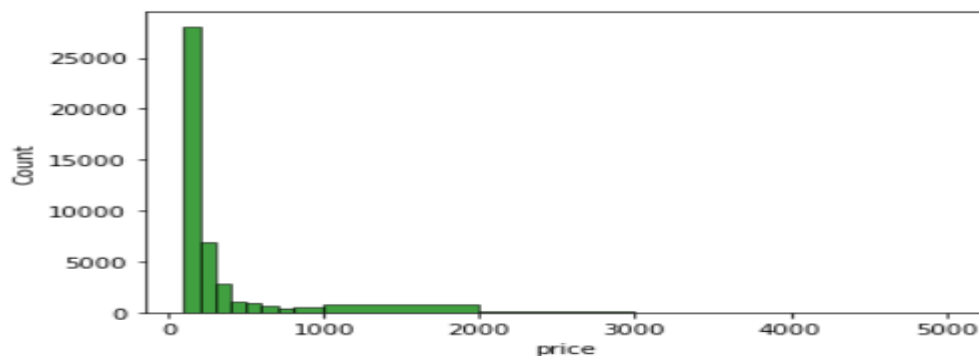
**Inclusion Criteria:**

Few custom columns need to be included that calculates number of days between estimated_delivery_date column and order_delivered_customer_date. Key features that are to be used in this analysis are order price, payment value, freight value, customer state, product category and payment installments.

**Preliminary Visualizations:**

**Univariate Analysis**

**1. Price distribution of products. (Price of product vs Count)**



On X-axis we have price of each product and on Y-axis we have counts of products which have same price.Price of orders anything above 1000 is considered as an outlier.

**2. Distribution of freight value. (Freight value of product vs count)**



The graph shows the distribution of fright value among orders. if an order has more than one item the freight value is splitted between items. The data from the graph is right skewed.

**3. Distribution of payment value (Payment Value of product vs count)**

This graph shows distribution of payment value. Almost most of the payment is under 500 and this distribution clearly has some outliers. So, anything above 1000 could be considered as an outlier.
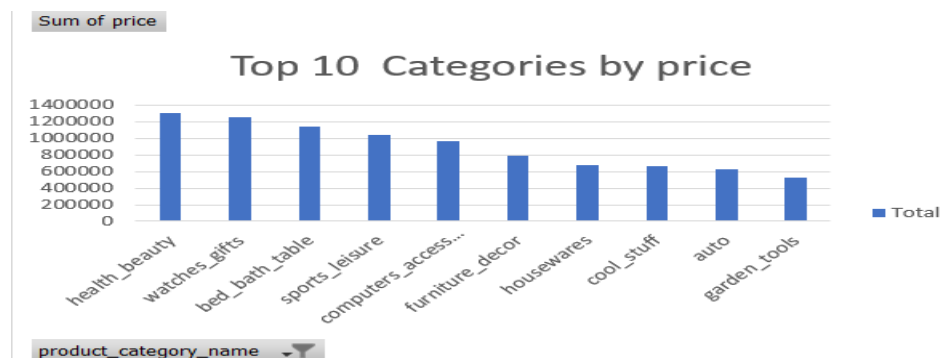
**4. Top 10 best selling product categories (Category of products vs count)**

Count of product_category_name

Top 10 Best selling Categories



Among all the categories these products are in top 10 mostly purchased by customers. We can observe the products of bed_bath_table is on the top 1 which are sold in more than double of the auto products which are in top 10.

**Bivariate Analysis**

**1. Top 10 categories by sum of their price (category of product vs price of product)**

Sum of price

Top 10 Categories by price



The total price of Health_beauty and watches_gifts purchased by customers in olist store is over 2.3M which generates highest revenue for the store.
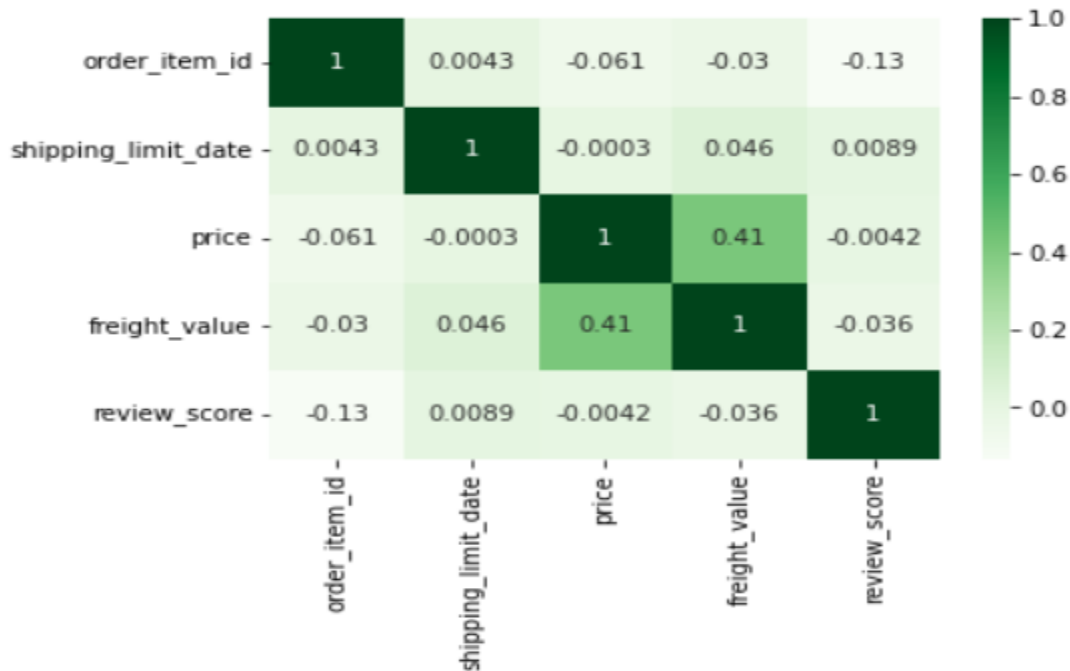
**Correlations:**

**1. Correlation between price and review score**

Most of the customers who had given their review score are the people who order value is less than 3000.Anything above 3000 maybe considered as outliers.

**2. Correlation matrix of fact table**

Correlation values ranges between -1 to +1 where -1 and +1 represents the strong correlation.



The above matrix shows correlation between key features that will be used in our analysis. It is clear from the matrix few variables are positively correlated and there is a negative correlation between few variables. The variable price and freight value has the second-best correlation.

# Data Transformation

**Merging Tables**: Data Merging was performed in Excel, Fact table was taken as a primary table to initiate merging. Second table is **olist_customers_dataset.csv**

**Step 1:**

First set of tables are fact table and olist_customers_dataset.csv.

Merging of the Fact table and the customer table was performed based on "customer id" column.

**Step 2:** The newly merged table is having the fact table information along with the customer details.
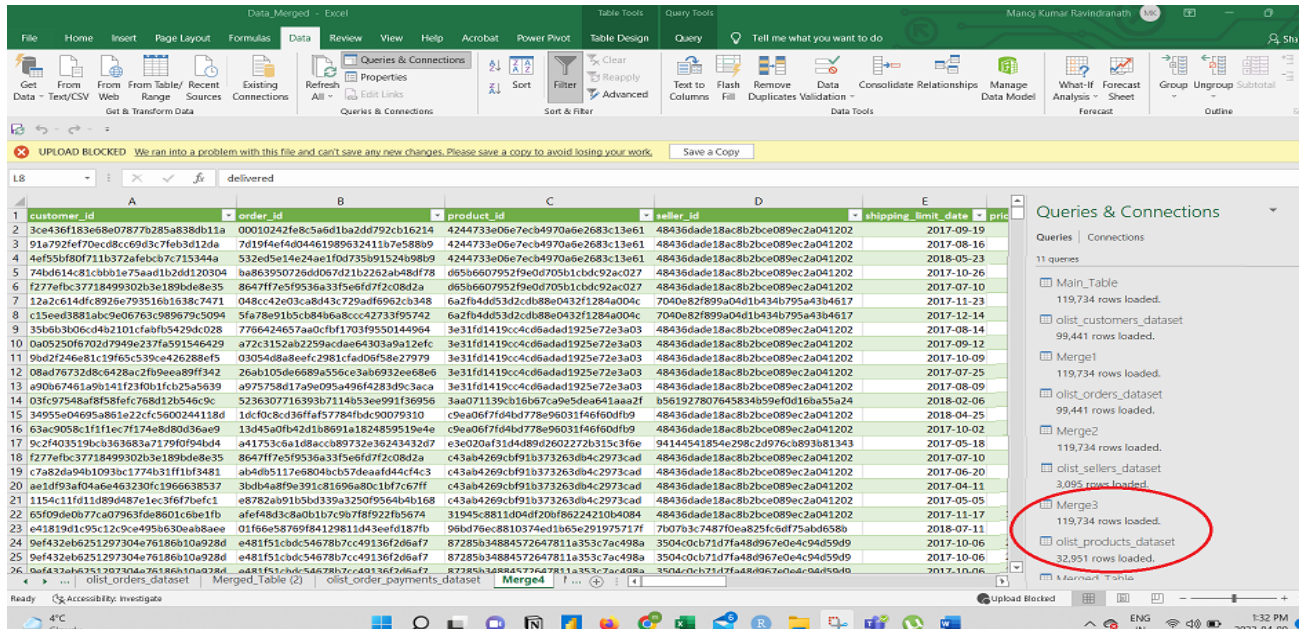
This table is now merged with the **Olist_orders_dataset.csv** based on order id column.
**Step 3:**

Naming the newly merged table as Merge2. This table is further merged with Olist_seller_dataset.csv based on "Seller id" Column.

**Step 4:**

The new table "Merge 3" is merged with Olist_products_dataset.csv based on the "product id" column.

**Step 5:** Finally, to get the product name translation from the product_category_name_translation.csv the final table is merged based on product name column.

Now the Merged tables has all the columns from the fact table and the other tables.

Additionally, to get the payment information from the Payments table the existing merged table is further merged with the Payment table.

The final table has all the information available in the olist_customers_dataset.csv, Olist_order_items_dataset.csv, Olist_order_payments_dataset.csv, Olist_orders_dataset.csv, Olist_products_dataset.csv,Olist_seller_dataset.csv, product_category_name_translation.csv.

# **Data Cleaning Using Python**

### **Step 1: Loading necessary packages**

Pandas and numpy is imported into python. A new python file was created in the same location where the csv file is located.

### **Step 2: Importing the csv file**

Import dataset using pd.read_csv and save it as df1.Check all the columns have the appropriate data type.

### **Step 3: Checking the null values**

Check the null values in the dataset. Since the null values account to less than 5% of the data set it doesn't affect the final analysis, Hence, dropping all the null values.

The following columns has null values:
- review_score
- olist_orders_dataset.order_approved_at
- olist_orders_dataset.order_delivered_carrier_date
- olist_orders_dataset.order_delivered_customer_date
- olist_orders_dataset.order_estimated_delivery_date
- product_category_name
- product_name_lenght

**Step 4: Checking duplicate observation**

All the duplicate observation were found and removed.

**Step 5: Dropping columns**

Columns like product_length_cm', 'product_height_cm', 'product_width_cm', product_name_lenght were dropped because they have the least correlation with our target variable.

**Step 6: Renaming the columns**

olist_orders_dataset.order_status": "order_status",
olist_orders_dataset.order_purchase_timestamp": "order_purchase_Time",
olist_orders_dataset.order_approved_at":"order_approved_at",
olist_orders_dataset.order_delivered_carrier_date":"order_delivered_carrier_date"
olist_orders_dataset.order_delivered_customer_date":"order_delivered_customer_date",
olist_orders_dataset.order_estimated_delivery_date":"order_estimated_delivery_date"}
**Step 7: Creating custom columns**

Based on the analysis questions few custom columns like "**delivery_delay", "Days_to_Deliver", "delayed"** have been created.
Where "delivery delay" has the information on difference between estimated delivery and the actual delivery.
Days to deliver: Difference between delivered date and order approved date.
Delayed: This column has information on delayed orders. This is a categorical column that has only two values called "Delayed" and "Not delayed"
**Step 8: Outlier Treatment**

Since we have only few and valid outliers in the numerical columns like price, freight_value, Payment_Value It doesn't affect the output of the Analysis. Hence it doesn't require any transformation

**Step 9: Exporting the into CSV file**

The cleaned data is exported as a csv file into the working location.
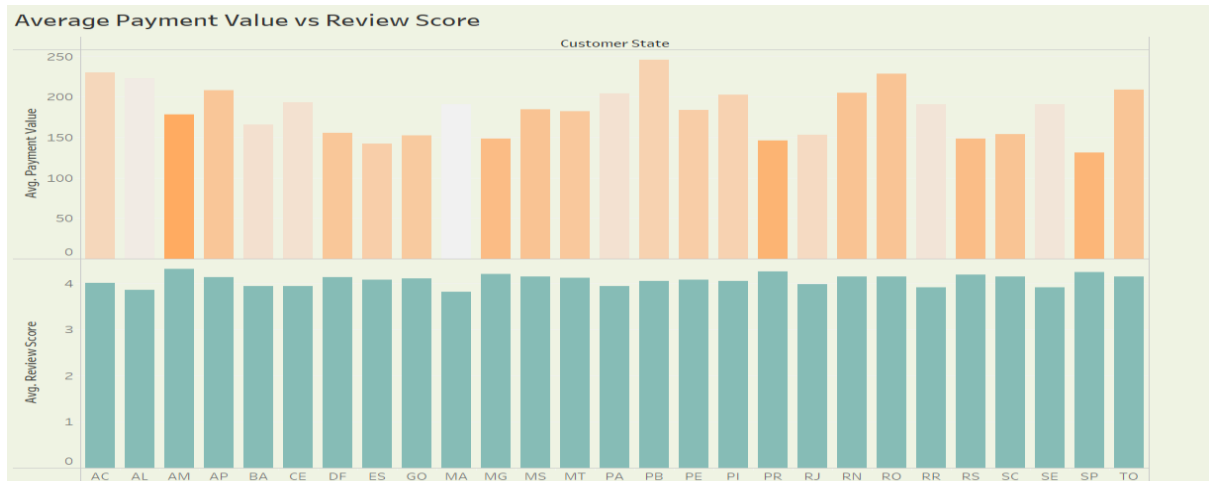
# Visualizations Using Tableau

Answering the questions form problem statement using Tableau visualizations.

**What is the average review score based on delivery delay?**

The above graph was plotted to find if the delay in order influences the review score. The average review score for each state has been plotted based on whether the order was delivered on time or not. It is evident from the graph that the review score has been affected by the delivery delays. The orange portion of the graph which shows the order delivered on time has an average review score of more than 5 whereas the blue portion that represents the delayed order has an average review score of 2. this clearly shows the customer is not happy with the delayed orders.
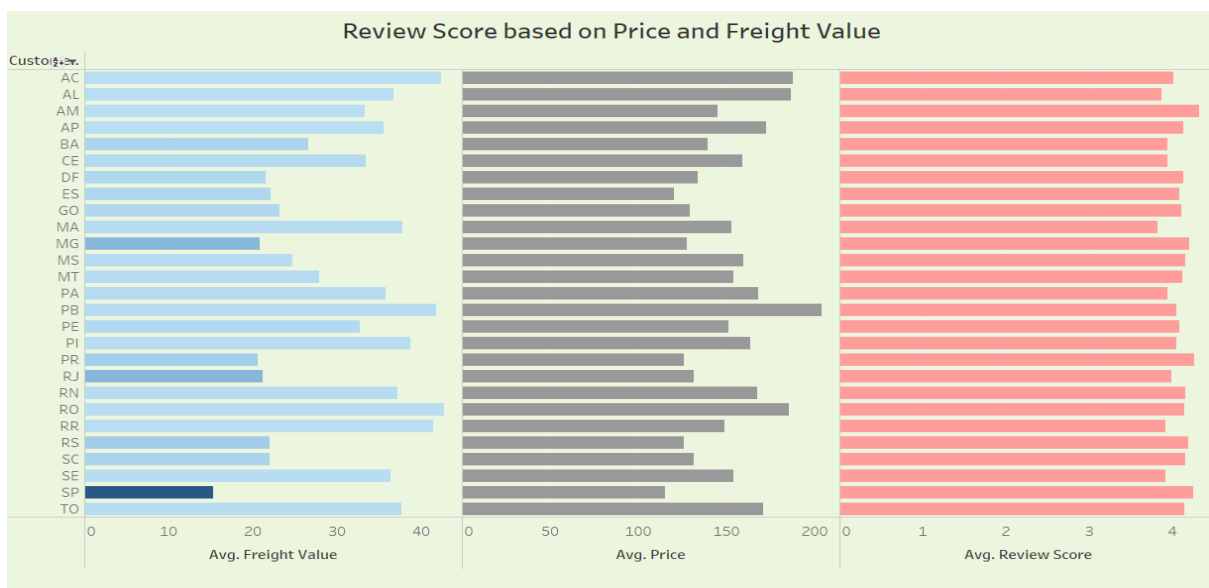
**Does the payment value affect the review score?**



The above graph was made with the intention to know if the payment value has any effect on the review score. The average payment value vs the review score has been plotted on province level.

While almost all provinces have average review score over 4, states like AL, BA, MA have review score less than 4. By looking at the graph the payment value has no significant effect because there is no big correlation between these two variables. There are different averages for the same payment range.
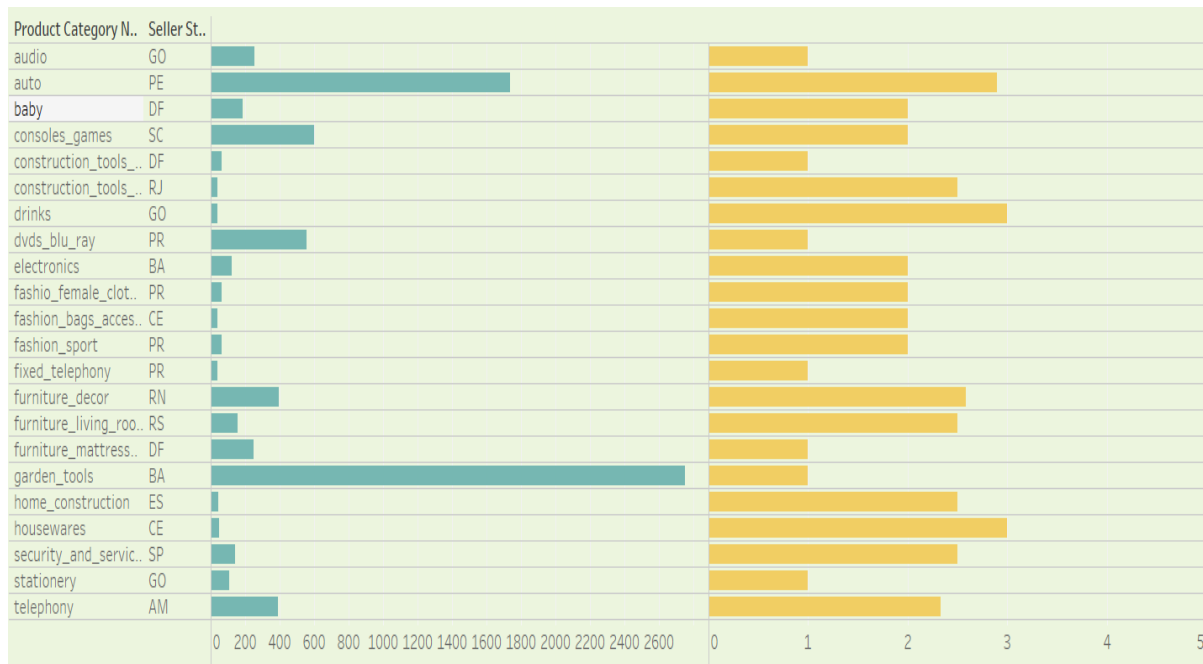
**What is the review score based on Price and freight value?**



Breaking down the payment value in the previous graph into freight value and average price and plotting against review score. When trying to find the relationship between freight value and the review score it
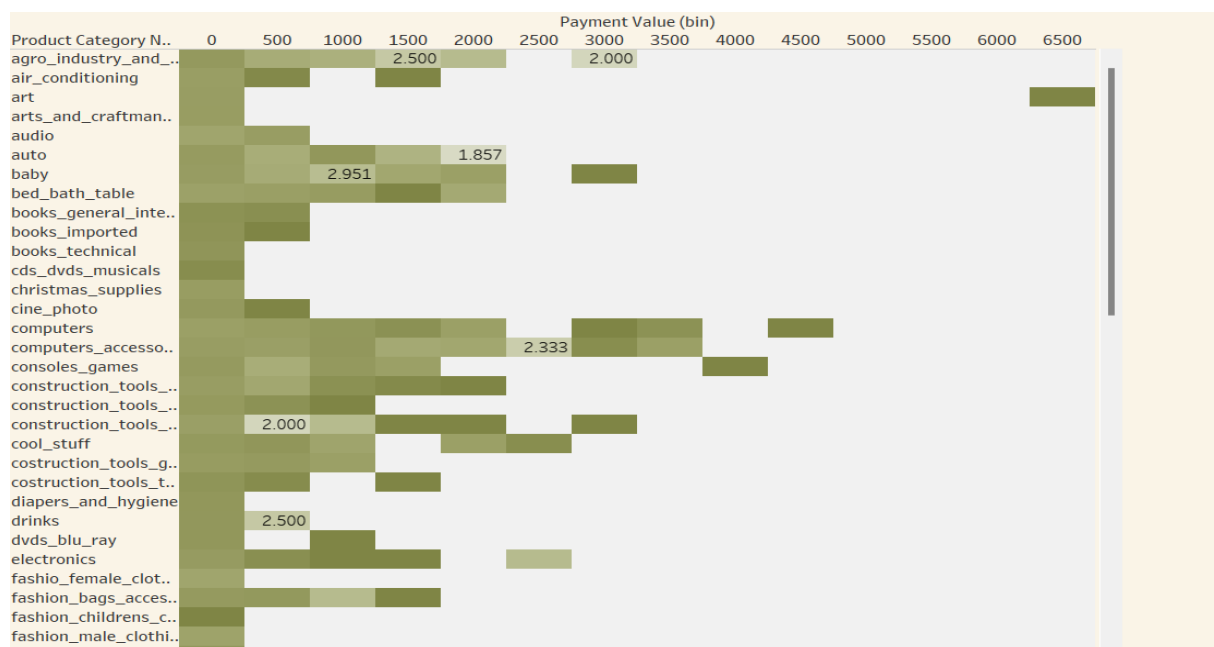
was found that the state with the least freight value has the highest average. And the review score is dropped below 4 in the states where the freight value is relatively higher. The "price" variable doesn't seem to have much influence on the review score.

**Does the product from a particular seller state affect the review score?**



A visual for product category from a particular seller state is plotted against the average review score. Out of all product category and seller state combination around 20 combinations have review score less than 3. It is found that certain products from a state didn't meet the minimum expectations of the customers. And one more interesting find is the state "PR" and "GO" have been mentioned many times in this graph. Hence products from these states needs a check.

**Does the price range of a particular product have an impact on customer behavior?**

The above graph shows the review score based on the product category and the amount spend by the customers. This graph gives a specific insight on the customer behaviour based on the amount they spend on a particular product. The darker the shade the higher the review. The light shades represent low review scores. Customers that spent around 1500 and 3000 on Argo industry and commerce product have given very low review scores. Similarly, Auto products that cost around 2000 didn't work well among customers. And Computer accessories that cost around 2500, Construction tools around 500, furniture décor around 1500 have got poor reviews.

**Conclusion:**

Customers are clearly not happy with the delayed orders, And the orders that were delivered on time is rewarded with high reviews.

Payment value has no direct influence in the review scores

The freight value which is included inside the payment value has an impact on the review scores. If the logistic partner has a way reduce this value, there can be a good change in the customer behavior.

There are states like "GO", "PR" that are selling products that are disliked by customers, a special attention has to be given to the product from these states.

A product in general has high review scores but drilling down with respect to the prices, products like agro, auto, baby, computer accessories have very poor review score , Hence this must be taken care of.

**References:**

Plotting of correlation matrix

How to Create a Correlation Matrix using Pandas - Data to Fish

Link to Dataset

https://www.kaggle.com/olistbr/brazilian-ecommerce?select=olist_customers_dataset.csv

Finding Relationships in Data

Python Guide: Finding Data Relationships with a Correlation Matrix | Pluralsight

Link to GitHub repository

https://github.com/Hackslash-DA/DB103-Data_Analysis_Project

Tableau Public Dashboard link

https://public.tableau.com/app/profile/elda3801/viz/103Final/Dashboard1