

Visualization Assignment 1

Student Name: Ngoc Chien Le

Student ID: 1057045

Tutor Name: Bruno Luis Mendivez Vasquez (Lab Monday 12pm)

Activity Number: 26

Contents

1.	Introduction.....	2
2.	Data Wrangling.....	2
3.	Data checking in detail.....	5
4.	Data Exploration	6
5.	Conclusion	8
6.	Reflection.....	8

1. Introduction

a. Problem description:

Taking the view as a manager of an e-commerce platform, this analysis will focus on the questions aiming to grasp the overview of the business as well as understand more about the sellers who are joining the platform, thanks to that we could find the way to improve seller performance, boost their cohesion to the e-commerce platform.

b. Question

Who is making money on our e-commerce platform?

- Who are those sellers? What do they sell? Where are they?
- Where are their customers?
- What are the marketing channels the company used to acquire them?

c. Motivation

From those the Info and analysis, we could:

- See the bigger picture from the process of seller recruitment into tracking sellers' performance.
- Up sales those sellers by providing them more services (sell strategy consultancy, advertisements.).
- Find more similar sellers (based on location, industry - could be similar or relatively alike, or marketing channels).

2. Data Wrangling

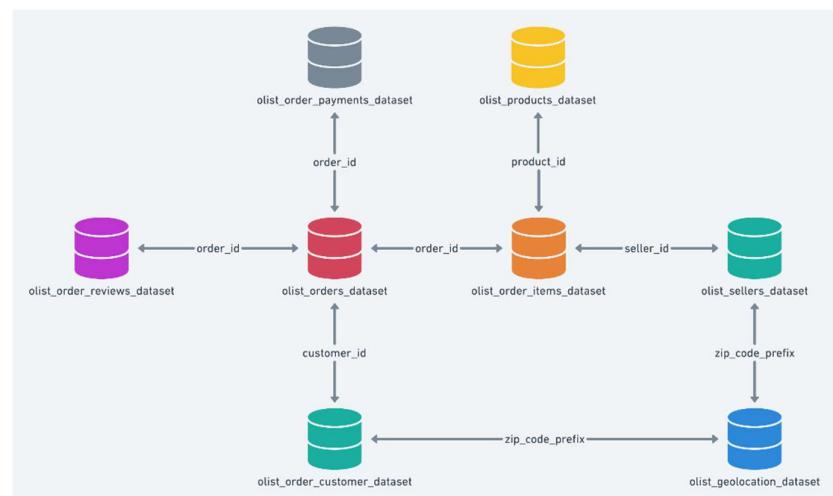
a. Data Overview

The data is from Olist, a Brazilian department store that operates in e-commerce segment, separated into multiple files therefore we need to merge them together based on database structure provided.

Data of 100k orders from 2016 to 2018

Link to the data source:

- [Olist - e-commerce](#)
- [Olist - Marketing funnel](#)

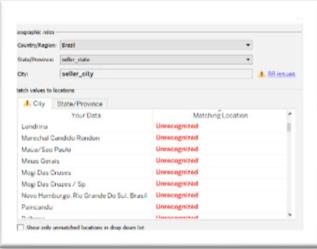


b. Data wrangling in detail

Problems	Actions	Methods	Software	Outcome																								
The data is separated as can be seen in data structure image above, included many columns that are redundant to answer our question	Merge datafiles and drop the unnecessary columns	<p>Based on analysing the question, we could filter the features (necessary vs unnecessary)</p> <p>Merge the files by look at the database structure and the key that used to connect files together</p>	Tool: Python Package: Pandas	Merged_ecommerce_data file																								
<u>Order Item Dataset</u> Instead of using one record as total number of a product, it has been separated into multiple records as an example below which led to many duplications with the data <div style="border: 1px solid black; padding: 5px; margin-top: 10px;">  <pre> 9924b6fc7c4eb0b0094318b105f 1 43423cff1de7fda0d90414ed38c11a73 b1f64f64d59aeb8b6911ab38803c57a9 9924b6fc7c4eb0b0094318b105f 2 43423cff1de7fda0d90414ed38c11a73 b1f64f64d59aeb8b6911ab38803c57a9 9924b6fc7c4eb0b0094318b105f 3 43423cff1de7fda03d0414ed38c11a73 b1f64f64d59aeb8b6911ab38803c57a9 </pre> </div> <p>order_item_id = sequential number - identifying number of items included in the same order (link to the explanation)</p>	Sort the data and remove the unnecessary rows, only keep the last row represents the total amount	Run the loop to check if product ID repeat in the 2 consecutive rows and then remove the previous one since the last one is last number (total amount) Note: after sorting and removing the duplication, saving to a new file to avoid rerunning In order to reflect the true meaning of the order_item_id in the modified data, it is going to be renamed into amount	Tool: Python Run a for loop	A new “order item Dataset” with no duplication <div style="border: 1px solid black; padding: 10px; margin-top: 10px;"> <pre> In [9]: #read the new modified dataset and check the duplication sorted_order_item = pd.read_csv('sorted_order_items.csv') sorted_order_item[sorted_order_item.duplicated(['order_id', 'product_id'])] Out[9]: index order_id order_item_id product_id seller_id shipping_limit_date price freight_value </pre> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th>index</th> <th>order_id</th> <th>order_item_id</th> <th>product_id</th> <th>seller_id</th> <th>shipping_limit_date</th> <th>price</th> <th>freight_value</th> </tr> </thead> <tbody> <tr> <td>0</td> <td>f30149f4a8882a08895b6a242aa0d612</td> <td>1</td> <td>00066f42aeeb9f3007548bb9d3f33c38</td> <td></td> <td></td> <td></td> <td></td> </tr> <tr> <td>1</td> <td>f5eda0ded77c1293b04c953138c8331d</td> <td>1</td> <td>00088930e925c41fd95ebfe695fd2655</td> <td></td> <td></td> <td></td> <td></td> </tr> </tbody> </table> </div>	index	order_id	order_item_id	product_id	seller_id	shipping_limit_date	price	freight_value	0	f30149f4a8882a08895b6a242aa0d612	1	00066f42aeeb9f3007548bb9d3f33c38					1	f5eda0ded77c1293b04c953138c8331d	1	00088930e925c41fd95ebfe695fd2655				
index	order_id	order_item_id	product_id	seller_id	shipping_limit_date	price	freight_value																					
0	f30149f4a8882a08895b6a242aa0d612	1	00066f42aeeb9f3007548bb9d3f33c38																									
1	f5eda0ded77c1293b04c953138c8331d	1	00088930e925c41fd95ebfe695fd2655																									

<p><u>Order Item Dataset</u></p> <p>Missing total value of invoice (or seller_revenue from the seller perspective) we need to do a calculation</p>	<p>Calculation</p>	<p>Calculate the invoice value by using the formula:</p> <p>Seller_revenue = amount * price</p>	<p>Tool: Python</p>	<p>Seller Revenue Column</p> <pre>#calculate revenue for each seller sorted_order_item["seller_revenue"] = sorted_order_item["amount"] * sorted_order_item["price"] sorted_order_item.head()</pre> <table border="1"> <thead> <tr> <th>product_id</th> <th>seller_id</th> <th>shipping_limit_date</th> <th>price</th> <th>freight_value</th> <th>seller_revenue</th> </tr> </thead> <tbody> <tr> <td>30066142aeab9f3007548bb9d313c38</td> <td>5670f4db5b62c43d542e1b2d5fb0c7c</td> <td>2018-05-24 18:58:59</td> <td>101.65</td> <td>18.59</td> <td>101.65</td> </tr> <tr> <td>30088930e925011d95ebfe6959d2655</td> <td>7142540dd4c91e2237ac07e91104eba2</td> <td>2017-12-18 19:32:19</td> <td>129.90</td> <td>13.93</td> <td>129.91</td> </tr> <tr> <td>30064066f7479715e4be8f1dd9112482</td> <td>4a3ca9315b744ce9fe9b9374361493884</td> <td>2017-12-29 16:12:38</td> <td>229.00</td> <td>13.10</td> <td>229.01</td> </tr> <tr> <td>100b8955cb9e00096488278317764d19</td> <td>40ec8ab6cdafbcc04f544d38c697da39a</td> <td>2018-08-16 13:35:21</td> <td>58.90</td> <td>19.60</td> <td>58.91</td> </tr> <tr> <td>10d9be29b5207b54e86aa1b1ac54872</td> <td>8ae520247981aa06bc94abddf5f46d34</td> <td>2018-04-09 10:09:40</td> <td>199.00</td> <td>19.27</td> <td>199.01</td> </tr> </tbody> </table>	product_id	seller_id	shipping_limit_date	price	freight_value	seller_revenue	30066142aeab9f3007548bb9d313c38	5670f4db5b62c43d542e1b2d5fb0c7c	2018-05-24 18:58:59	101.65	18.59	101.65	30088930e925011d95ebfe6959d2655	7142540dd4c91e2237ac07e91104eba2	2017-12-18 19:32:19	129.90	13.93	129.91	30064066f7479715e4be8f1dd9112482	4a3ca9315b744ce9fe9b9374361493884	2017-12-29 16:12:38	229.00	13.10	229.01	100b8955cb9e00096488278317764d19	40ec8ab6cdafbcc04f544d38c697da39a	2018-08-16 13:35:21	58.90	19.60	58.91	10d9be29b5207b54e86aa1b1ac54872	8ae520247981aa06bc94abddf5f46d34	2018-04-09 10:09:40	199.00	19.27	199.01
product_id	seller_id	shipping_limit_date	price	freight_value	seller_revenue																																			
30066142aeab9f3007548bb9d313c38	5670f4db5b62c43d542e1b2d5fb0c7c	2018-05-24 18:58:59	101.65	18.59	101.65																																			
30088930e925011d95ebfe6959d2655	7142540dd4c91e2237ac07e91104eba2	2017-12-18 19:32:19	129.90	13.93	129.91																																			
30064066f7479715e4be8f1dd9112482	4a3ca9315b744ce9fe9b9374361493884	2017-12-29 16:12:38	229.00	13.10	229.01																																			
100b8955cb9e00096488278317764d19	40ec8ab6cdafbcc04f544d38c697da39a	2018-08-16 13:35:21	58.90	19.60	58.91																																			
10d9be29b5207b54e86aa1b1ac54872	8ae520247981aa06bc94abddf5f46d34	2018-04-09 10:09:40	199.00	19.27	199.01																																			
<p><u>Merged marketing data</u></p> <p>Including missing values and a vague category(unknow)</p>	<p>Fill the missing values</p> <pre>merged_marketing_data.info() <class 'pandas.core.frame.DataFrame'> Int64Index: 842 entries, 0 to 841 Data columns (total 6 columns): # Column Non-Null Count Dtype --- 0 mql_id 842 non-null object 1 seller_id 842 non-null object 2 business_segment 841 non-null object 3 business_type 832 non-null object 4 landing_page_id 842 non-null object 5 origin 828 non-null object dtypes: object(6) memory usage: 46.0+ KB merged_marketing_data.origin.value_counts() organic_search 271 paid_search 195 unknown 179 social 75 direct_traffic 56 referral 24 email 15 display 6 other 4 other_publicities 3 Name: origin, dtype: int64</pre>	<p>179 observations have unknown origin. So, we will turn all the missing value at origin feature into unknown value</p>	<p>Tool: Python</p>	<p>All the missing values in origin of marketing channel fill with "unknown" value</p> <pre>merged_marketing_data.origin.value_counts() organic_search 271 paid_search 195 unknown 193 social 75 direct_traffic 56 referral 24 email 15 display 6 other 4 other_publicities 3 Name: origin, dtype: int64</pre>																																				
<p><u>Merged ecommerce data</u> (data file attained after merging other data files together)</p> <p>The “order purchase date” is not in a good format to analyse the trend by month and year</p>	<p>Convert to a year – month format</p>	<p>Step 1: Convert the column to datetime datatype</p> <p>Step 2: Using the map(lambda date: 100 * date.year + date.month) to convert to the new format</p> <p>Step 3: Sum over the month</p>	<p>Tool: Python</p> <p>Package: Pandas</p>	<p>We could calculate the sum of value by month and year</p> <table border="1"> <thead> <tr> <th>order_year_and_month</th> <th>seller_revenue</th> </tr> </thead> <tbody> <tr> <td>0</td> <td>201609</td> <td>215.28</td> </tr> <tr> <td>1</td> <td>201610</td> <td>32407.19</td> </tr> <tr> <td>2</td> <td>201701</td> <td>54771.28</td> </tr> <tr> <td>3</td> <td>201702</td> <td>117664.95</td> </tr> </tbody> </table>	order_year_and_month	seller_revenue	0	201609	215.28	1	201610	32407.19	2	201701	54771.28	3	201702	117664.95																						
order_year_and_month	seller_revenue																																							
0	201609	215.28																																						
1	201610	32407.19																																						
2	201701	54771.28																																						
3	201702	117664.95																																						

3. Data checking in detail

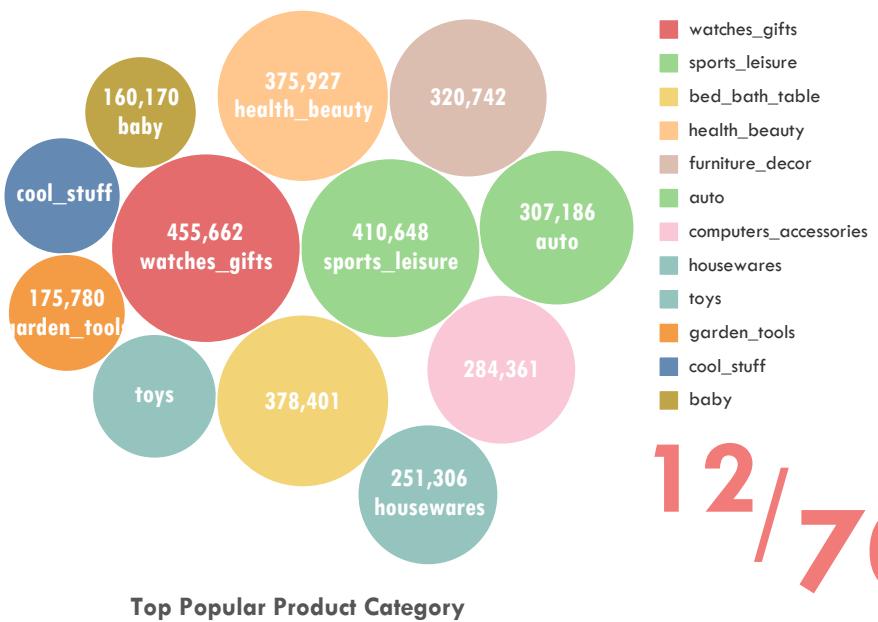
Problems	Actions	Methods (How - Why)	Software	Outcome																						
<p><u>Olist Seller Dataset</u></p> <p>Containing improper location format that seem fine at first during the cleaning process, however, cannot be recognized in Tableau geographic map</p> <p>Example:</p> <ul style="list-style-type: none"> +/ Using district name instead of city +/ Incorrect city name, typo +/Using state name instead city name 	<p>Matching location</p> 	<p>Google the location to identify the mistakes</p> <p>District name replaced by city name</p> <p>Matching with the default location in Tableau</p> <p>Unidentified name will be matched by using the nearby city or provide the latitude + longitude</p> <p>English name turns into Latino name</p>	<p>Tool: Tableau</p>	<p>Location format is corrected and recognized by Tableau</p> <table border="1"> <thead> <tr> <th>Your Data</th> <th>Matching Location</th> </tr> </thead> <tbody> <tr> <td>Andira-Pr</td> <td>Andirá</td> </tr> <tr> <td>Andradas</td> <td>Andradas</td> </tr> <tr> <td>Angra Dos Reis Rj</td> <td>Angra Dos Reis</td> </tr> <tr> <td>Auriflama/Sp</td> <td>Auriflama</td> </tr> <tr> <td>Balenario Camboriú</td> <td>Balneário Camboriú</td> </tr> <tr> <td>Picarras</td> <td>Balneário Piçarras</td> </tr> <tr> <td>Barbacena/ Minas Gerais</td> <td>Barbacena</td> </tr> <tr> <td>Barro Alto</td> <td>Barro Alto</td> </tr> <tr> <td>Mogi Das Cruzes / Sp</td> <td>Mogi Das Cruzes</td> </tr> <tr> <td>Novo Hamburgo. Rio Grande Do Sul. Brasil</td> <td>Novo Hamburgo</td> </tr> </tbody> </table>	Your Data	Matching Location	Andira-Pr	Andirá	Andradas	Andradas	Angra Dos Reis Rj	Angra Dos Reis	Auriflama/Sp	Auriflama	Balenario Camboriú	Balneário Camboriú	Picarras	Balneário Piçarras	Barbacena/ Minas Gerais	Barbacena	Barro Alto	Barro Alto	Mogi Das Cruzes / Sp	Mogi Das Cruzes	Novo Hamburgo. Rio Grande Do Sul. Brasil	Novo Hamburgo
Your Data	Matching Location																									
Andira-Pr	Andirá																									
Andradas	Andradas																									
Angra Dos Reis Rj	Angra Dos Reis																									
Auriflama/Sp	Auriflama																									
Balenario Camboriú	Balneário Camboriú																									
Picarras	Balneário Piçarras																									
Barbacena/ Minas Gerais	Barbacena																									
Barro Alto	Barro Alto																									
Mogi Das Cruzes / Sp	Mogi Das Cruzes																									
Novo Hamburgo. Rio Grande Do Sul. Brasil	Novo Hamburgo																									

4. Data Exploration

What I do	Visualization	Tool	Outcome
Dashboard 1: "Product Category By Year Table" Visualize the product category by year to see the change of them year by year	Text table	Tableau	Many product categories were not in 2016, had been added in 2017, and some has been removed in 2018 such as fashion_children_clothers , cds_dvds. Comment: Those changes might reflect the changes of Olist (e-commerce company) business strategy in diversifying the product category in order to recruit more shops and attract more customers
Dashboard 1: "Top Popular Product Category" Using the groupby method in python to group the revenue by product categories Sort the product category by descending	Bubble chart with the presence of top 12 product categories ranked by revenue	Tableau And Python (pandas)	Not all the product Category perform the same, it could be seen that top 12 categories was all appeared in 2016 list.
Dashboard 1: "Relationship between Active Customer and Seller Revenue" Using the groupby method in python to compute the sum of customer and revenue by month during time period	Combine Charts: Line chart + Area Chart	Tableau And Python (pandas)	It seem like the extending of product category significantly impacted to the customer attration. The active customer maintaining a quite stable growth trend, which boost the revenue increase simutaneously.
Dashboard 2: "Who is making money" – Answer the question which are the top sellers? Where are they? What do they sell? <u>Apply K-NN algorithm:</u> <ol style="list-style-type: none"> 1. Using the elbow method to find the optimal clusters value 2. Fit the sum of seller revenue - segment the seller by their total revenue 3. Using kmeans.cluster_centers_ to see the average of total revenue in each cluster 	Bar charts Symbol map	Tableau - Python – sklearn library	The purpose of this dashboard is to answer the question according to the seller as mentioned at the beginning: It could be seen that South regions of Brazil where the large cities and capital place has a dense distribution of seller. In fact, top sellers occupied about 3% in total number of sellers also located in this area (especially at Rio de Janeiro and Sao Paulo)

Filter and Group the total value by product category for each cluster	Pie Chart		The Top seller, in cluster 3, earning money by specialized in watches_gifts category while the lower rank one provide many different product types
Dashboard 3: “Customer distribution by Segment” – Answer the question: Where are the purchasers of top seller?	Symbol map Horizontal bar	Tableau	The purchaser spread over the coastal side of Brazil where have a high population density. However , the most active one are located around Sao Paulo and Rio de Janeiro Unsurprisingly, the majority of top customers come from both cluster 1 and 3 spend money on watches gifts
Answer the question: What are the marketing channels the company used to acquire them? After merging e-commerce data and the marketing data, I used value_counts() function in pandas to count the origin observation (marketing channel that used to recruit sellers) for each cluster	Dataframe	Python	It could be seen that the higher rank sellers do not join the platform from email, social, display marketing, so we could reduce the budget for those marketing, instead put more money on SEO(Search Engine Optimization) and SEM(Search Engine Marketing) to increase the organic_search and paid search

Dashboard 1: Overview

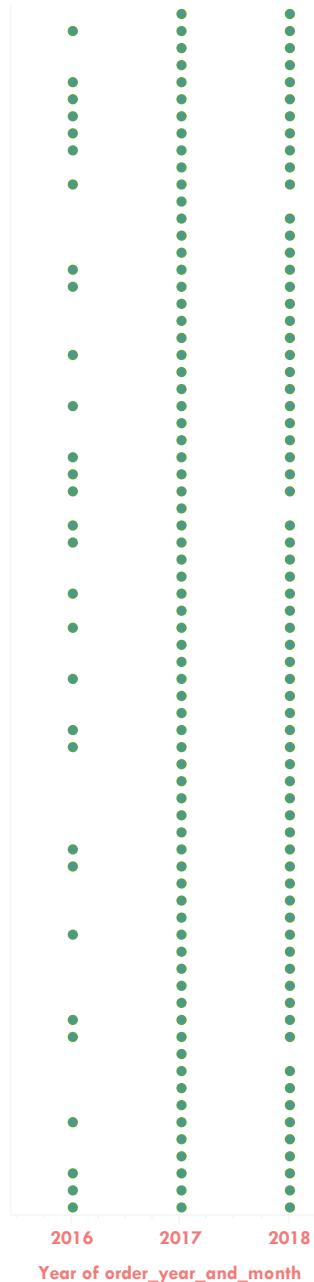


12/70 Product Category

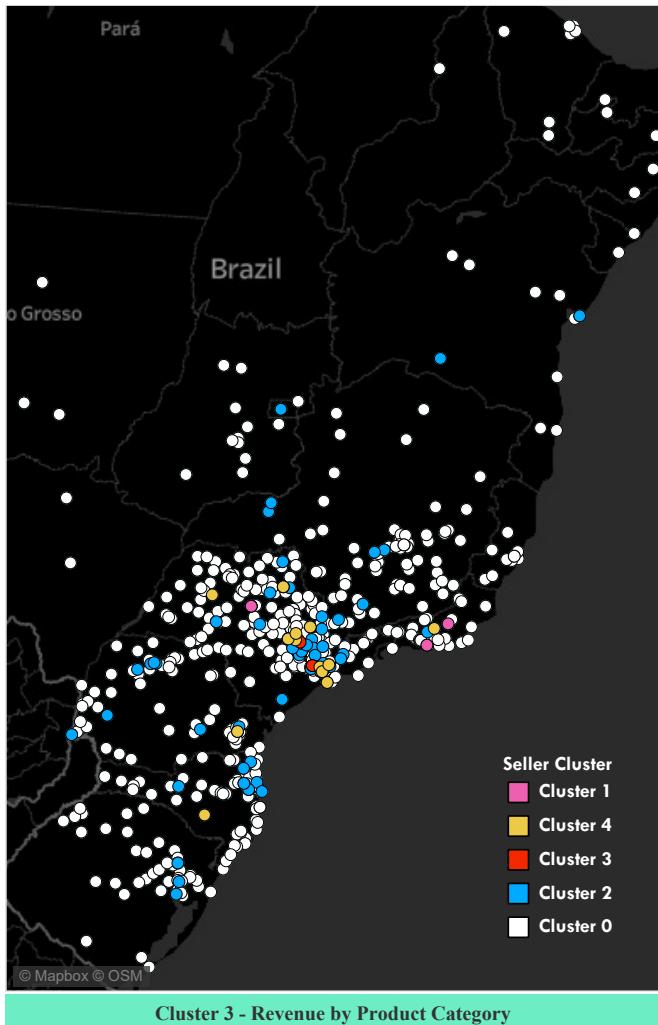


Product Category By Year

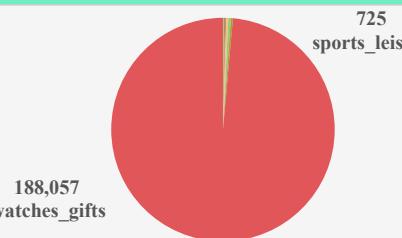
agro_industry_and_commerce
air_conditioning
art
arts_and_craftsmanship
audio
auto
baby
bed_bath_table
books_general_interest
books_imported
books_technical
cds_dvds_musicals
christmas_supplies
cine_photo
computers
computers_accessories
consoles_games
construction_tools_construction
construction_tools_lights
construction_tools_safety
cool_stuff
construction_tools_garden
construction_tools_tools
diapers_and_hygiene
drinks
dvds_blu_ray
electronics
fashio_female_clothing
fashion_bags_accessories
fashion_childrens_clothes
fashion_male_clothing
fashion_shoes
fashion_sport
fashion_underwear_beach
fixed_telephony
flowers
food
food_drink
furniture_bedroom
furniture_decor
furniture_living_room
furniture_mattress_and_upholstery
garden_tools
health_beauty
home_appliances
home_appliances_2
home_comfort_2
home_comfort
home_construction
housewares
industry_commerce_and_business
kitchen_dining_laundry_garden_furniture
la_cuisine
luggage_accessories
market_place
music
musical_instruments
office_furniture
party_supplies
perfumery
pet_shop
security_and_services
signalling_and_security
small_appliances
small_appliances_home_oven_and_coffee
sports_leisure
stationery
tablets_printing_image
telephony
toys
watches_gifts



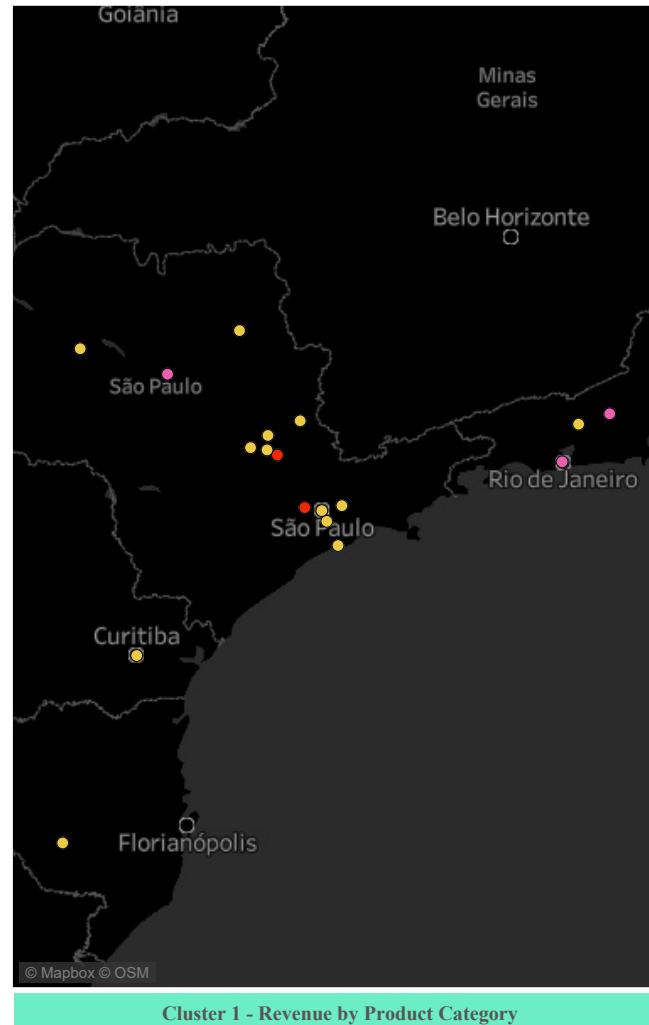
Seller Distribution - Brazil



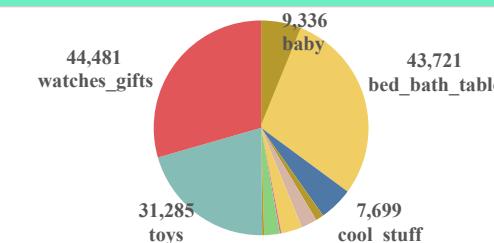
Cluster 3 - Revenue by Product Category



Top Seller Distribution



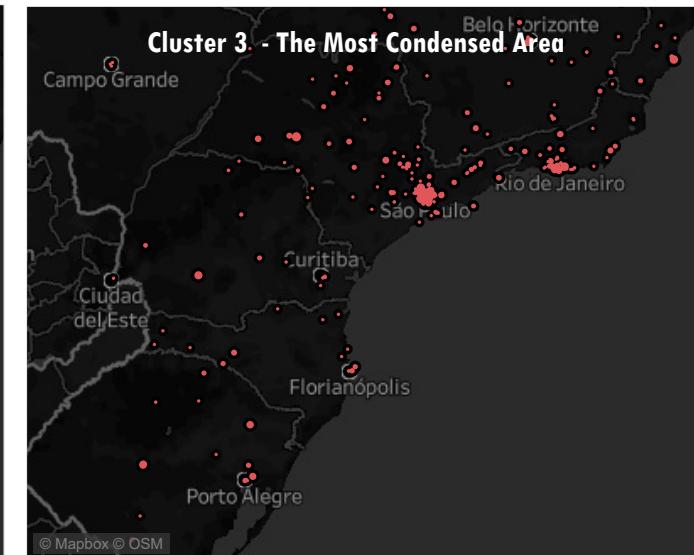
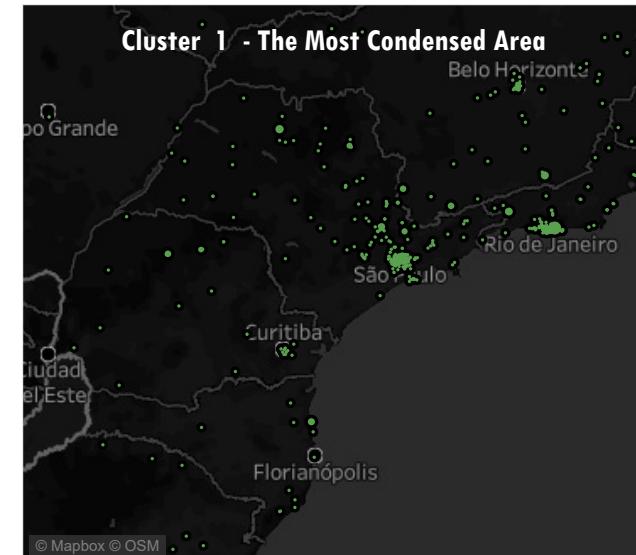
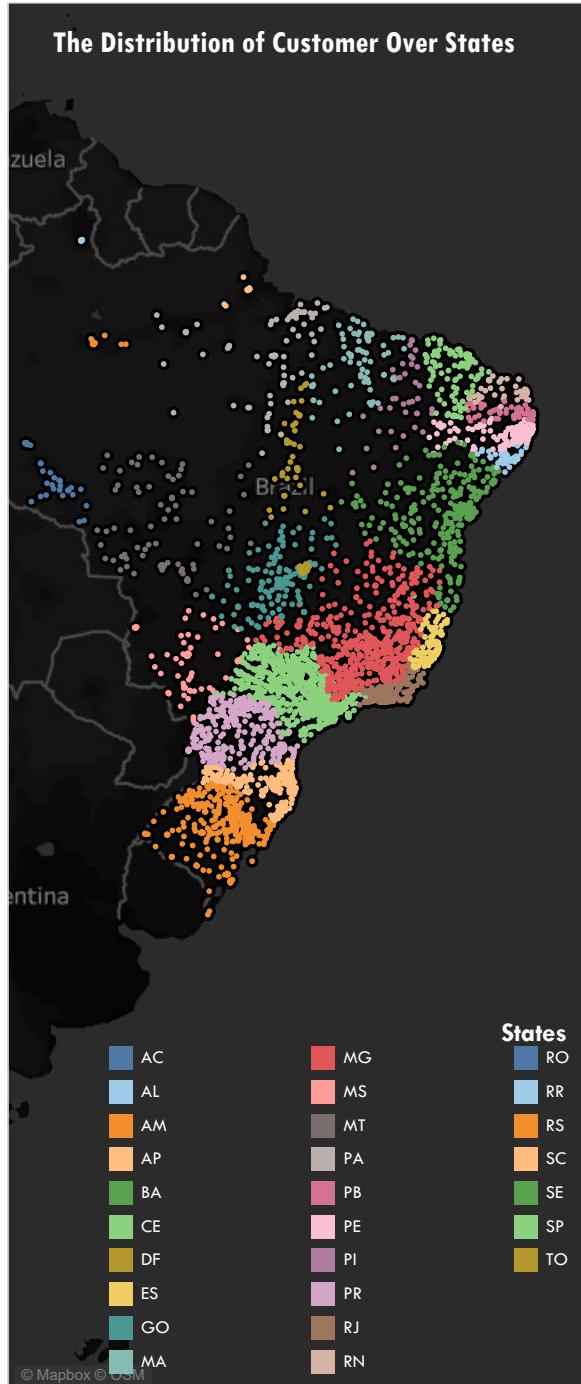
Cluster 1 - Revenue by Product Category



Cluster Segment Overview

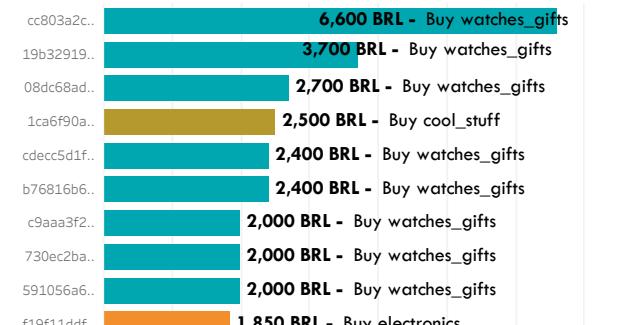


Dashboard 2
"Who is making money"



Top Ten Customers and Their Buying Category in Cluster 1 and 3

Customer ID - Cluster 1



Customer ID - Cluster 3



Dashboard 3
Customer Distribution by Segment

5. Conclusion

What I learnt from the data:

Expanding the product category draw more customers to the platform, so the company should maintain introduce the new category as well as improve the quality of product selling in the platform based on the customer reviews

Except the top 12 product categories, these other category performances are not as good as expected. To diversify product, we could search for the shops who has advantages in other categories as well. Firstly, we could start from the big states such as Rio or Sao Paulo where our top sellers located to find the potential candidates. The marketing method could be used are Search Engine Optimization and Search Engine Marketing, besides cut off other inefficient marketing channels such as email or display.

Watch gifts is the most popular product and the customer who spending the most are in the also gather around the big cities and states in Southern of Brazil. Obviously, the marketing budget should be allocated more to this area instead of spreading widely to another city and states.

6. Reflection

What I learnt in this project:

1. Understanding the meaning of each feature is a crucial step. In this data, amount of product is record as a sequence that created many duplications that does not help to answer our questions
2. The importance of data cleaning. The “city” feature includes many improper data formats and some of them is difficult to recognized such as name of district had been used instead of city name, or some of location cannot be found on the Tableau map
3. Using different type of charts and graphs in Tableau to obtain the insights

What in hindsight you might have done differently: select a smaller number of clusters maybe 3.