

Keystone Project

# Customer Satisfaction Prediction in E-commerce



Team C: Carol Wang, Yue Qiu, Jihua Chen, Yujie Zhang, Zinan Li



## Keystone Project

- Problem & Objectives
- EDA
- Feature Engineering
- Models
- Evaluation

# **Problem definition**

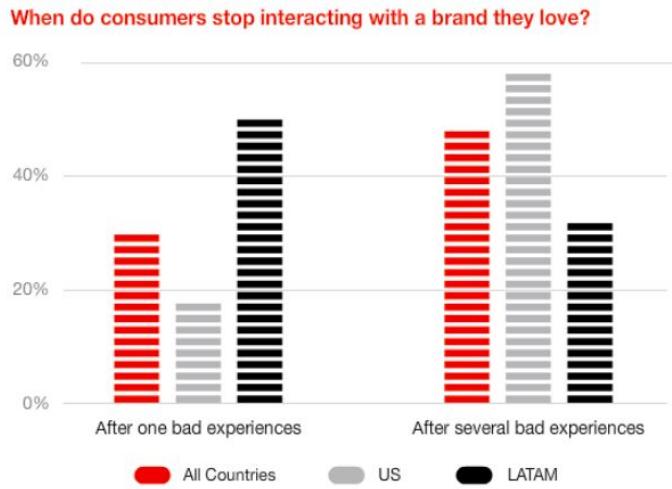
The scale of online retail business is constantly expanding with customers' increasing demand for online shopping. However, some customers did not leave a review which is challenging for the ecommerce company to keep track of each product of each seller.

# **Objectives**

Build a customer satisfaction classifier to help the ecommerce company to predict a review score

- Long term: To optimize the strategy with business and marketing team with top important features for a greater service and user experience
- Short term: To accurately predict the review score when a customer did not leave one; some measures may be implemented in a timely manner to reduce the bad reviews before customers give one.

# Why does customer satisfaction matter?



## Business impact:

- 17% personal consumption is buying online
- \$33.8 B online spending in 2021
- \$315 spending per consumer per year

## Project goal:

- Prediction on customer behavior
- Improve bottlenecks
- Optimize inventory and supply chain

# Ecommerce Sales Funnel

A generalized sales cycle for ecommerce stores



## 1. Visits Store

Visitor lands your website. Using a tracking pixel, you can begin re-marketing ads.

## 2. Views Product

Visitor has shown interest in a product. Pop up live chat window to answer any questions.

## 3. Starts Checkout

If you've captured an email address, send visitor reminders/coupons to complete purchase.

## 4. Offer Upsells

Before finalizing the purchase, show other/related products they might be interested in.

## 5. Complete Purchase

With a purchase made, continue sending marketing emails and/or coupons to encourage repeat purchases.

# TOC of EDA Section

1. Single Distribution
  1. Correlation: Time
  2. Correlation: Revenue
  3. Correlation: Score
  4. Correlation: Geo
1. Threshold for Satisfaction

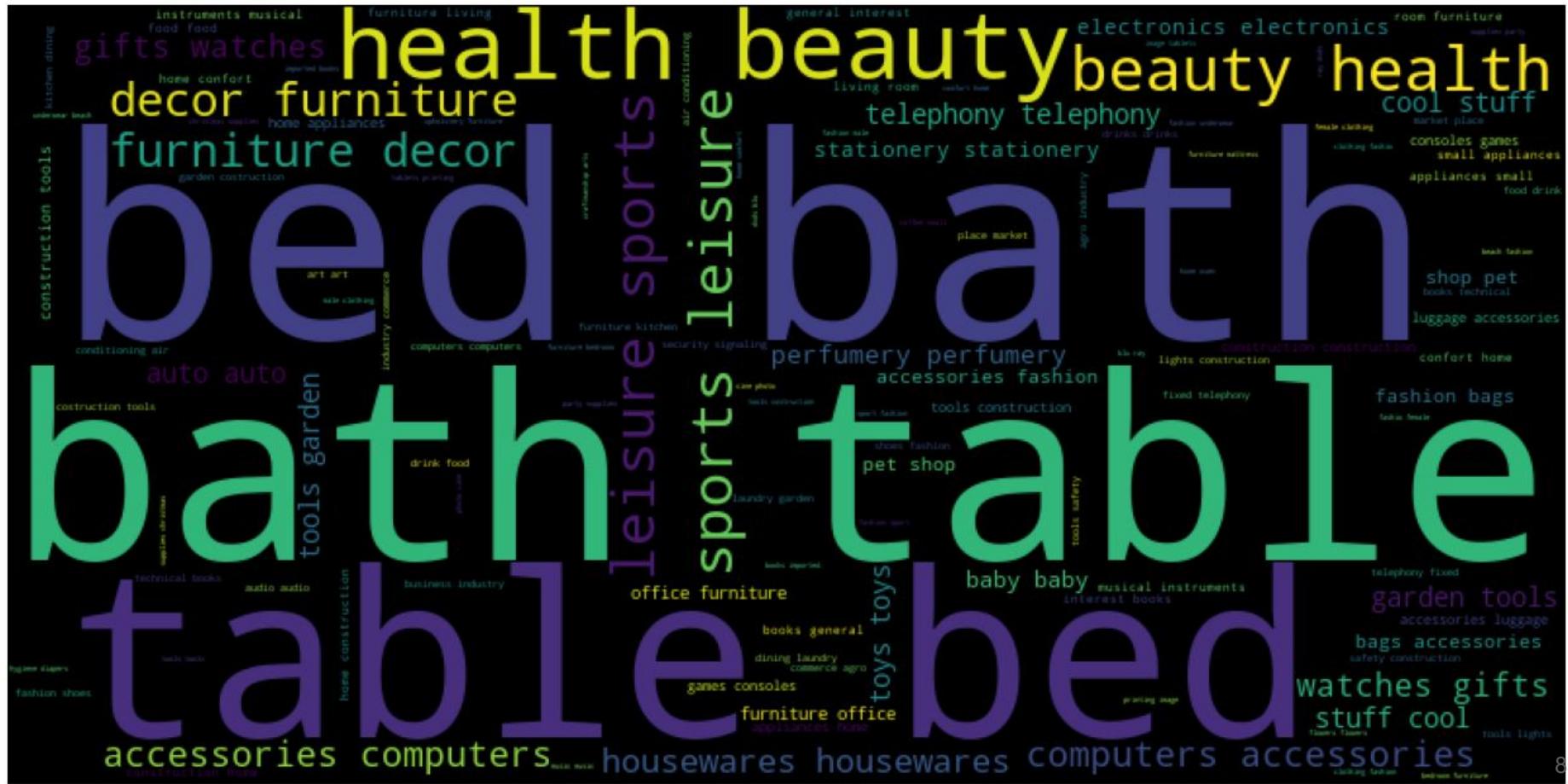
# 1. Single Distributions

Score and Order Status

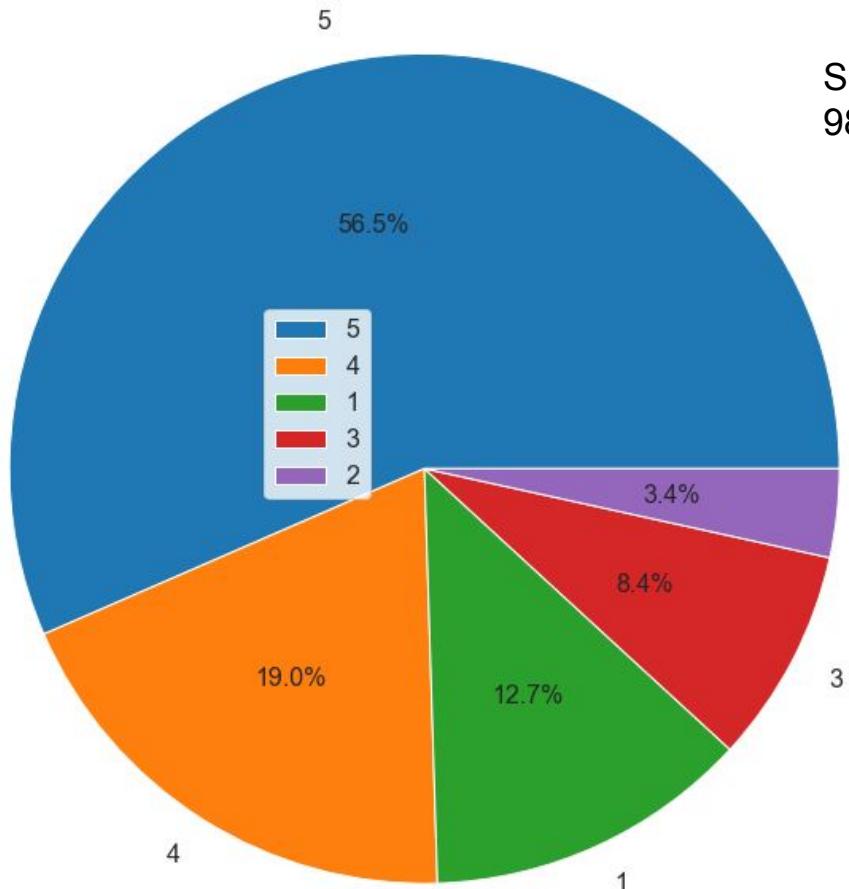
Seller State

Many other features are discussed last time

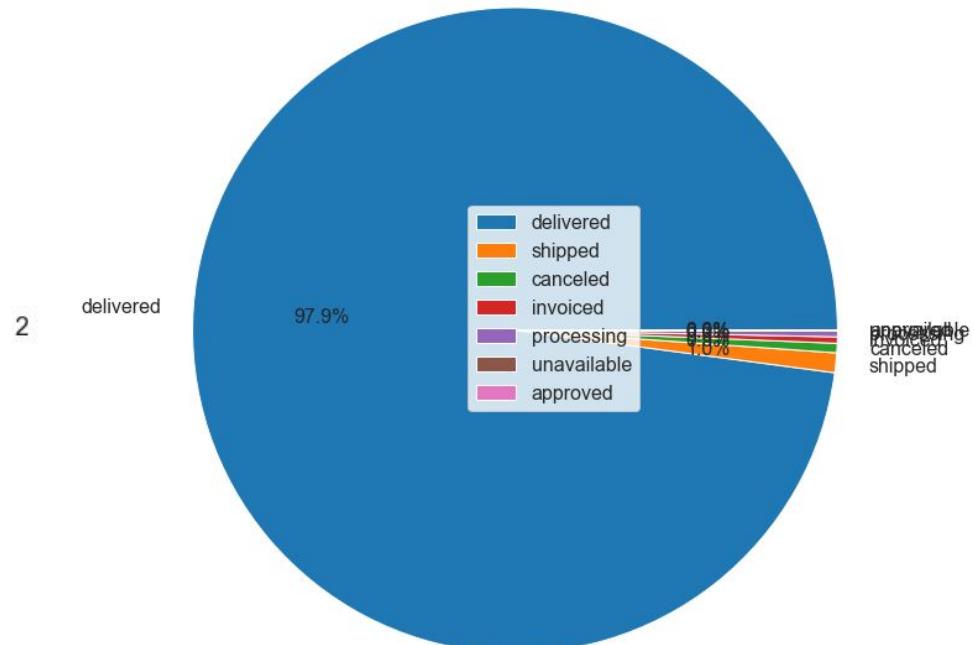
# Word Cloud of Product Category



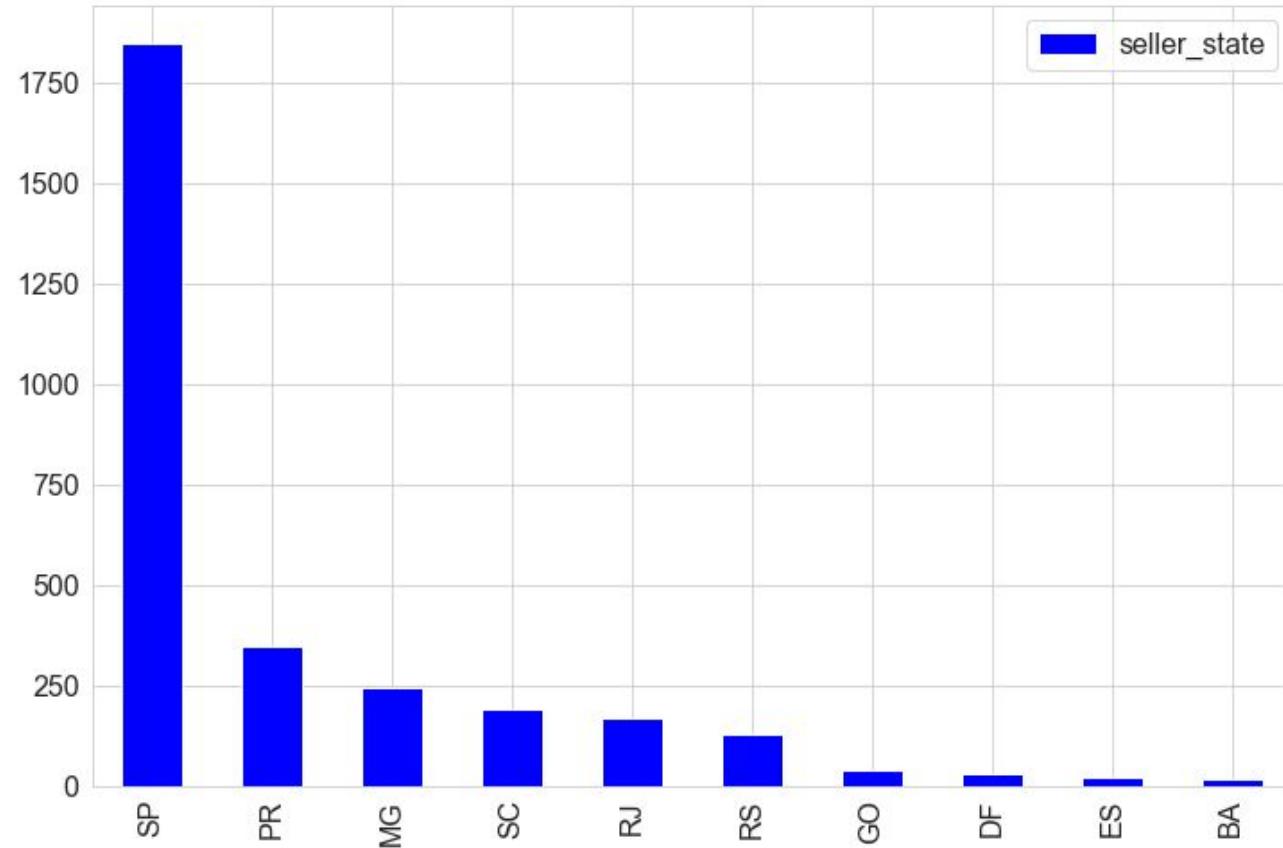
# Rating and Order Status



Skewed data with low ratings.  
98% order delivery rate.



# Seller Count vs State



## 2. Correlations: Time

Sales: 2017 vs. 2018

Order Time vs State

Order vs. Hours

Order vs. Month and Weekday

Score vs. Est. Time

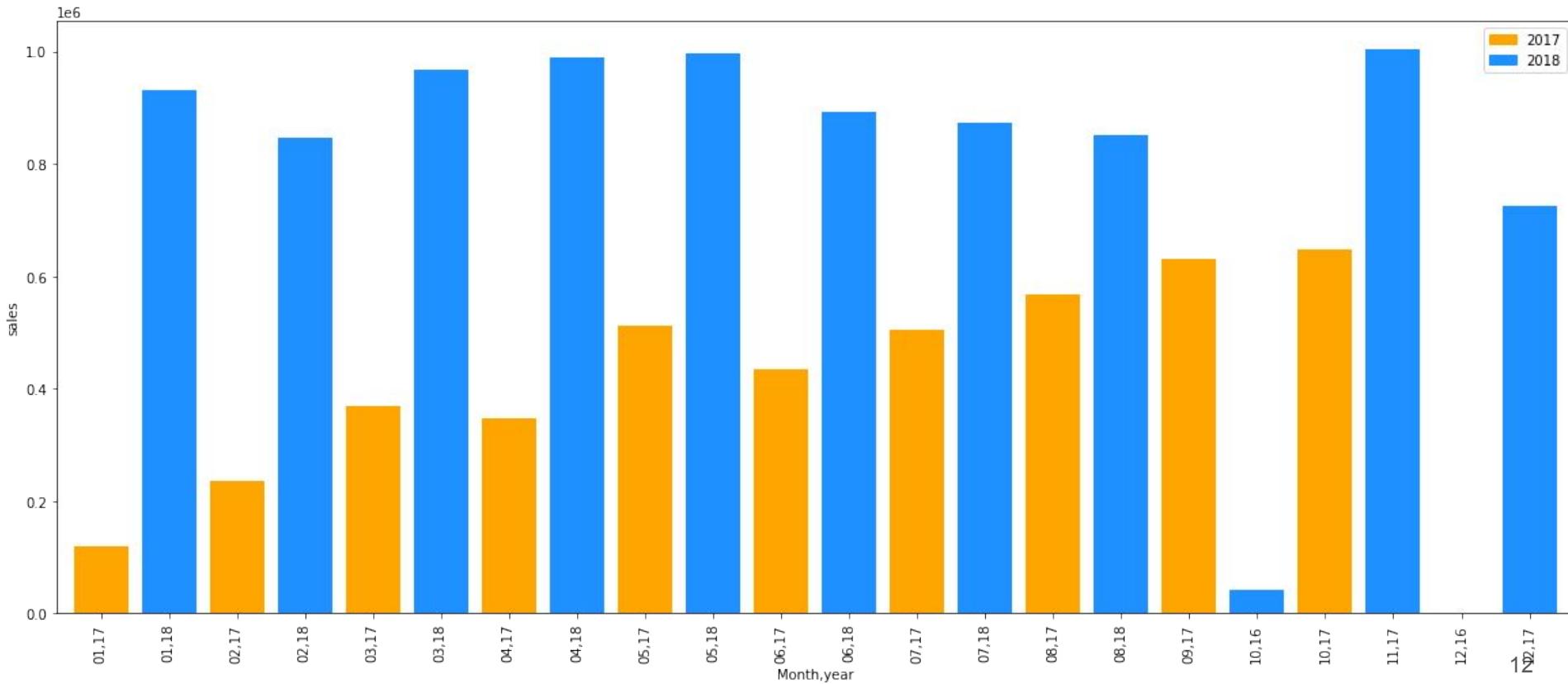
Score vs. Act. Time

Score vs. Diff btw Est. & Act. Time

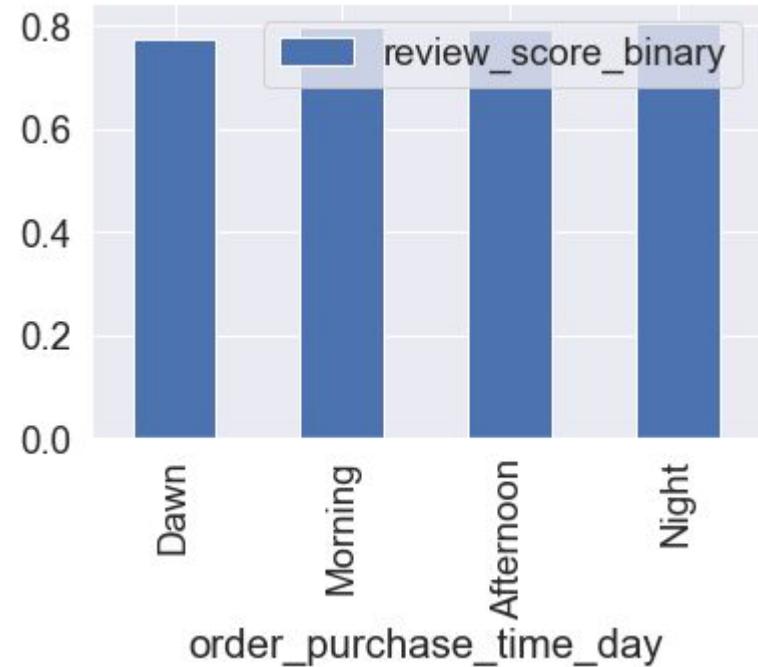
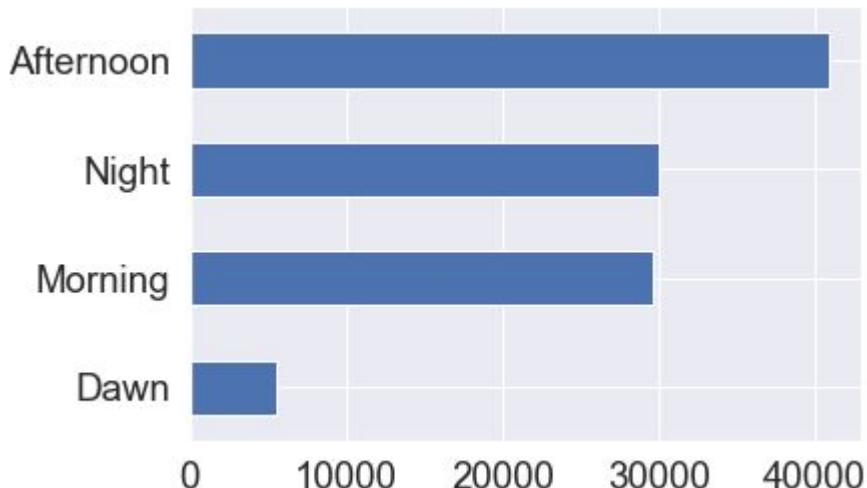
Score vs. Diff btw Purch. & Approv.

Score vs. Diff btw Purch. & CXR.

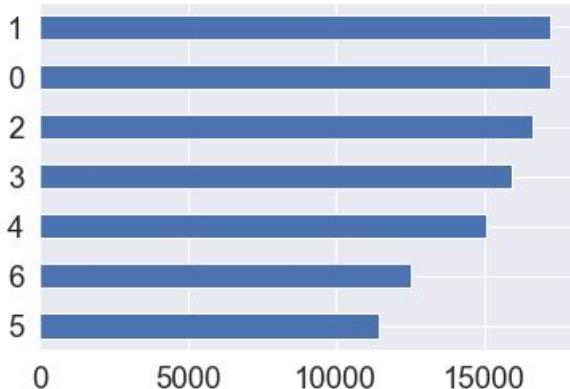
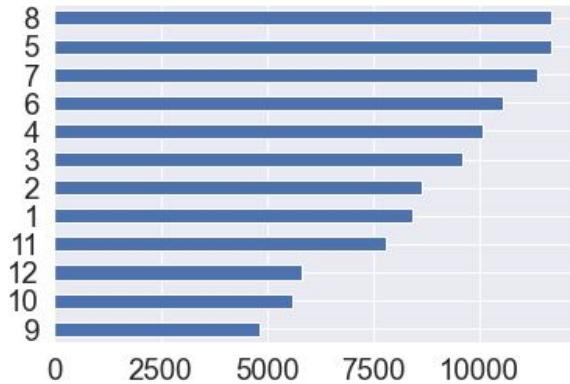
# Sales: 2017 vs 2018 Monthly Data



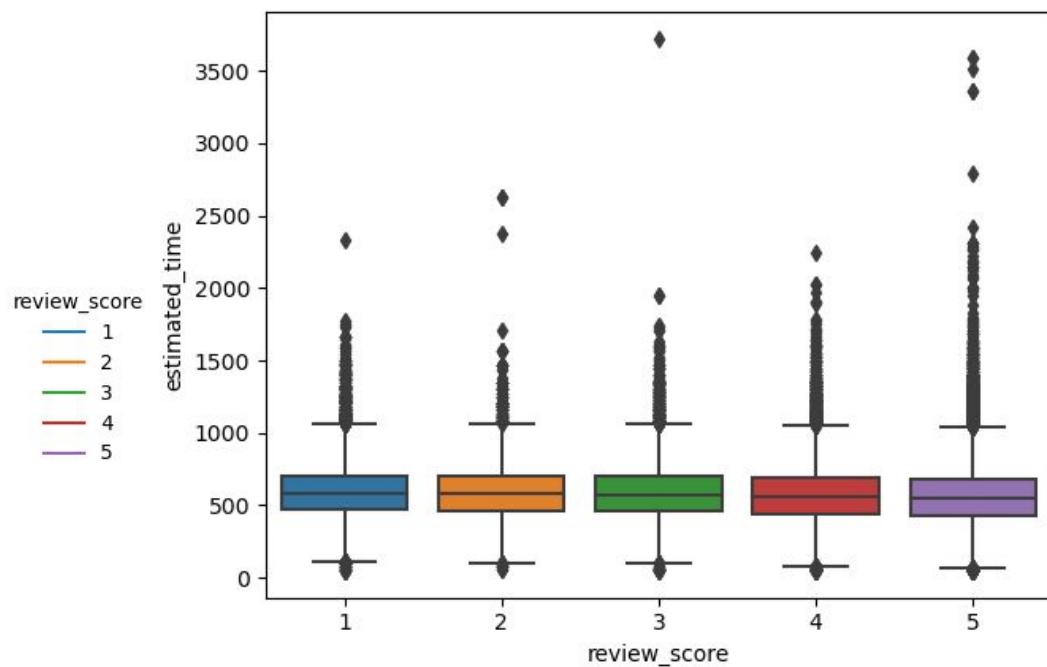
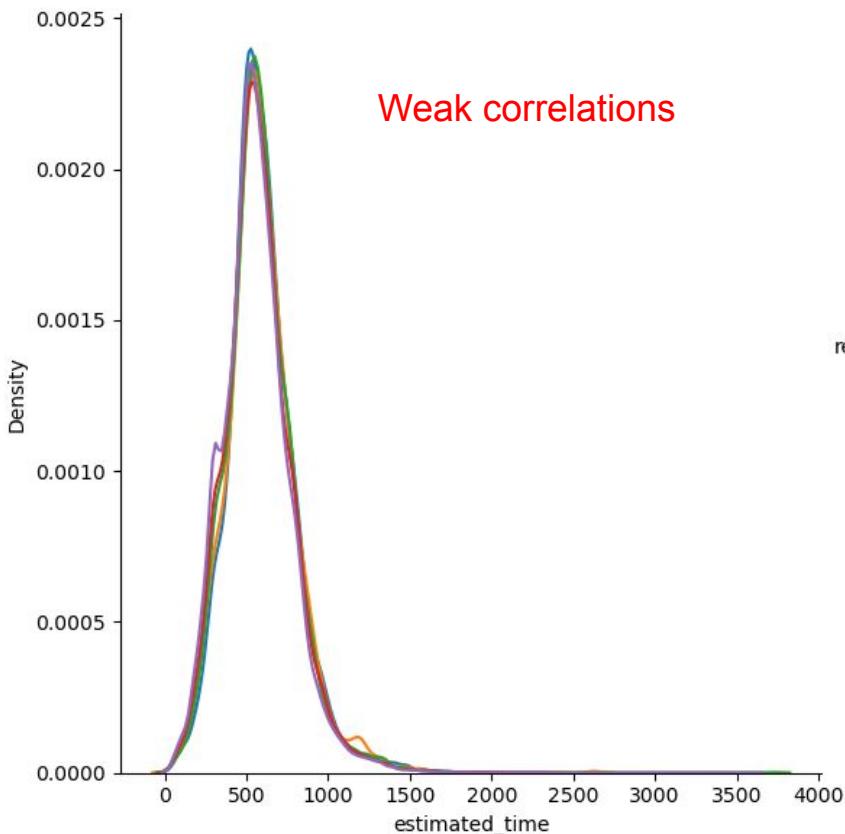
# Order vs. Hours



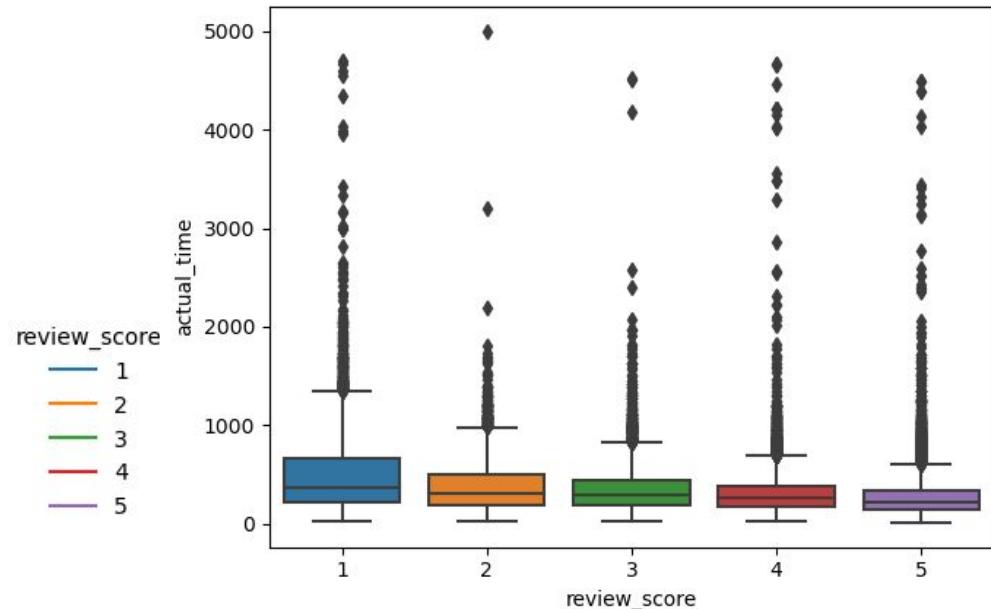
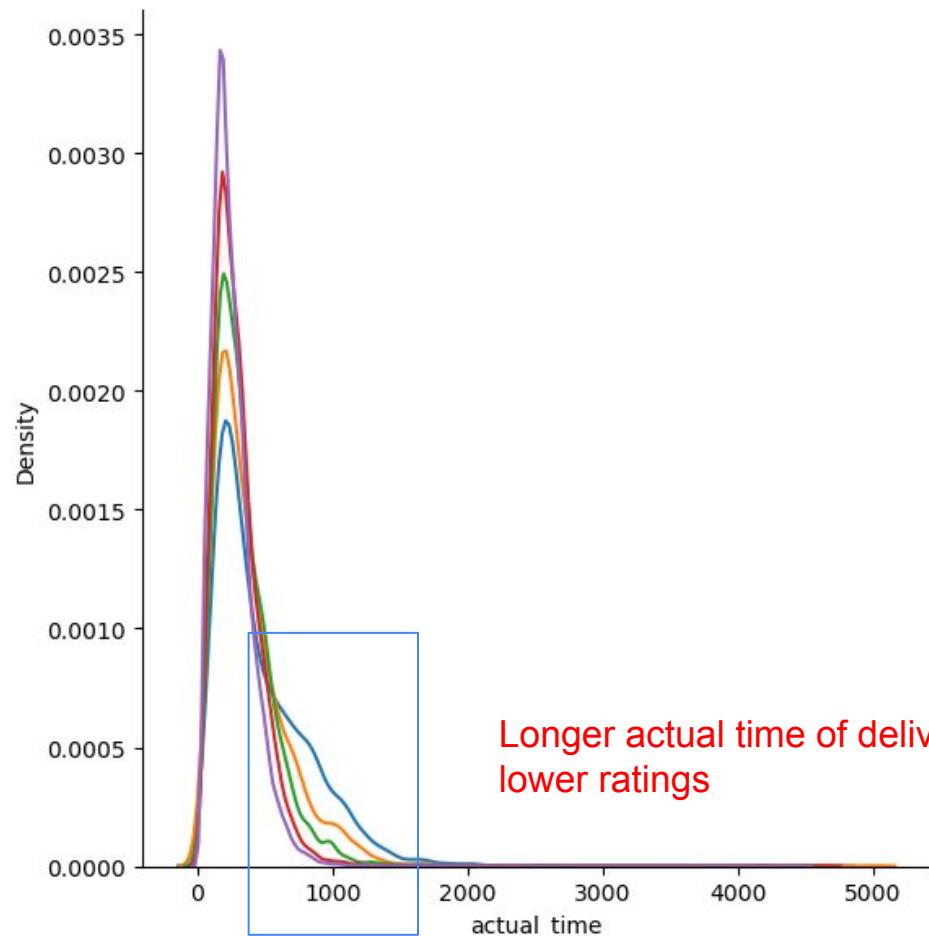
# Order vs. Month and Weekday



# Review Score by Estimated Delivery Time



# Review Score by Actual Delivery Time



### 3. Correlations: Revenue

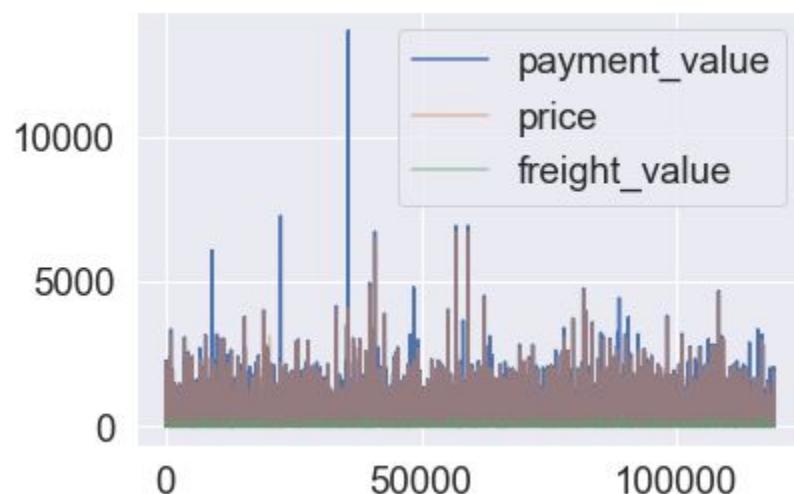
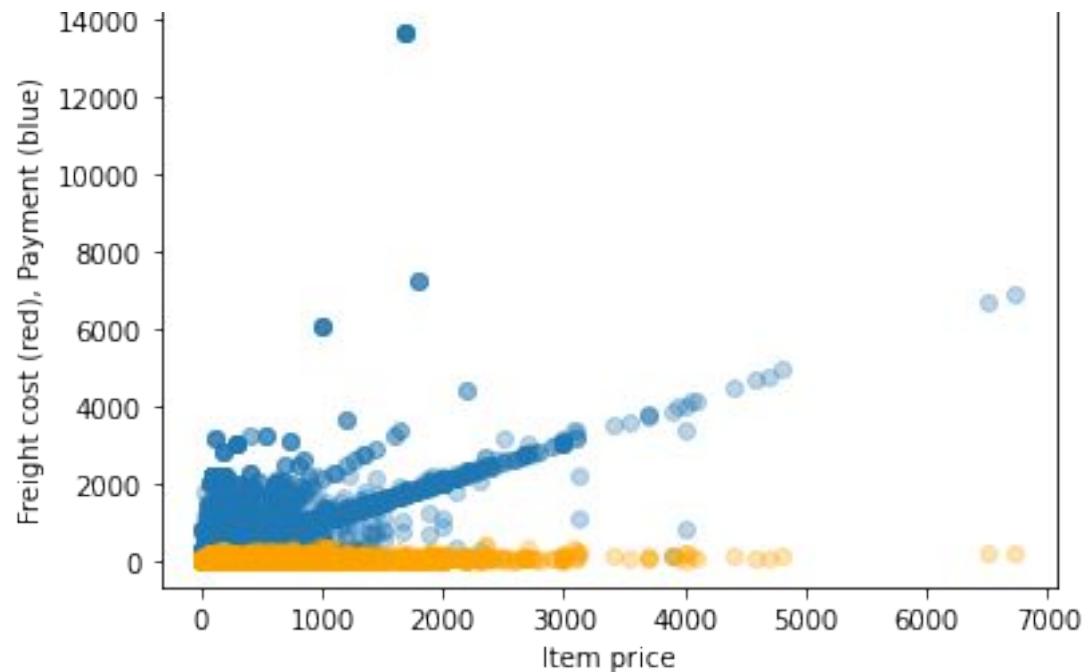
Price vs Freight cost and Payment

Category by counts or revenue

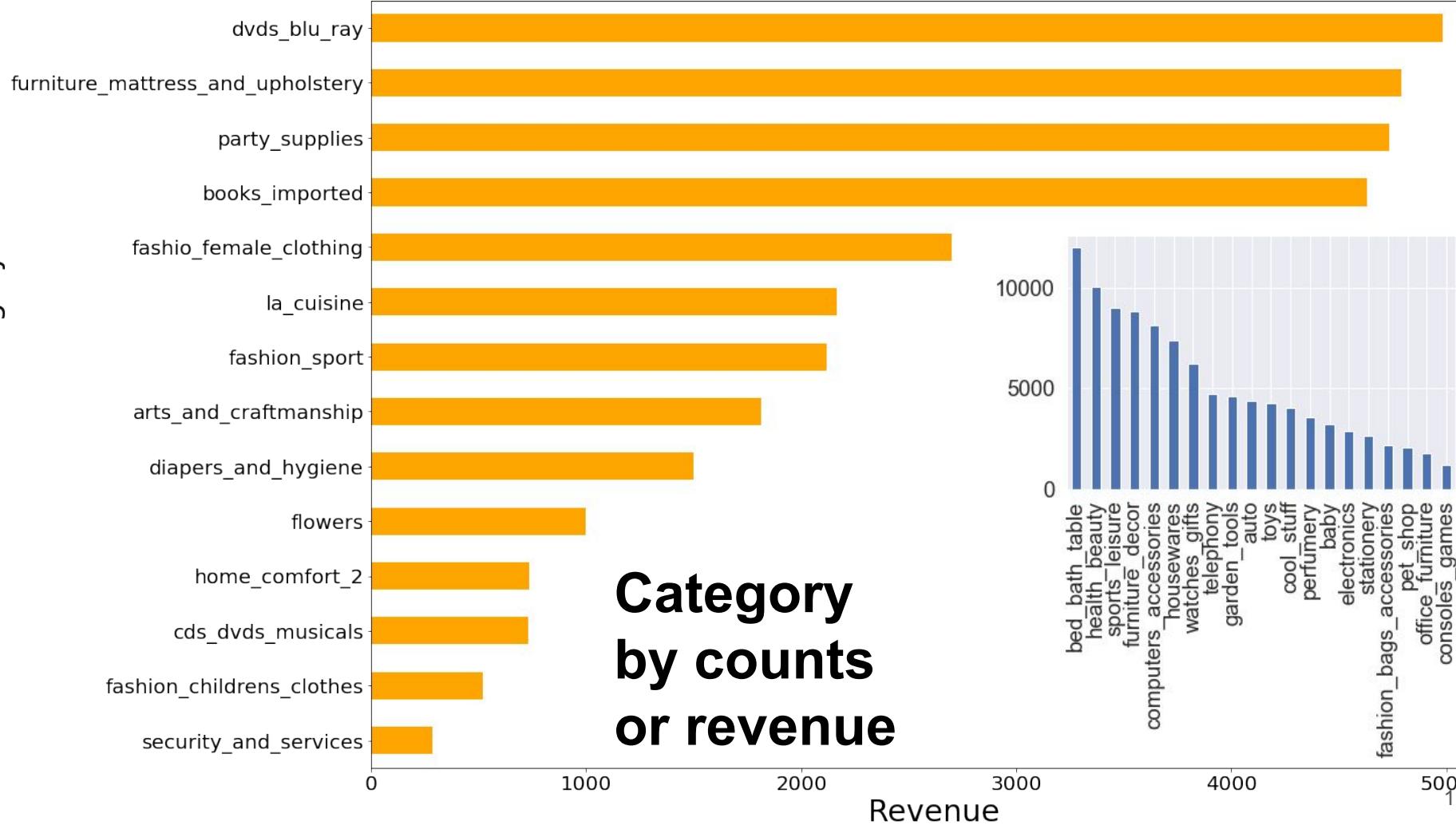
Freight Value vs. Product Weight

Score by Price

# Price vs Freight cost and Payment



Product Category



# 4. Correlations: Score

Score vs. Price

Score vs. Customer State

Score vs. Payment

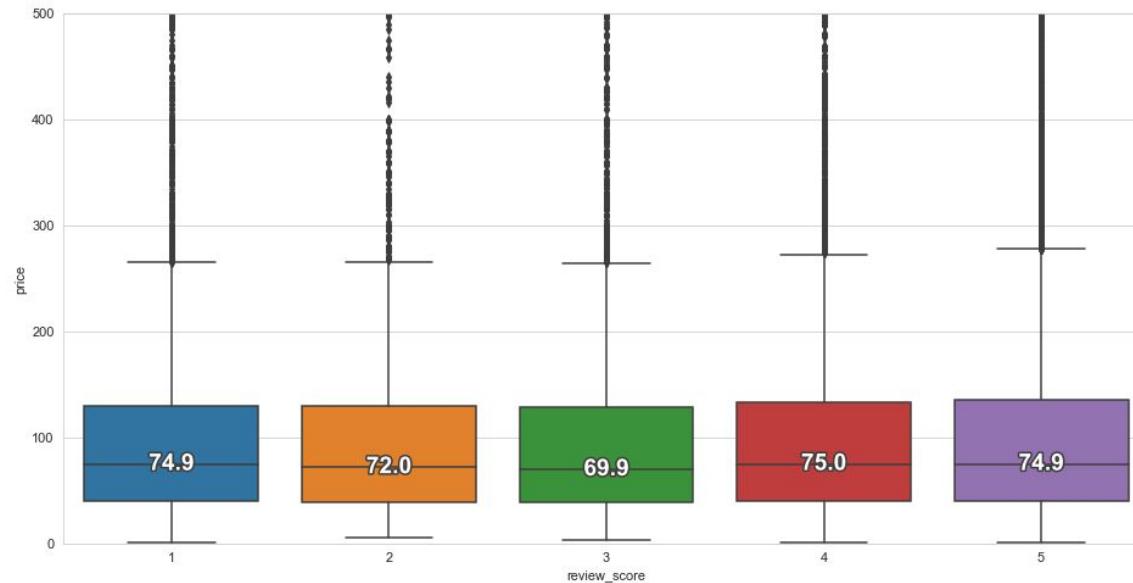
Score vs. Payment Installments

Score vs. Payment Sequential

Score vs. Category and Status

Score vs. Description and Name

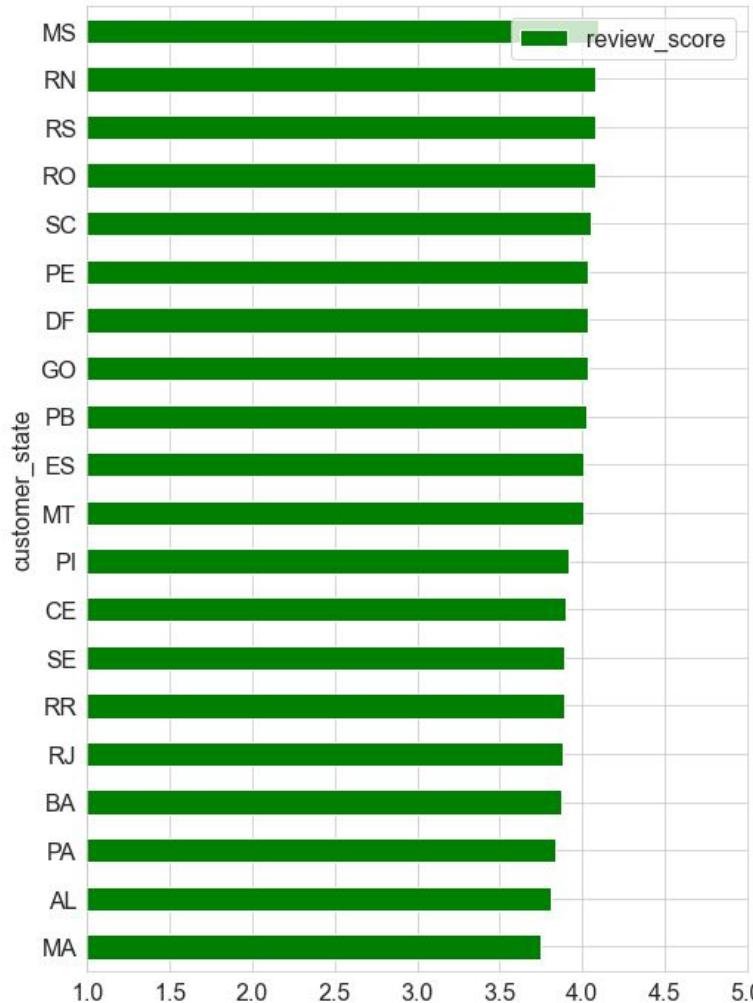
# Score by price

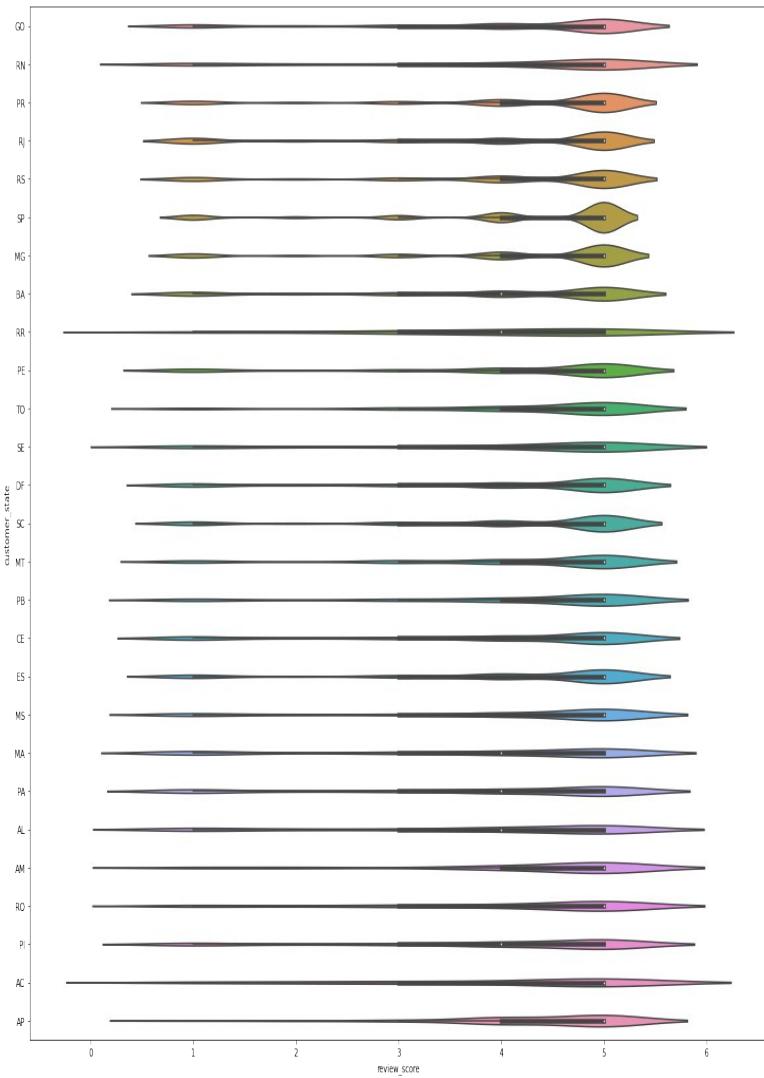


Weak correlations

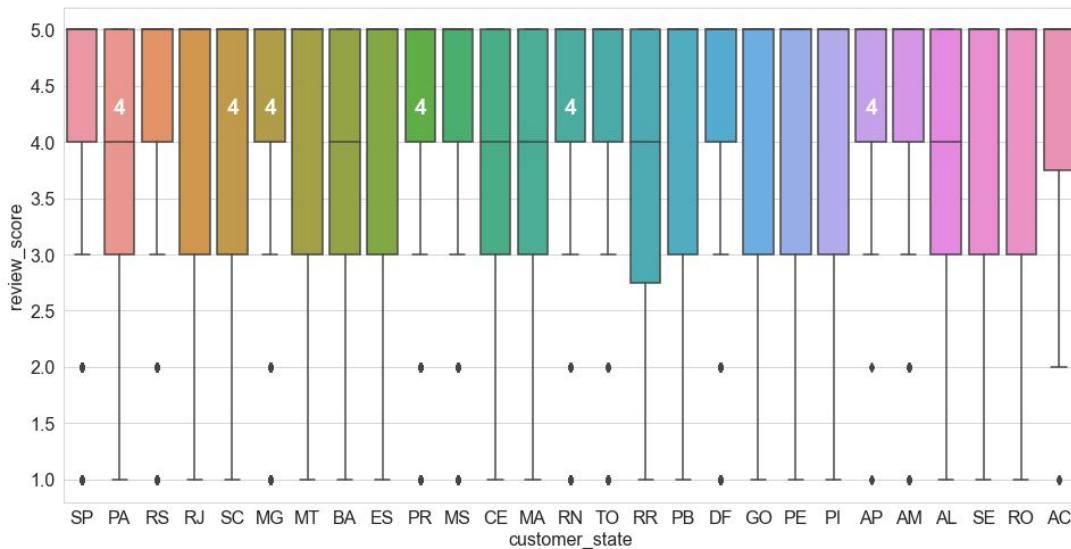
# Score vs. Customer State

Some correlations to the customer state





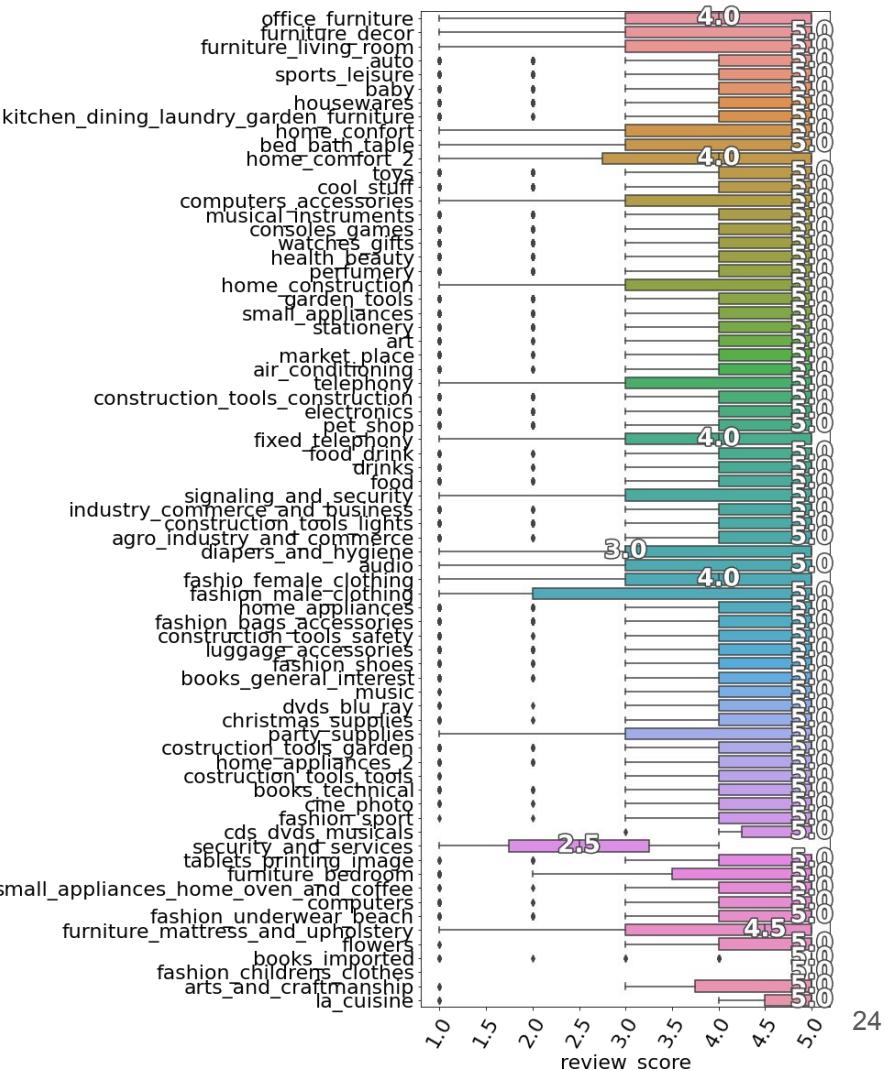
# Score vs. Customer State



# Score vs. Category

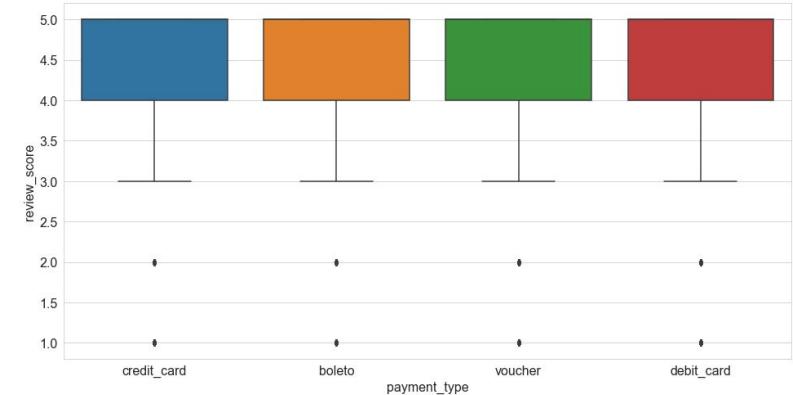
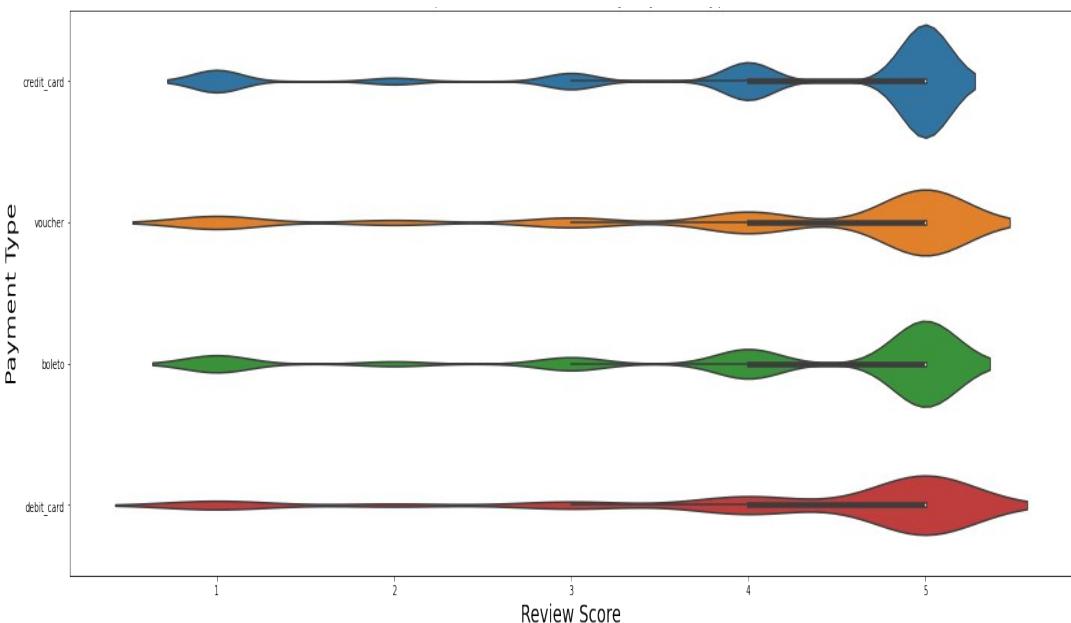
Certain product category has lower review scores in general

- Diapers and hygiene
- Security and services
- Fashion-male clothing



# Score vs. Payment Type

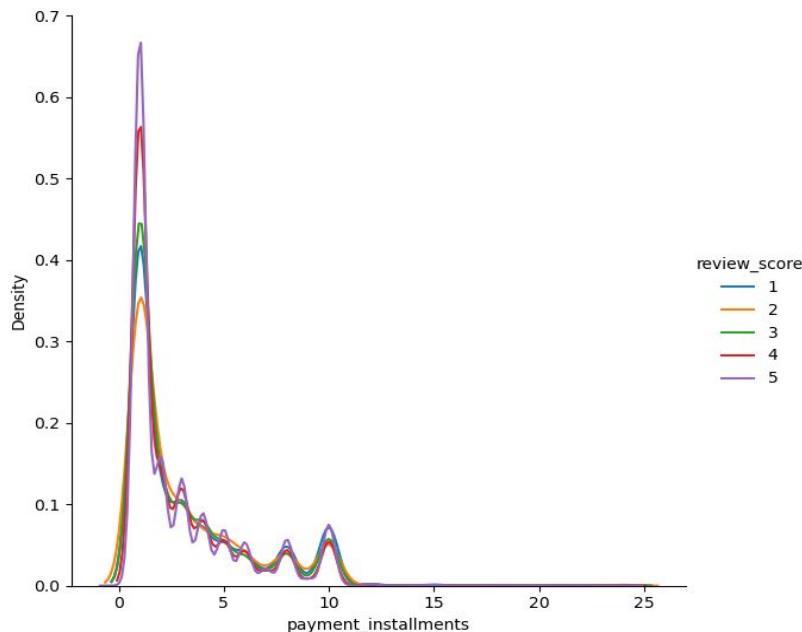
More concentrated 5 for two of the payment.



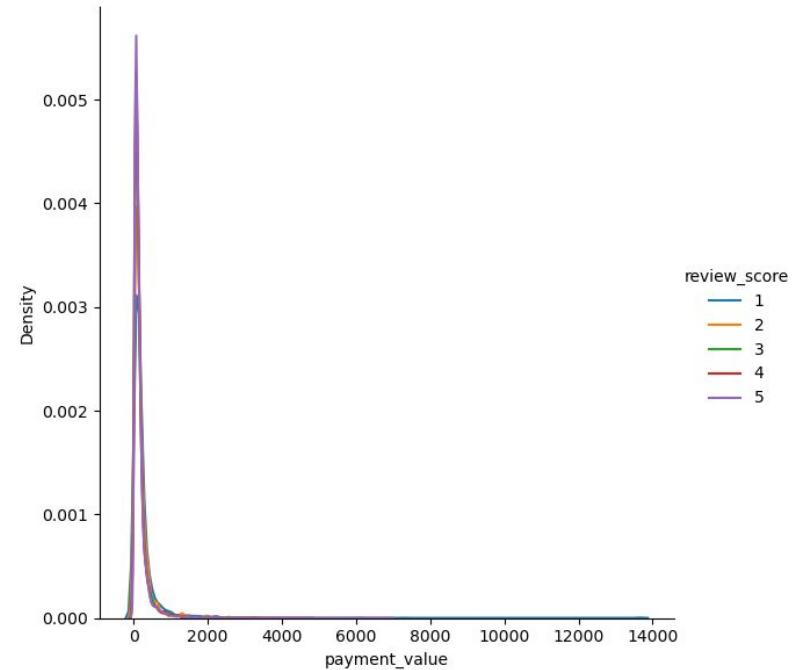
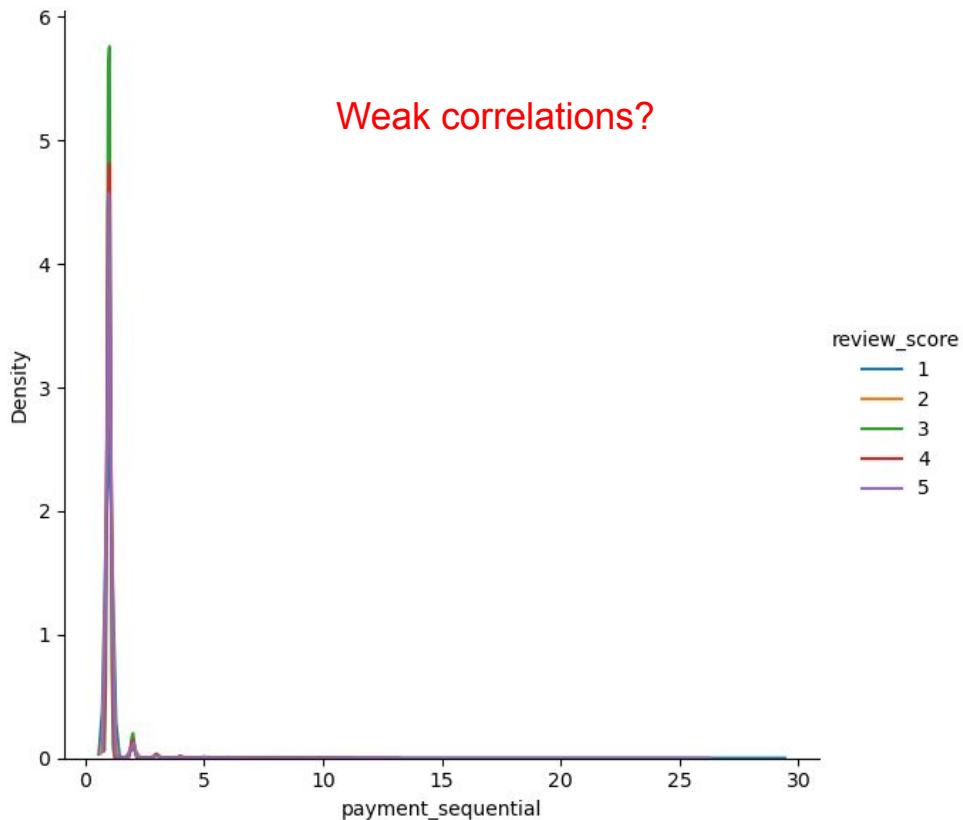
# Score vs. Payment Installment

Generally weak correlation

When installment > 20, score trends towards the lower end

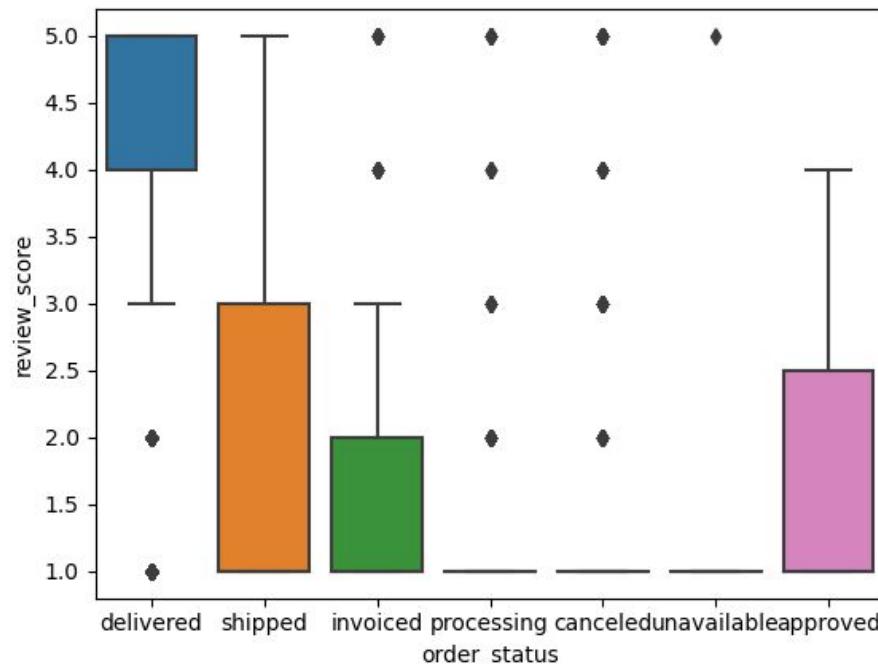


# Score vs. Payment Sequential



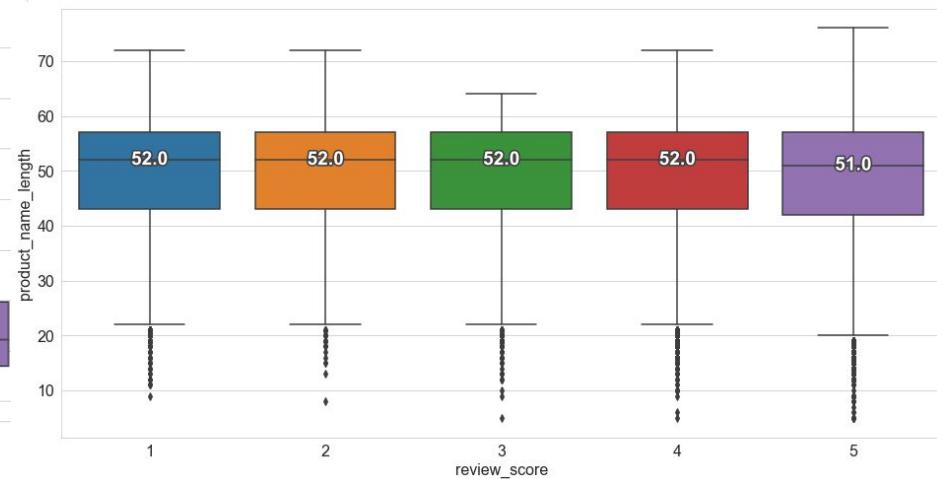
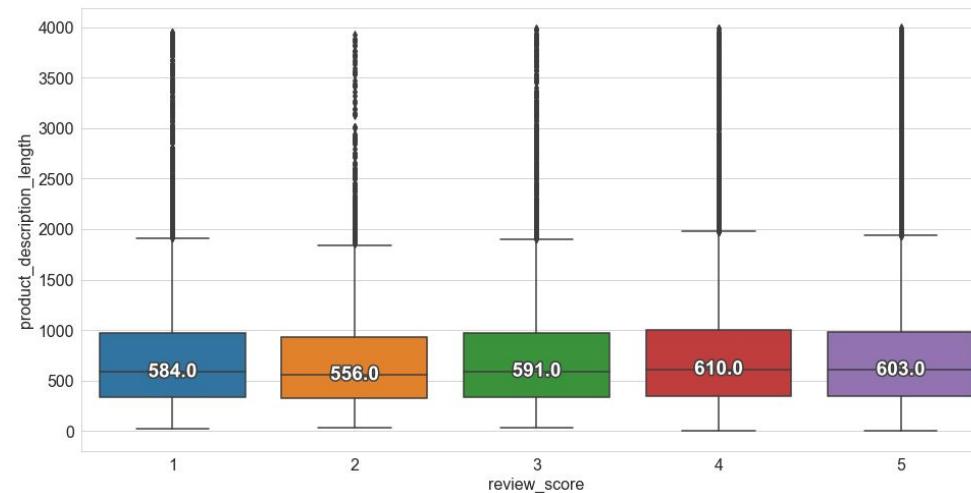
# Score vs. Order Status

Strong Correlations!

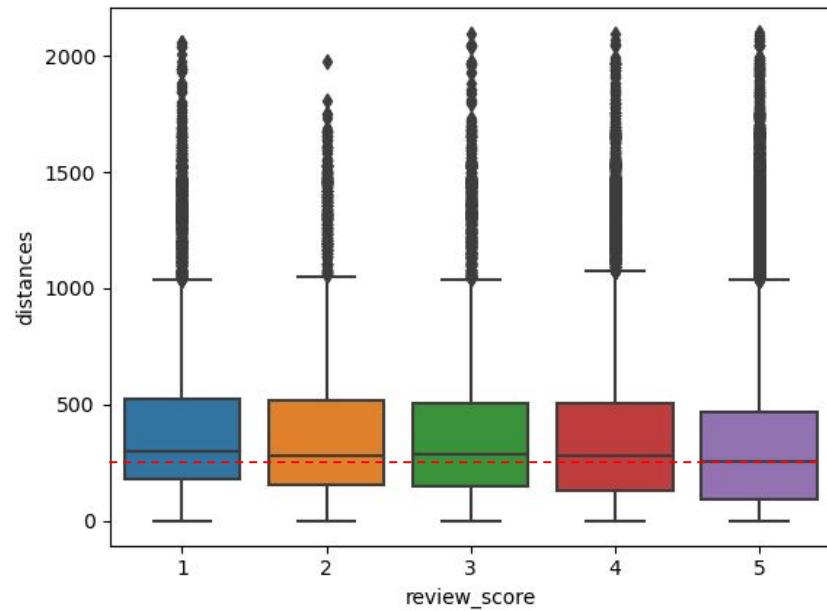
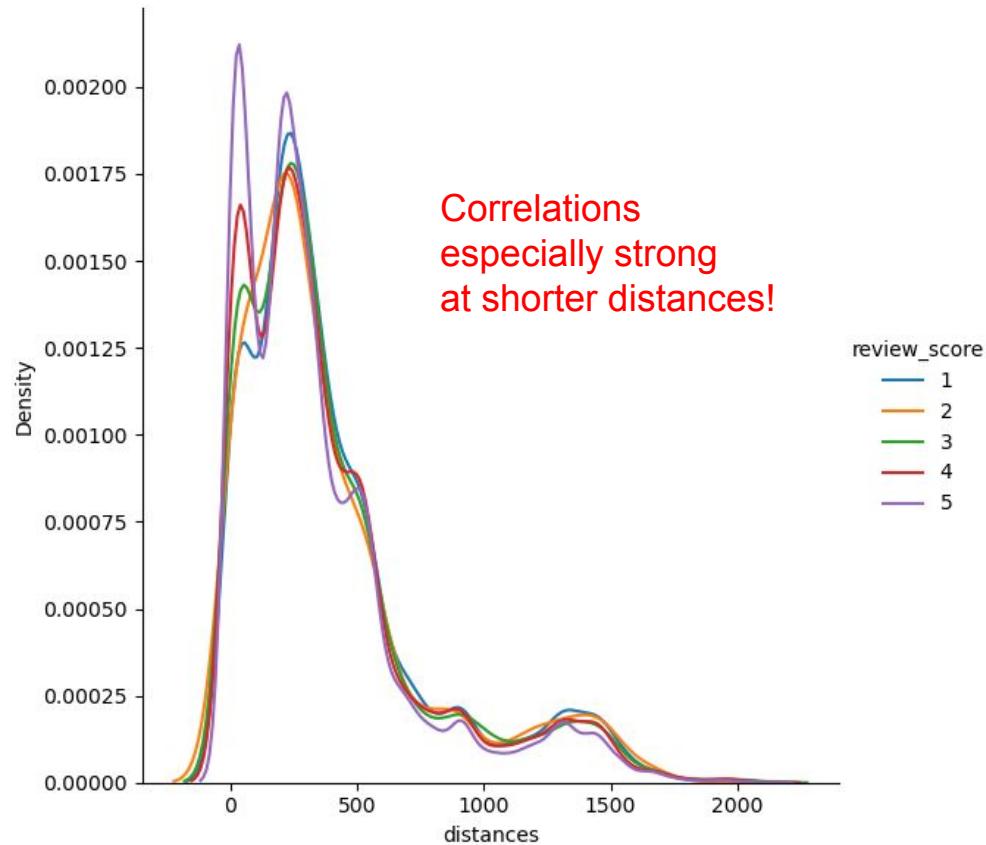


# Score vs. Description and Name

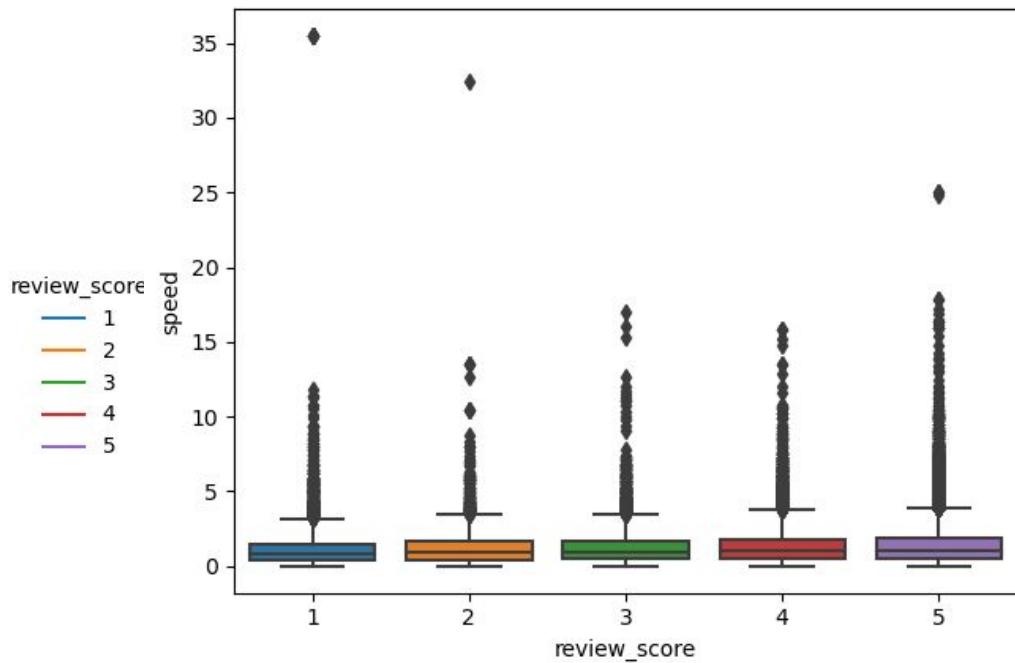
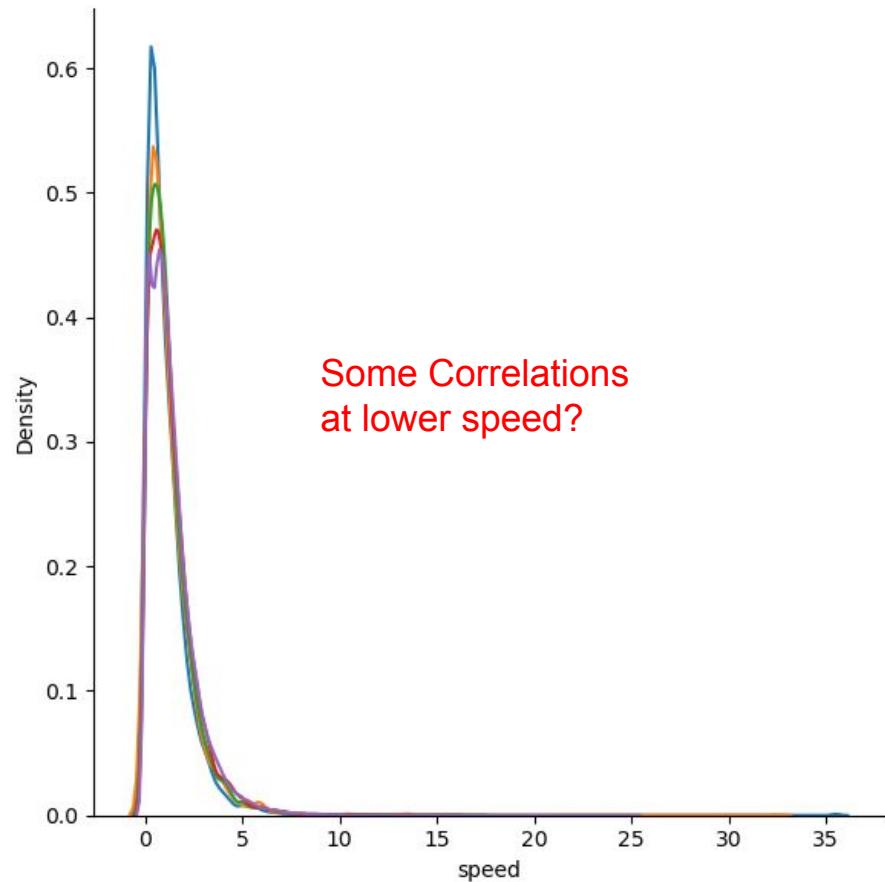
- Longer description has higher frequency of score 4 and 5.
- Product name length has no correlation with score



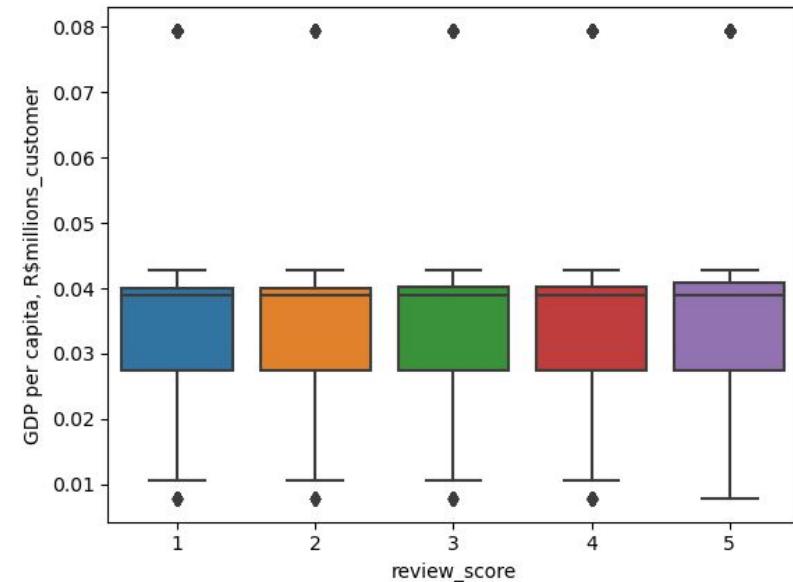
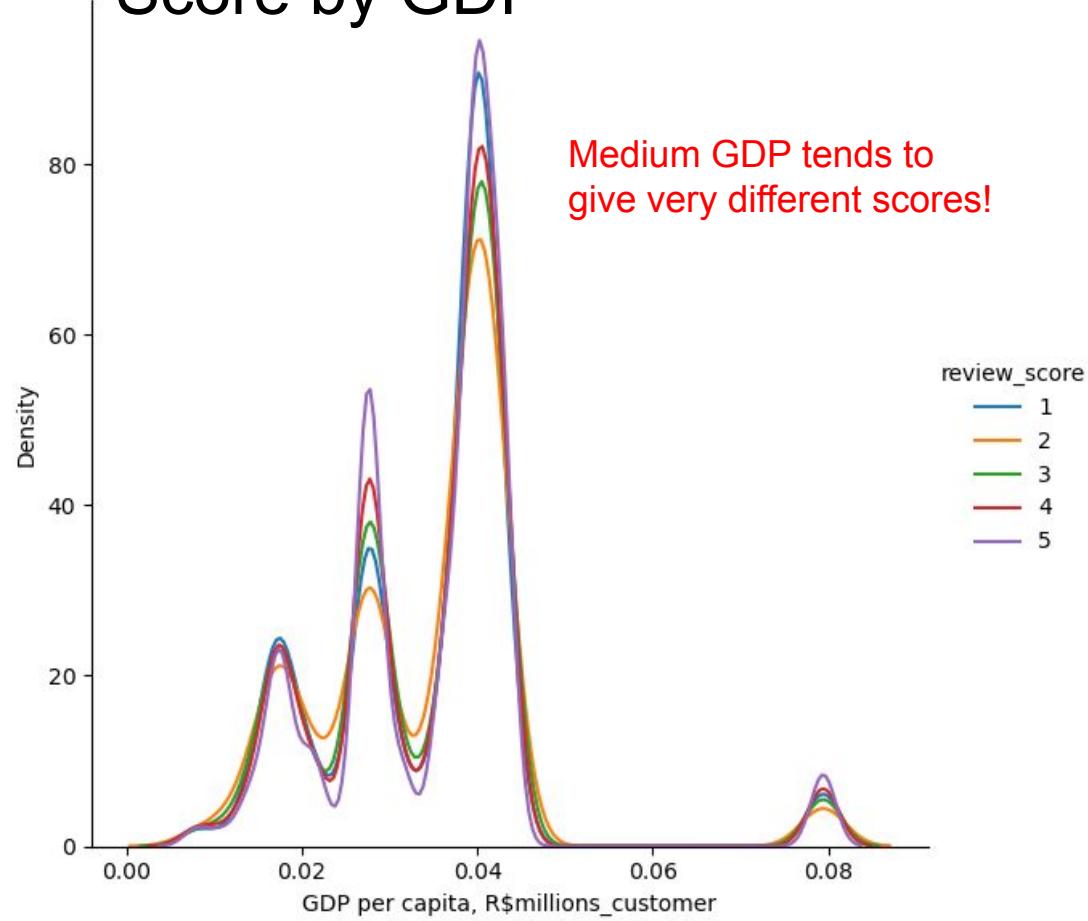
# Score by distance



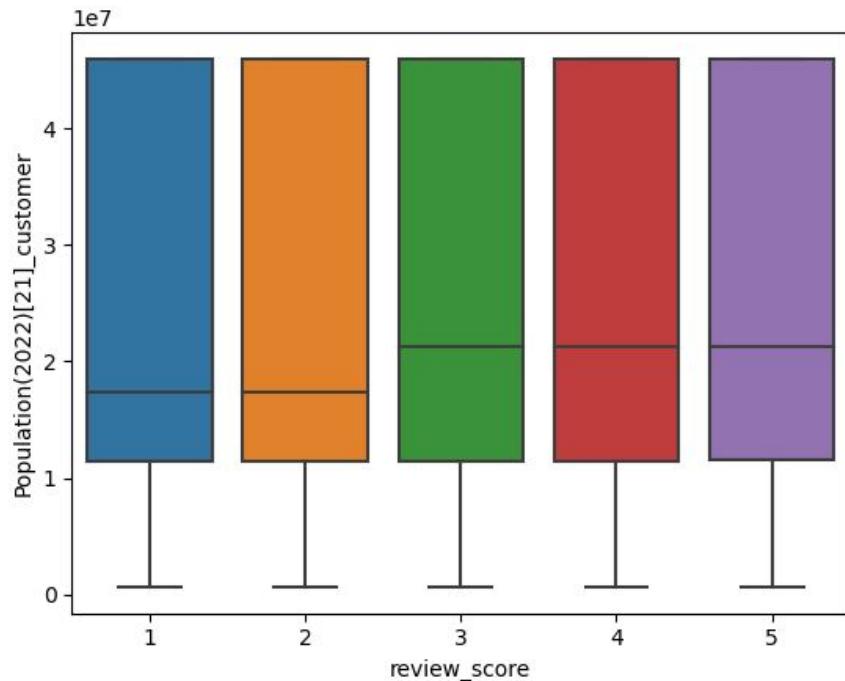
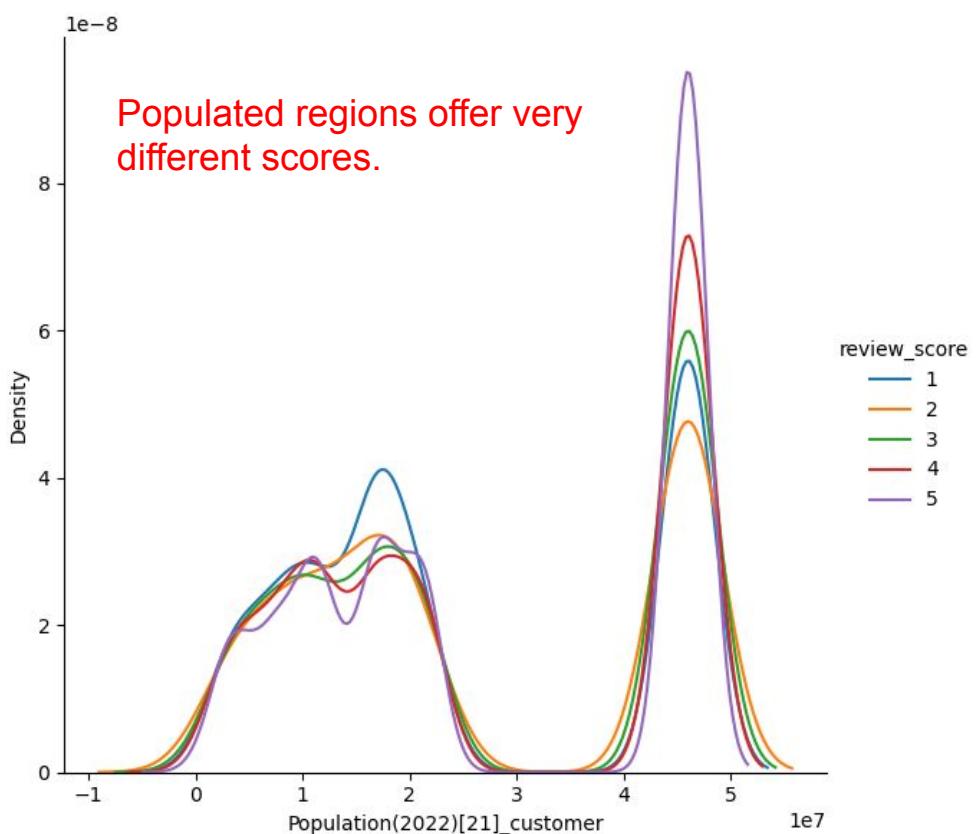
# Score by speed

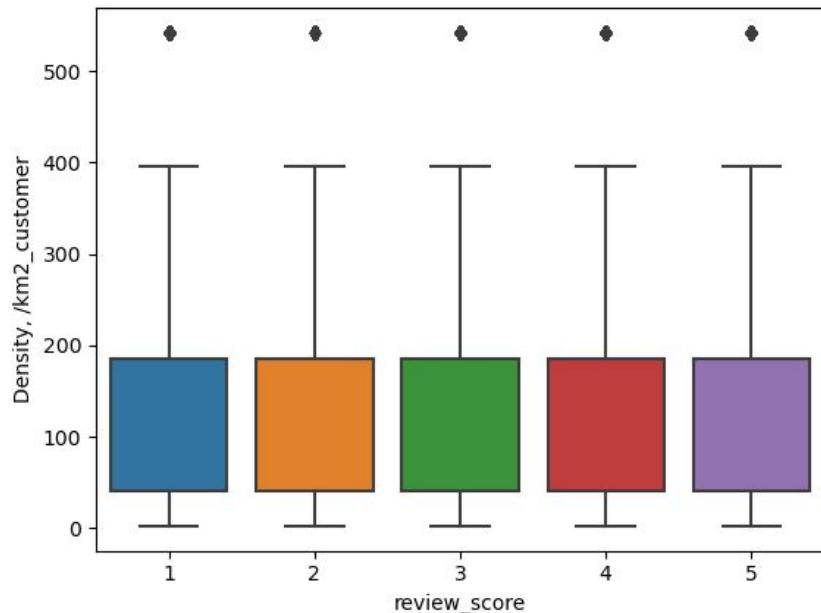
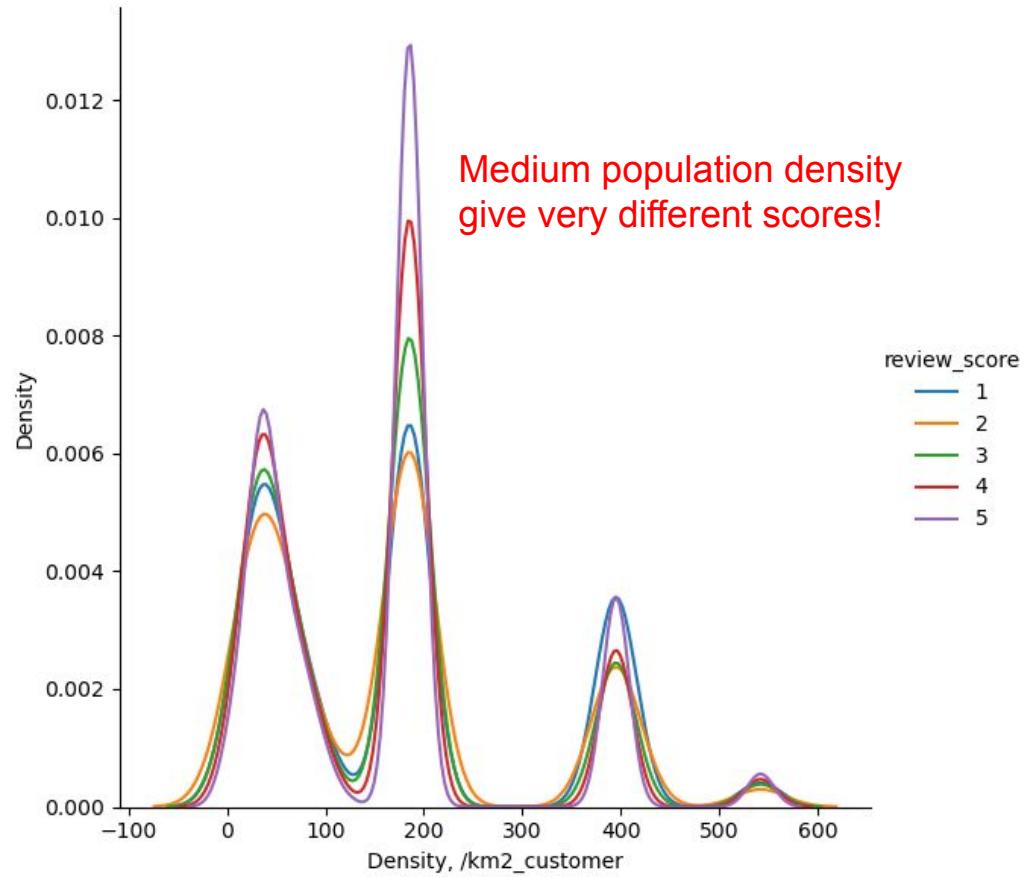


# Score by GDP



# Score by Population





# 5. Correlations: Geo

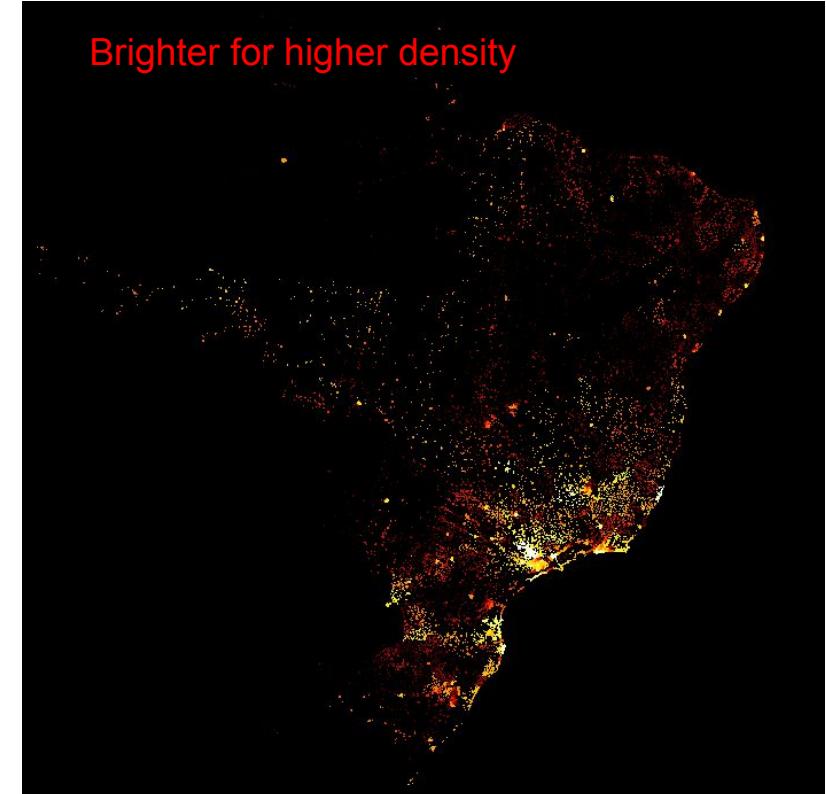
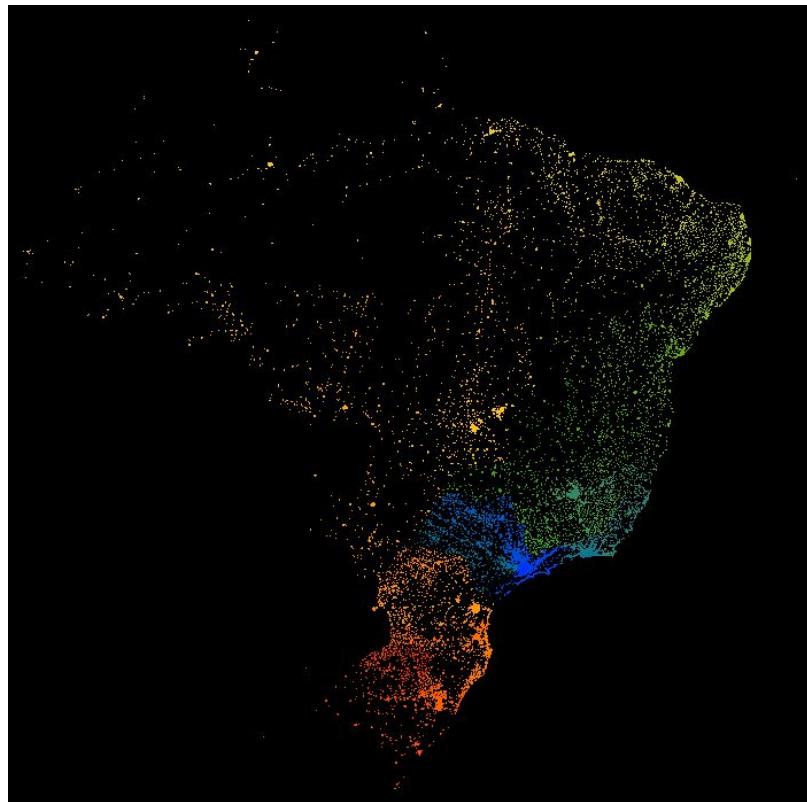
Zip Code and Revenue

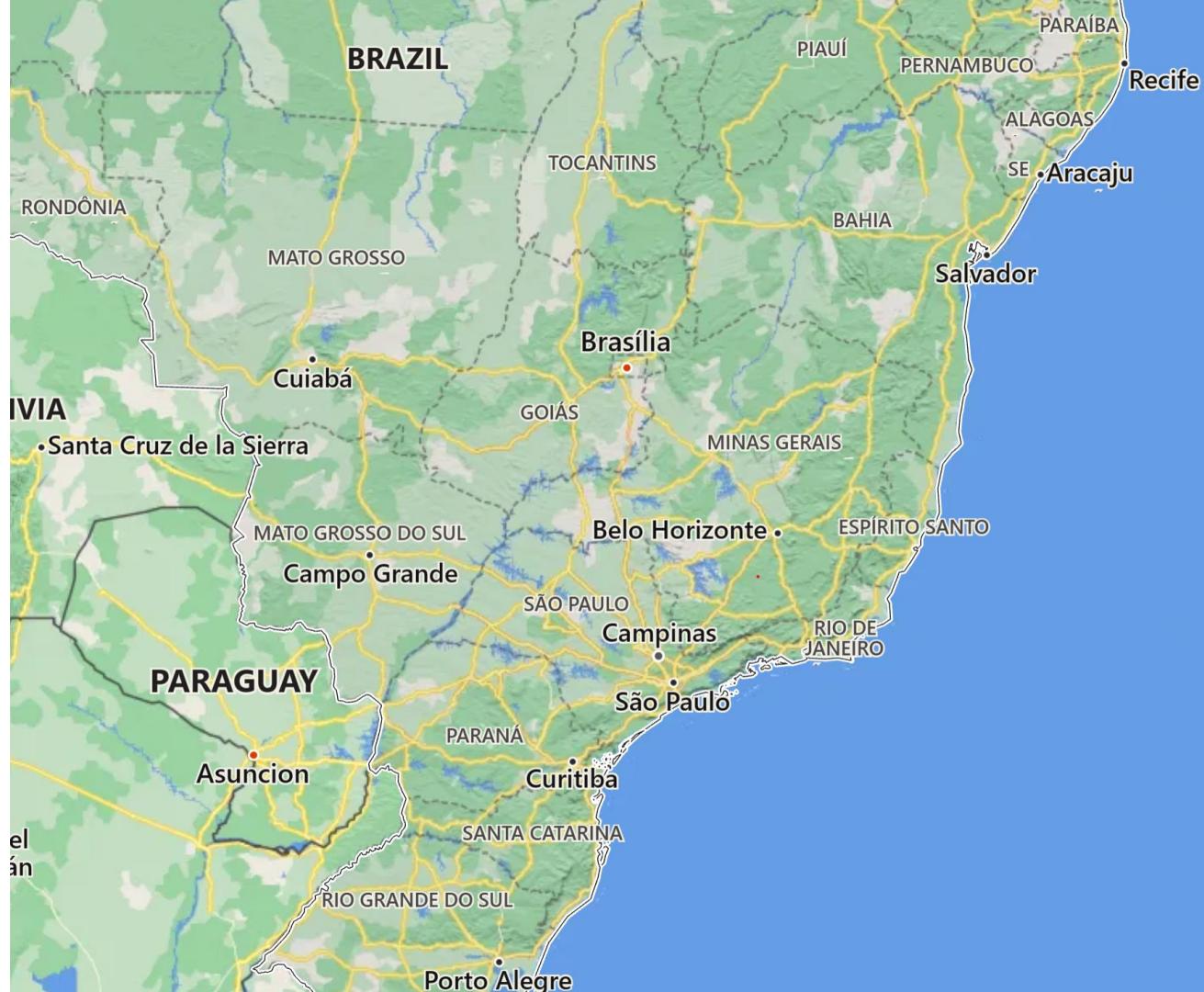
Average Ticket and Freight Ratio

Deliver Time and Delay

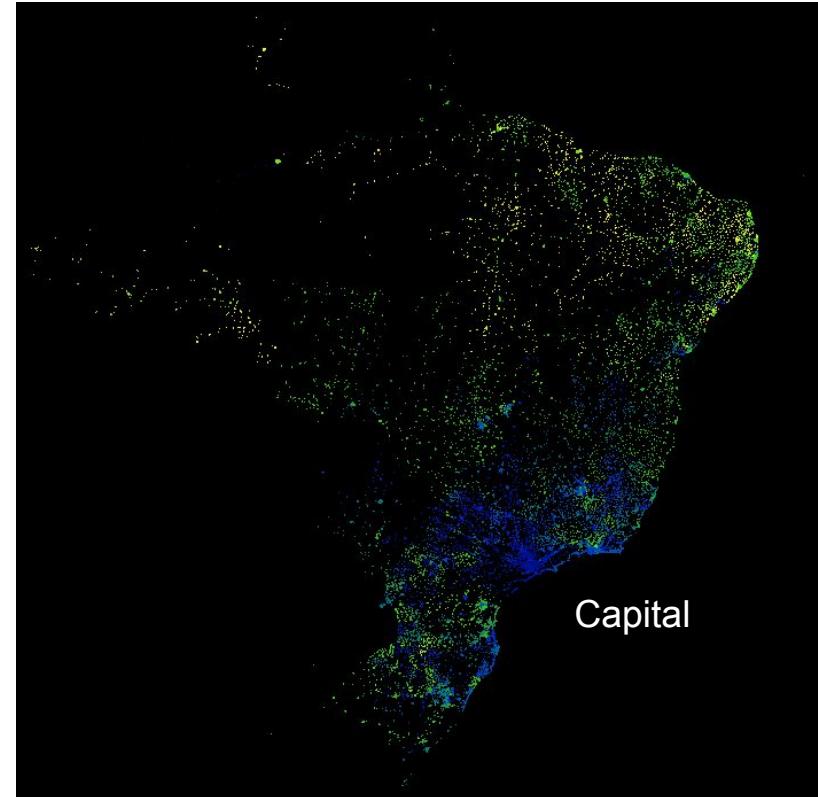
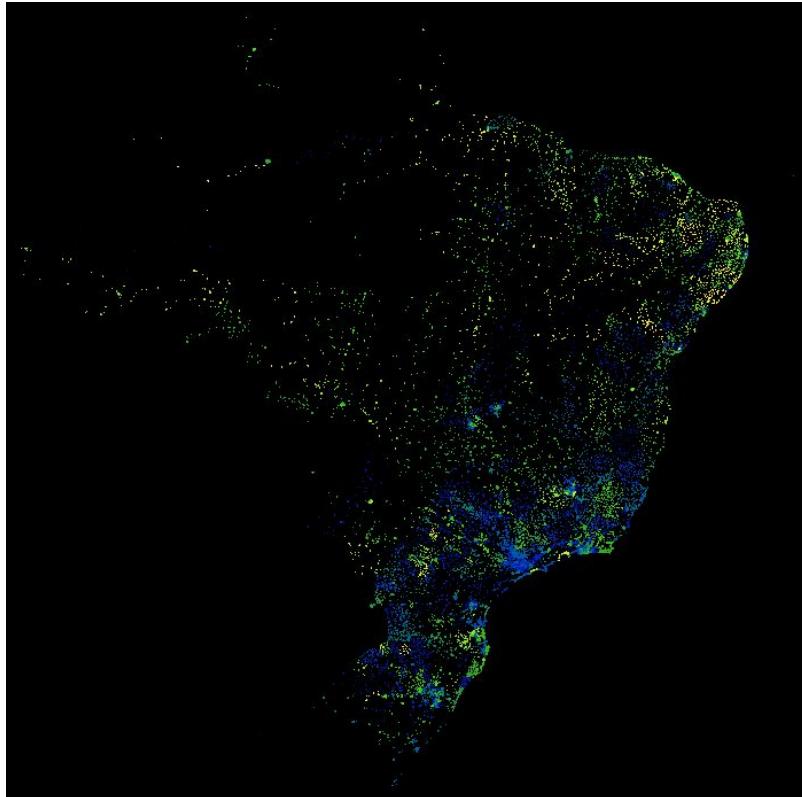
Average and Individual Review Score

# Zip code and Revenue

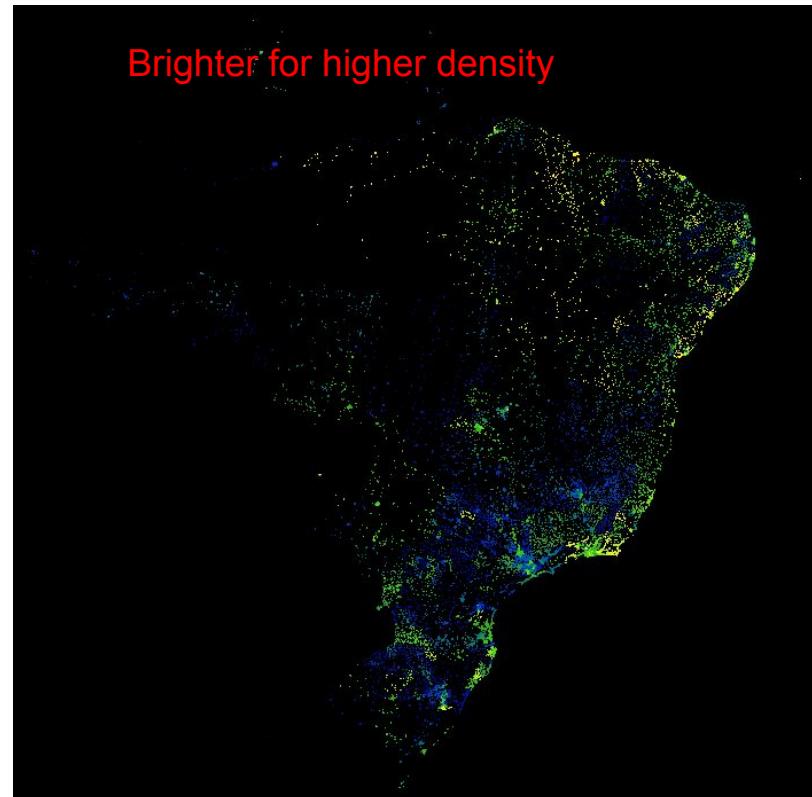
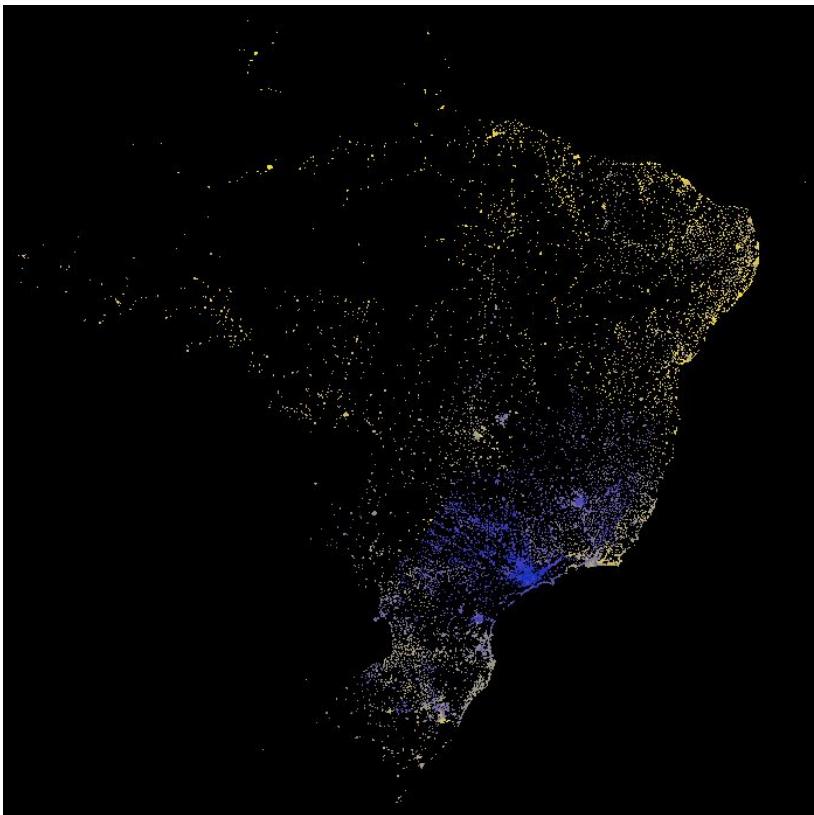




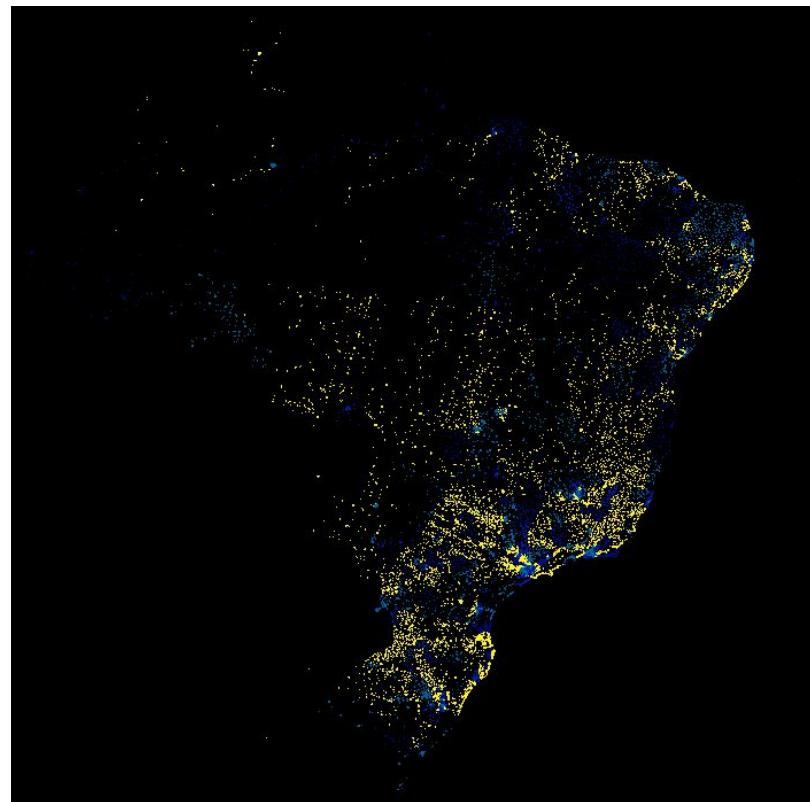
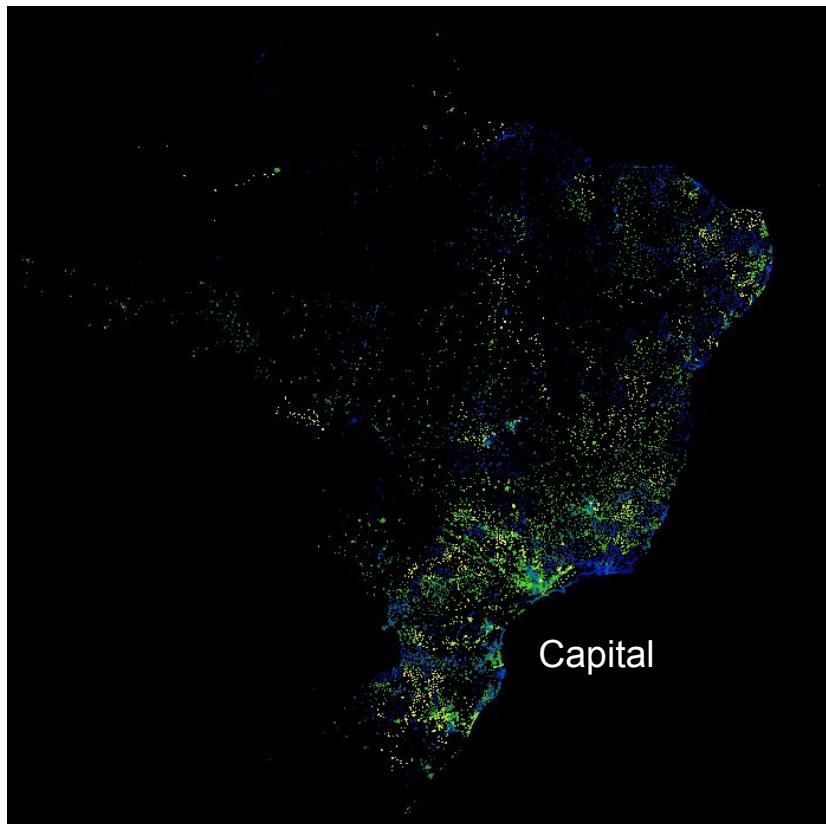
# Average Ticket and Freight Ratio



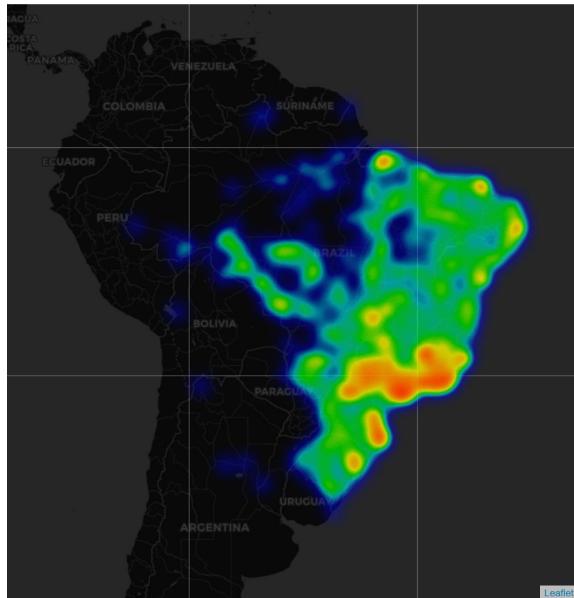
# Deliver Time and Delay



# Average and Individual Review Score



# Folium for time lapse and map

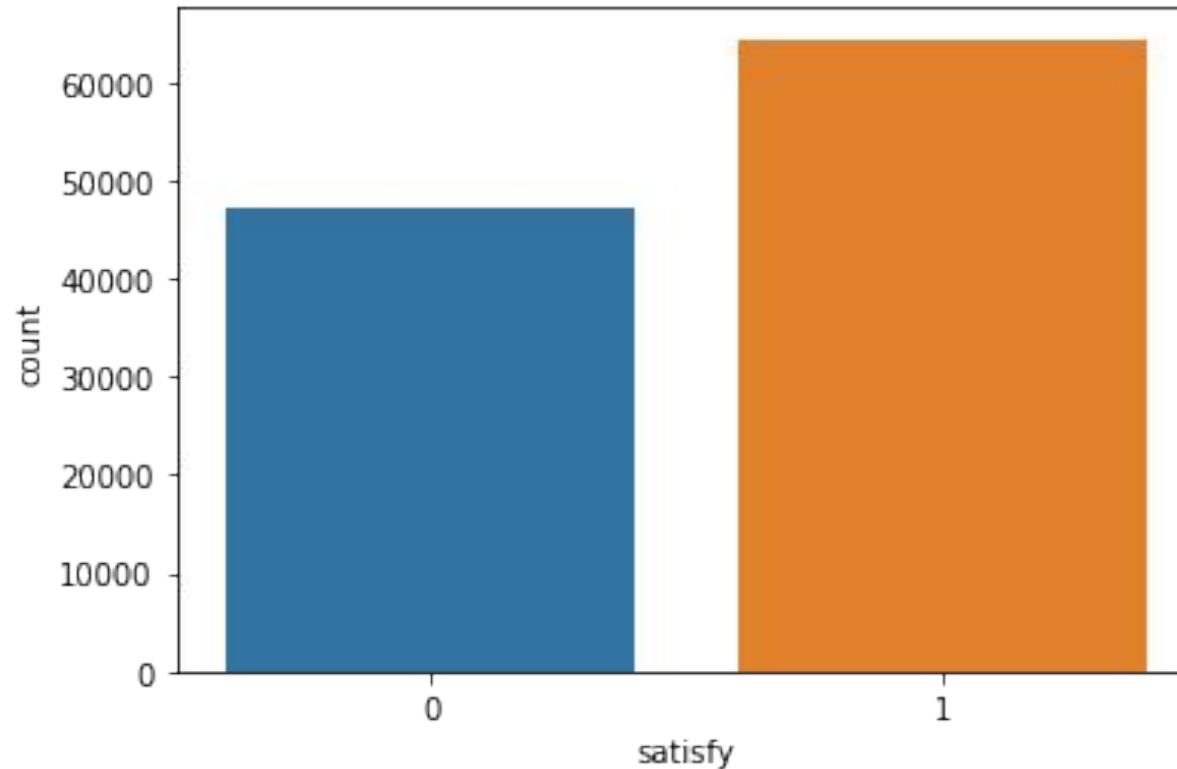


Youtube Link: [https://www.youtube.com/watch?v=0xSW-\\_N7-OA](https://www.youtube.com/watch?v=0xSW-_N7-OA)

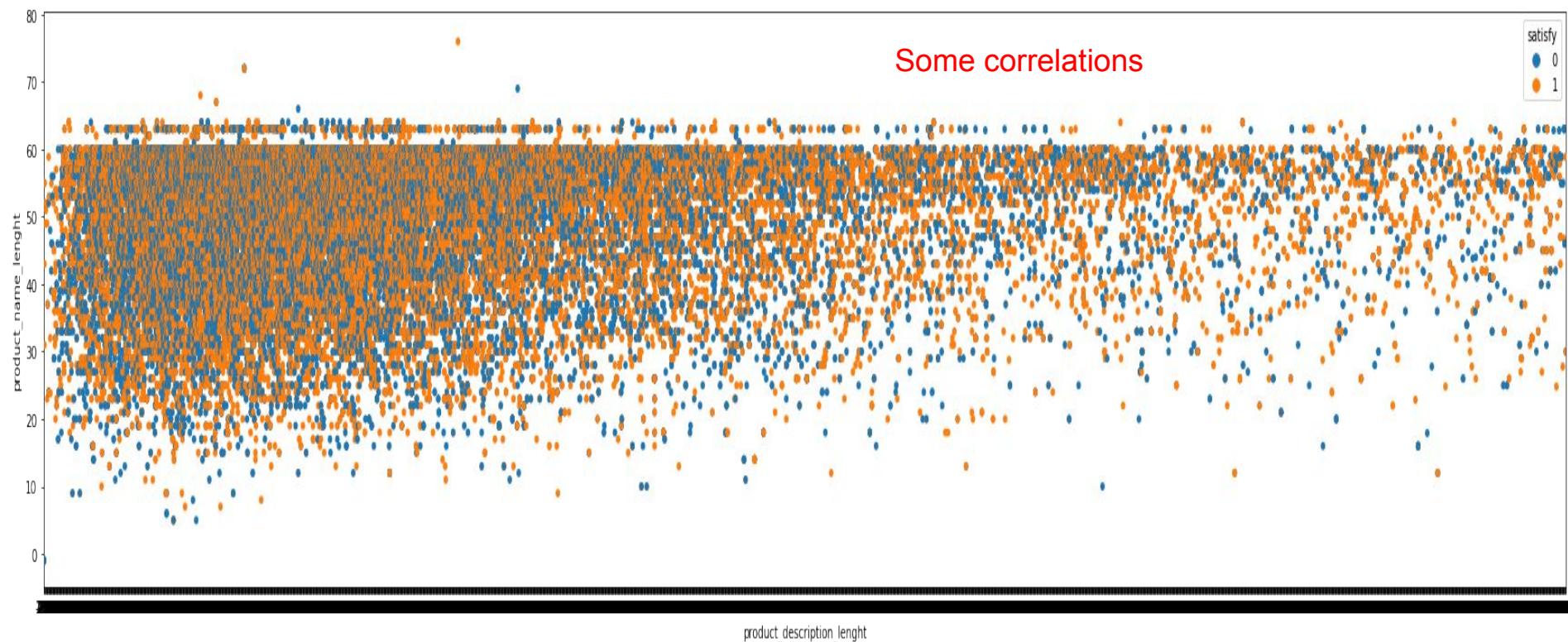
Reference: [E-Commerce Sentiment Analysis: EDA + Viz + NLP ↗ | Kaggle](#)

# 6. Threshold for Satisfaction

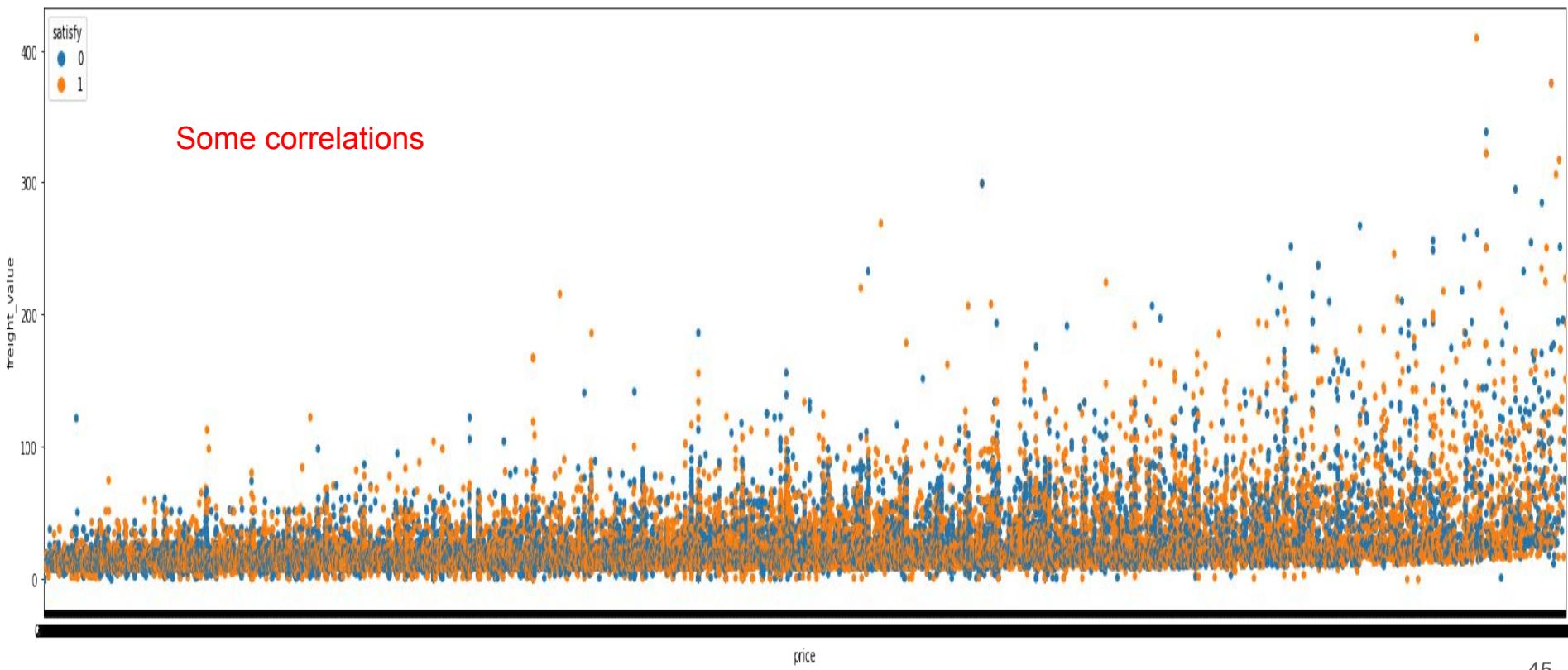
## Target Distribution (Threshold =4)



# Satisfaction by length of product name and description

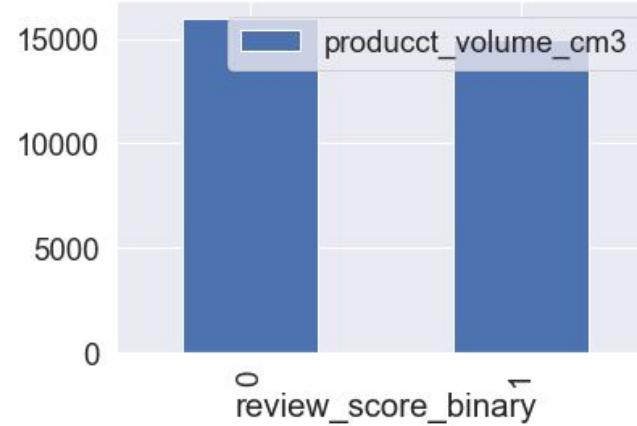
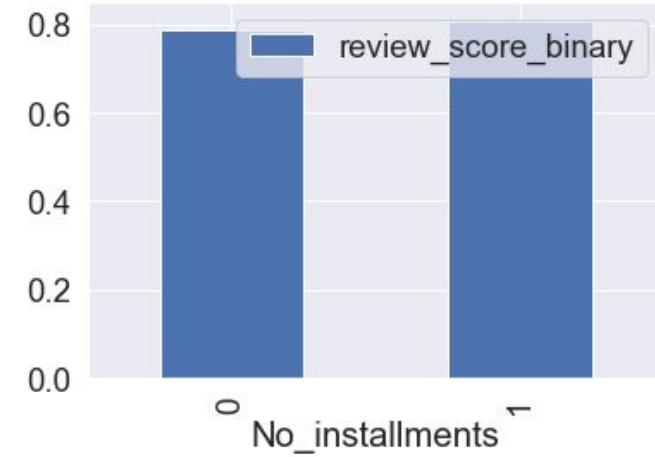


# Satisfaction by price and freight value

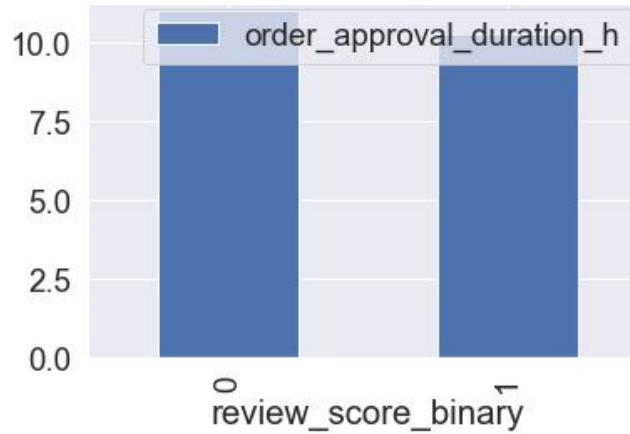
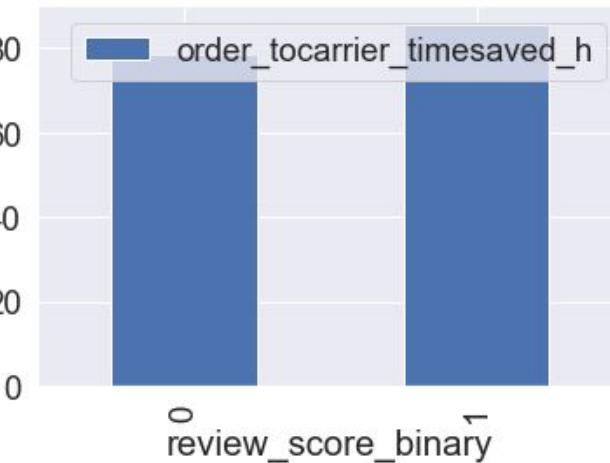
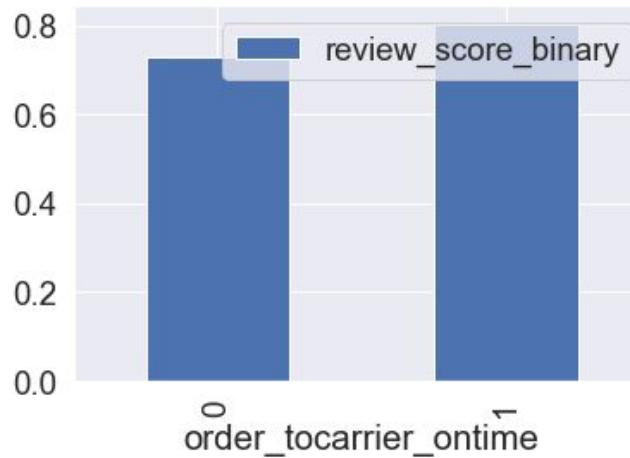
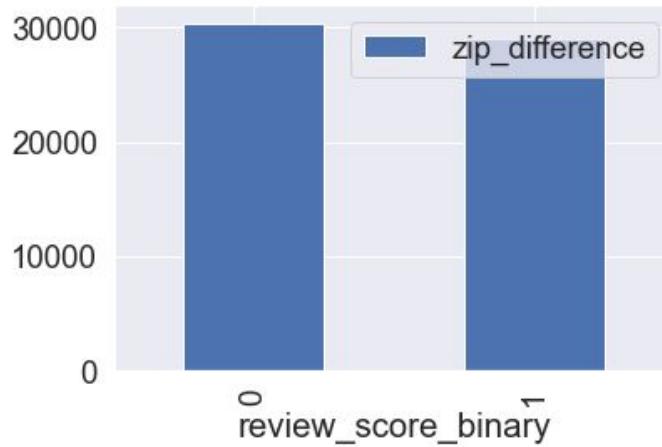




Satisfaction distribution  
for payment\_installment,  
product volume



## Even satisfaction distribution for other features



# Feature Selection

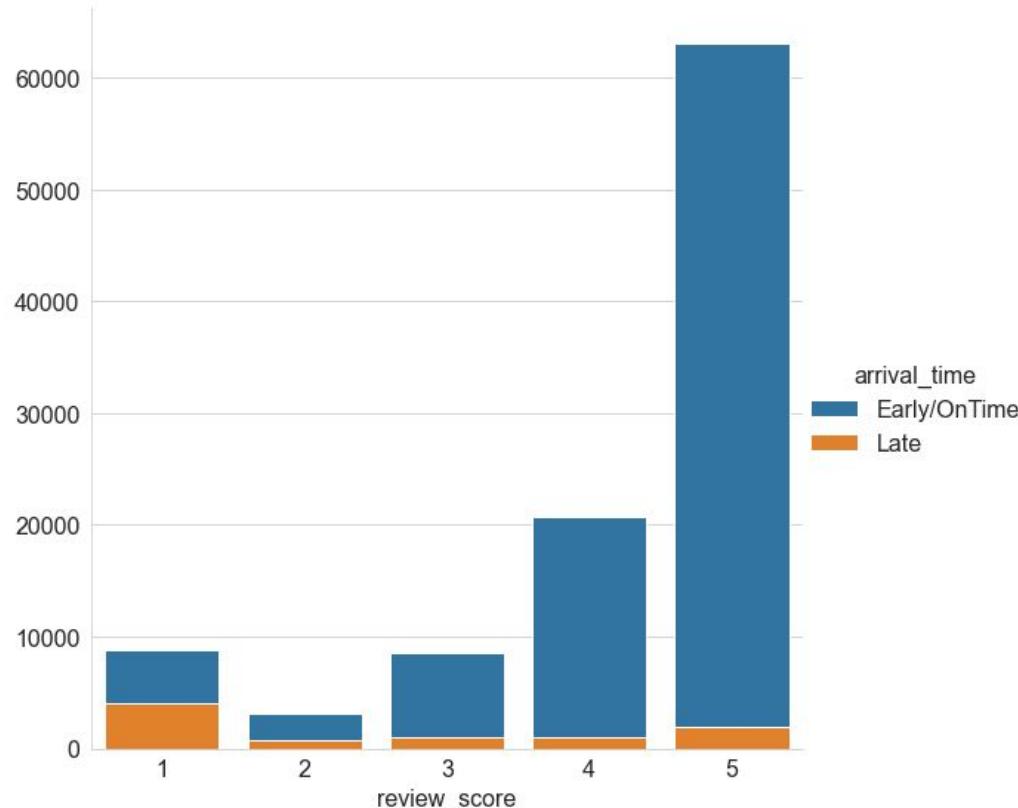
# Data Cleaning

1. Drop ID related feature
2. Make new features for time
3. Drop missing (110k to 90k)
4. CountVectorizer for categorical
5. Dropped outliers such as 60-day delivery
6. Drop timestamps after new time-difference feature generated

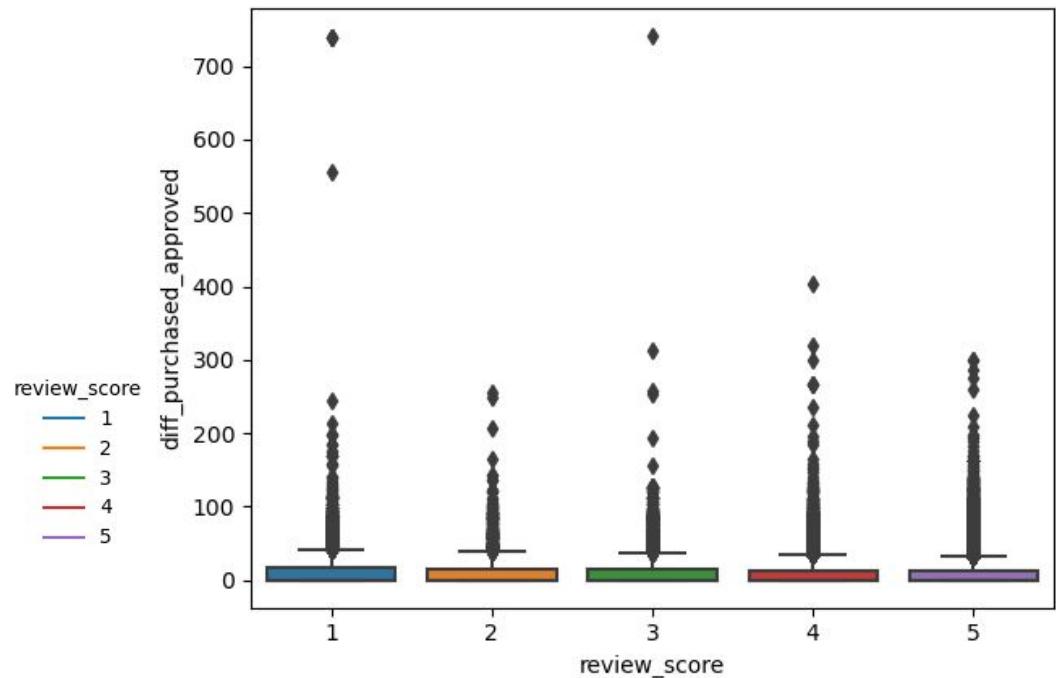
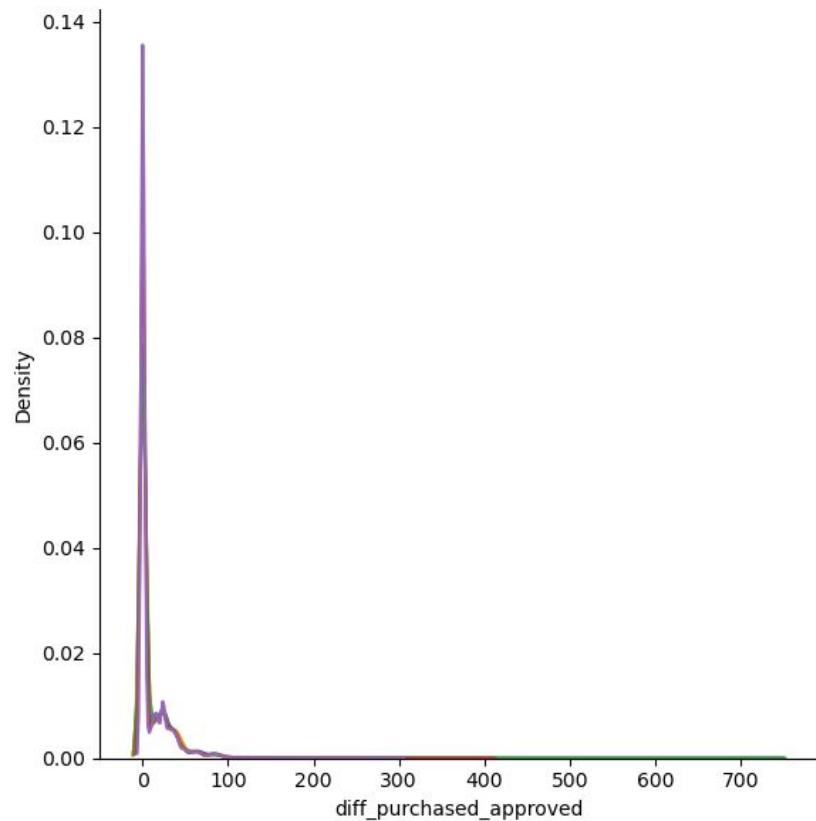
## *Other dropped features for high F1 models*

```
[ 'review_comment_title', 'review_comment_message', 'product_category_name',
  'product_weight_g', 'review_creation_date',
  'product_length_cm', 'product_height_cm', 'product_width_cm', 'seller_city',
  'review_answer_timestamp',
  'geolocation_lat_y', 'geolocation_lng_y', 'geolocation_city_y', 'geolocation_state_y',
  'review_id', 'order_approved_at', 'order_status',
  'order_id', 'customer_id', 'order_item_id', 'geolocation_lat_x',
  'geolocation_lng_x', 'geolocation_city_x', 'geolocation_state_x' ]
```

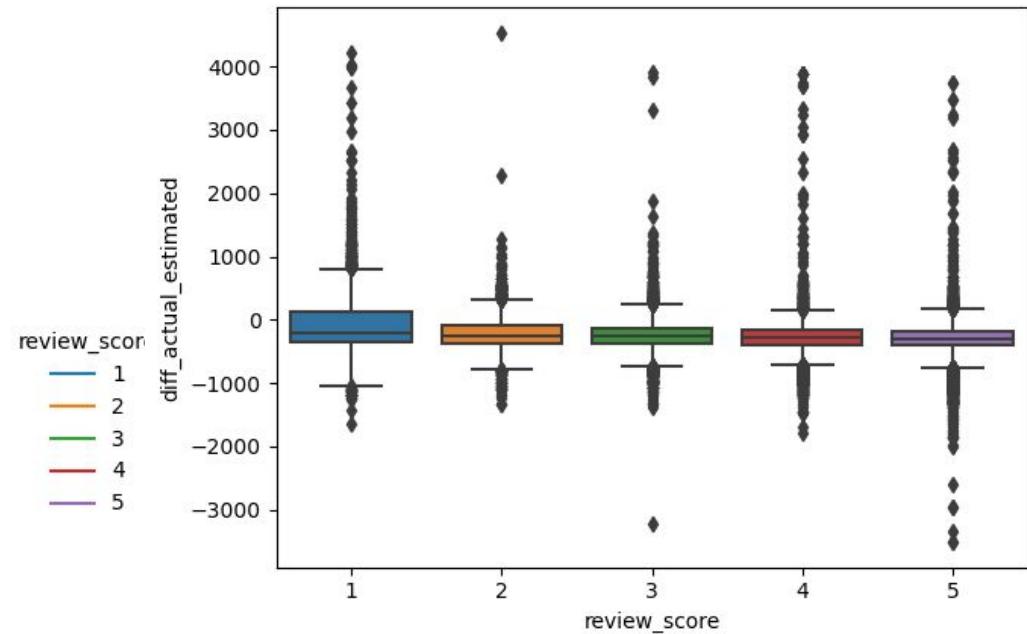
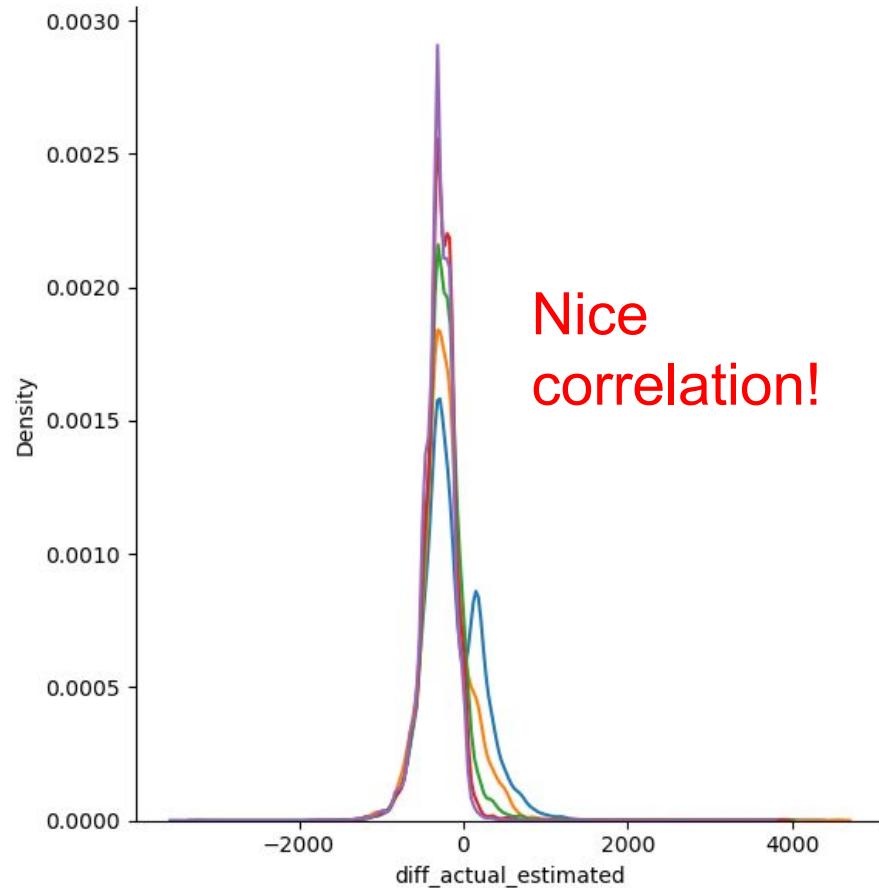
# New Time Feature: delivery time



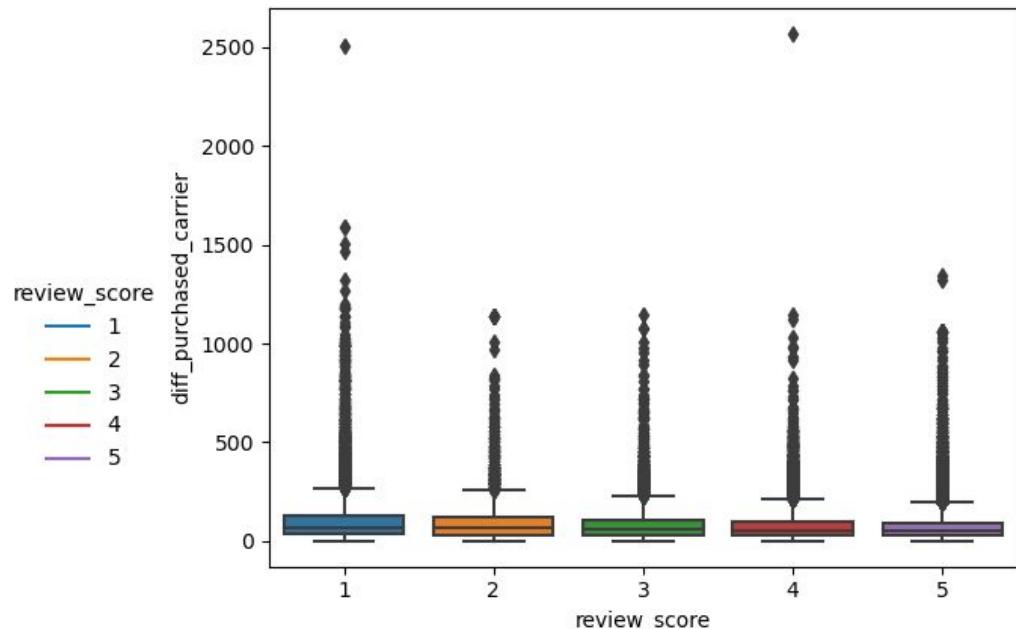
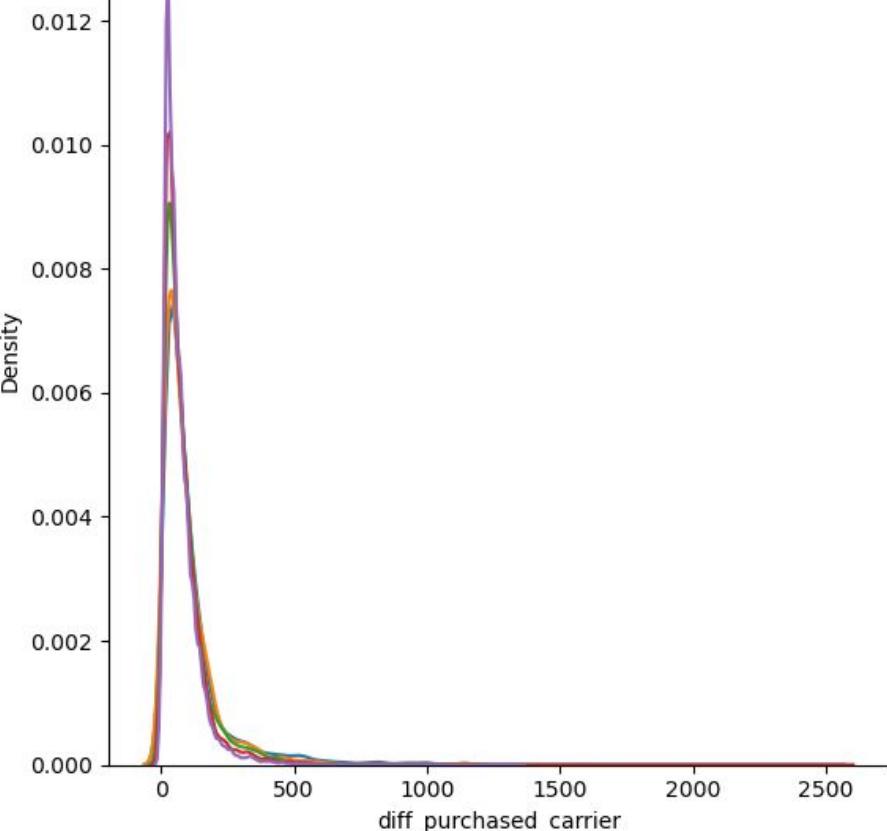
# New Time Feature 2



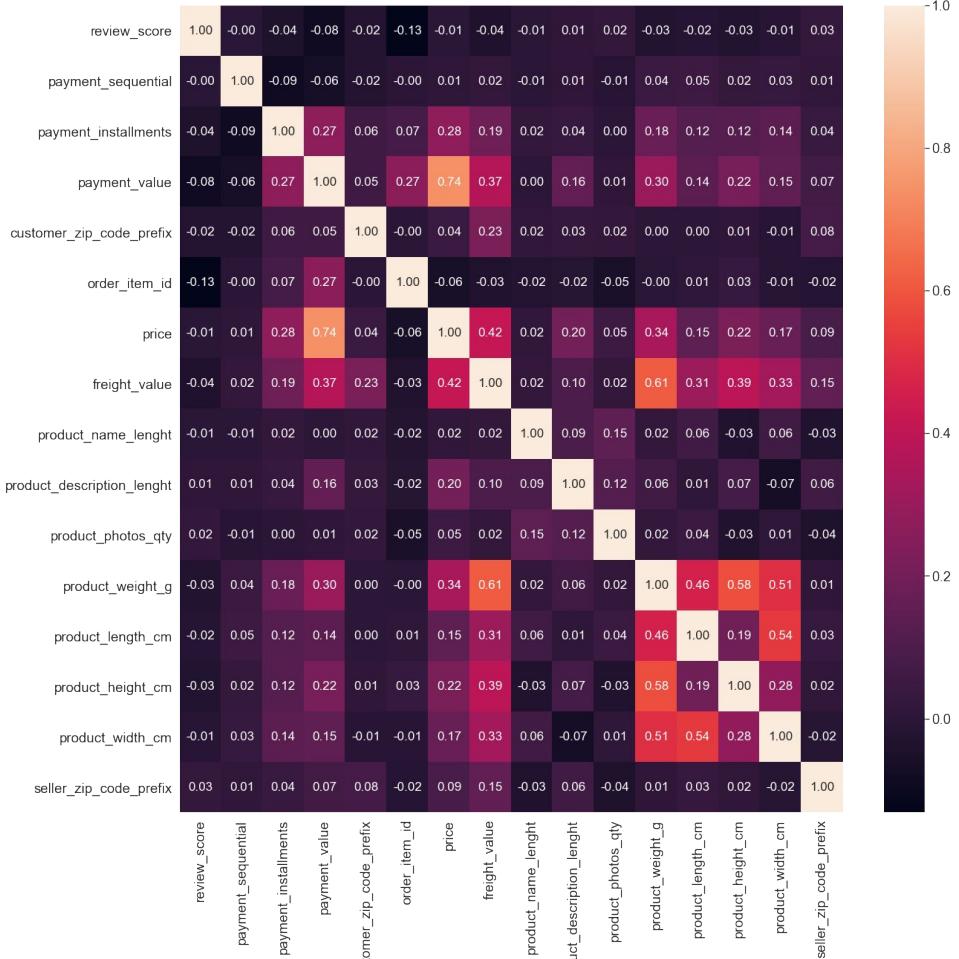
# New Time Feature 3



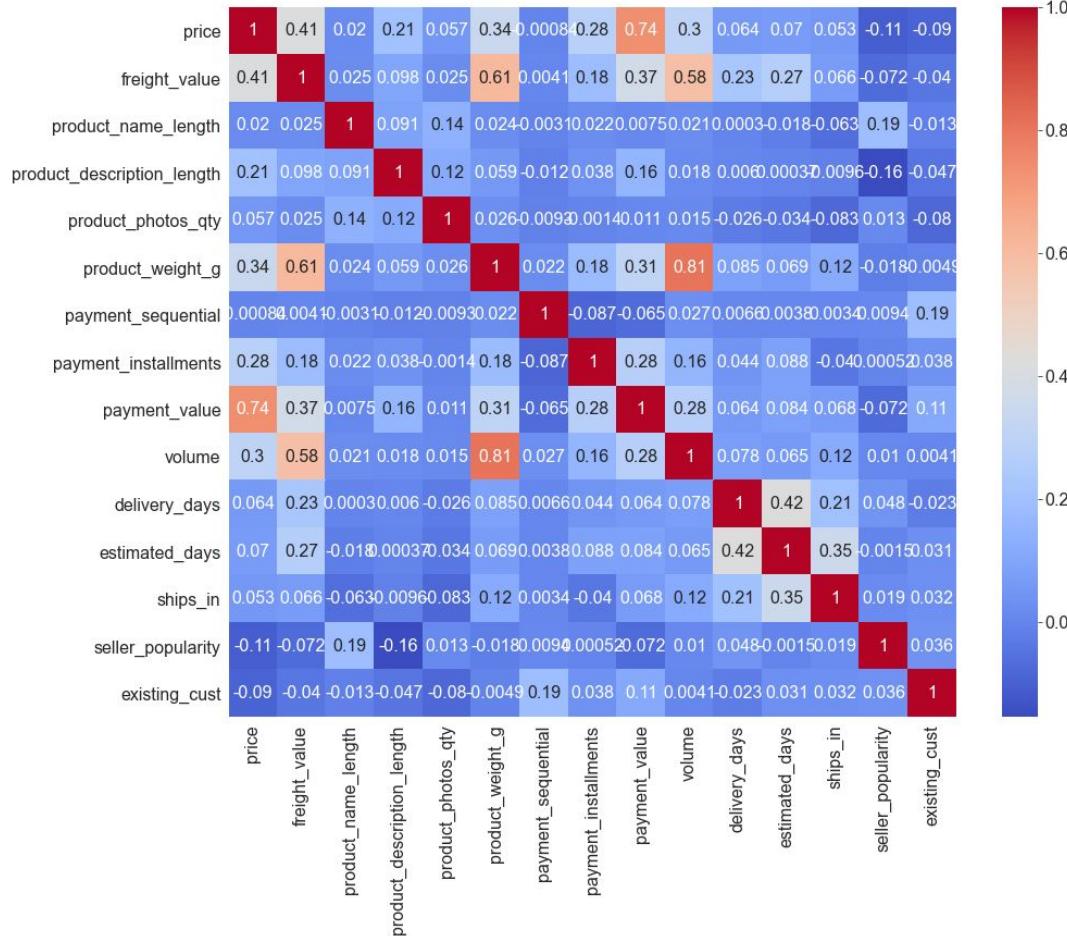
# New Time Feature 4



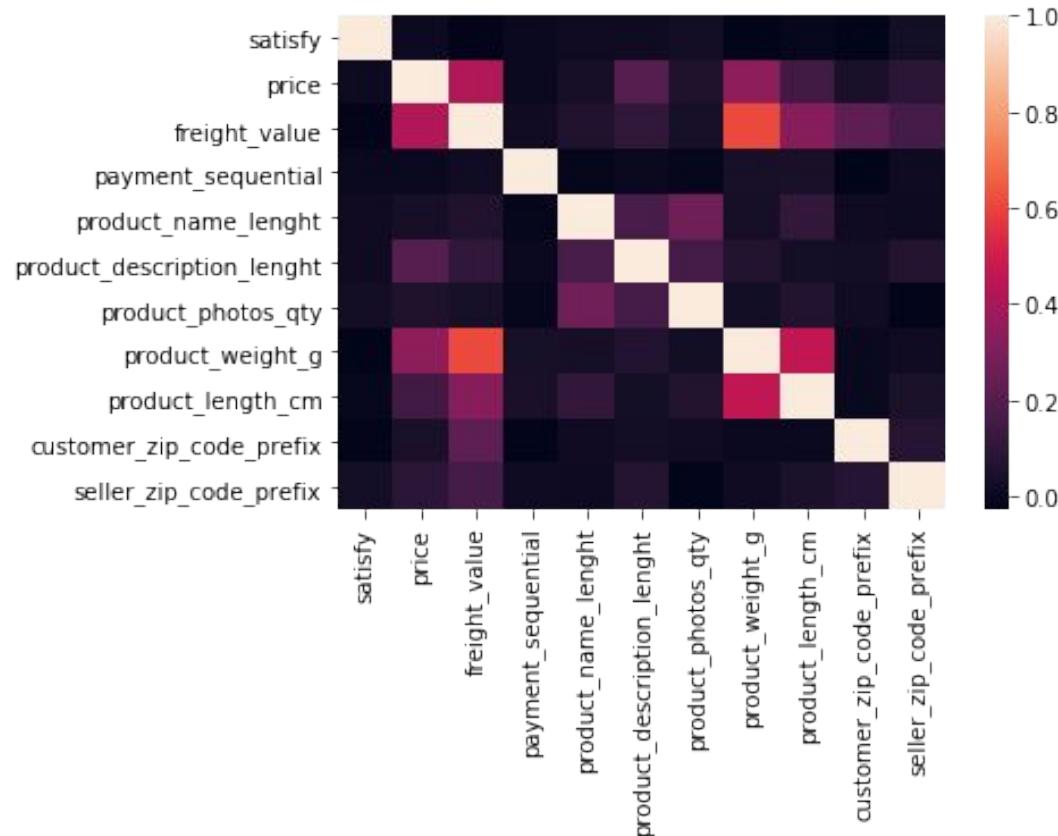
# Numerical Features before feature engineering (w/o target)



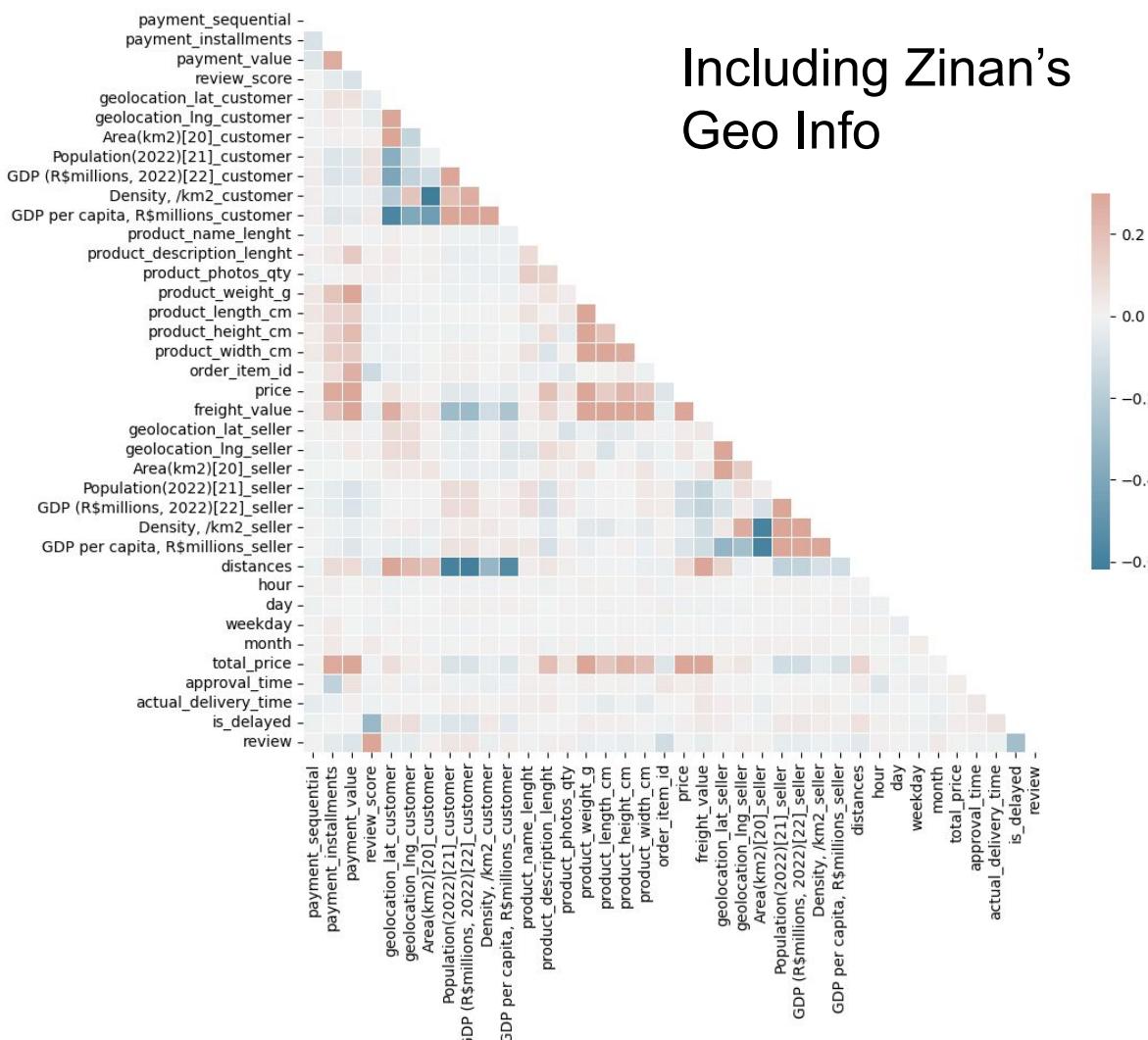
# after adding seller popularity and existing customer



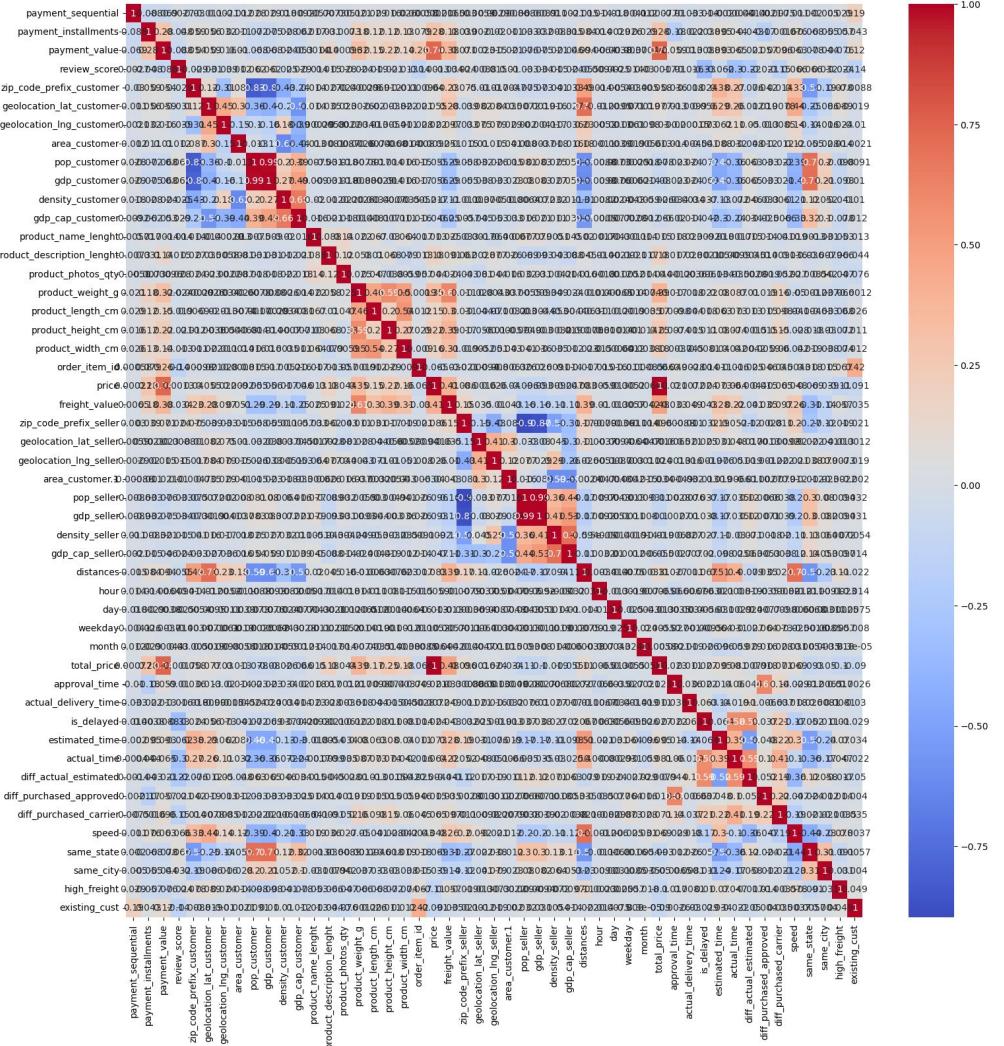
# Numerical Features before feature engineering (with target)



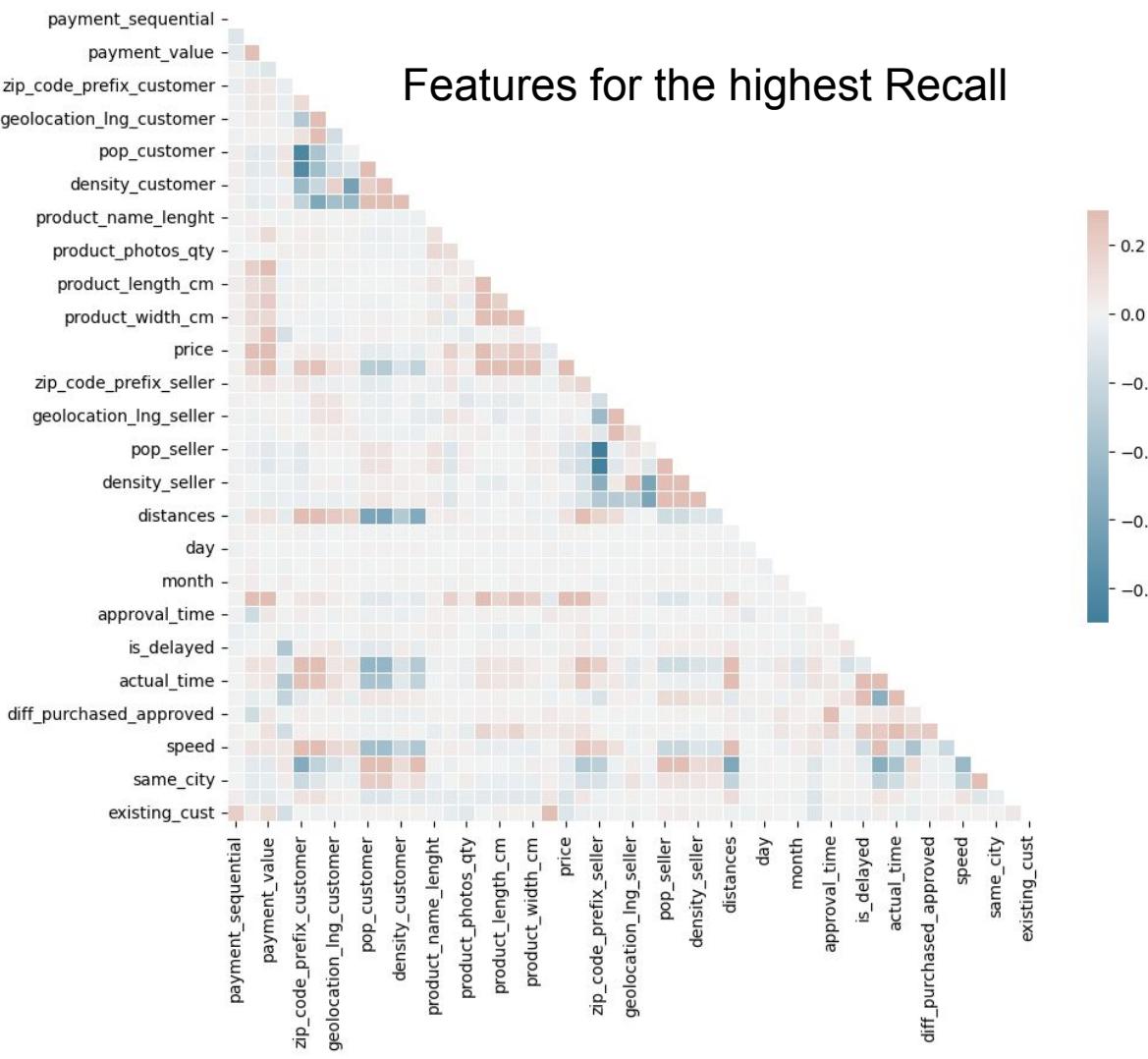
# Including Zinan's Geo Info



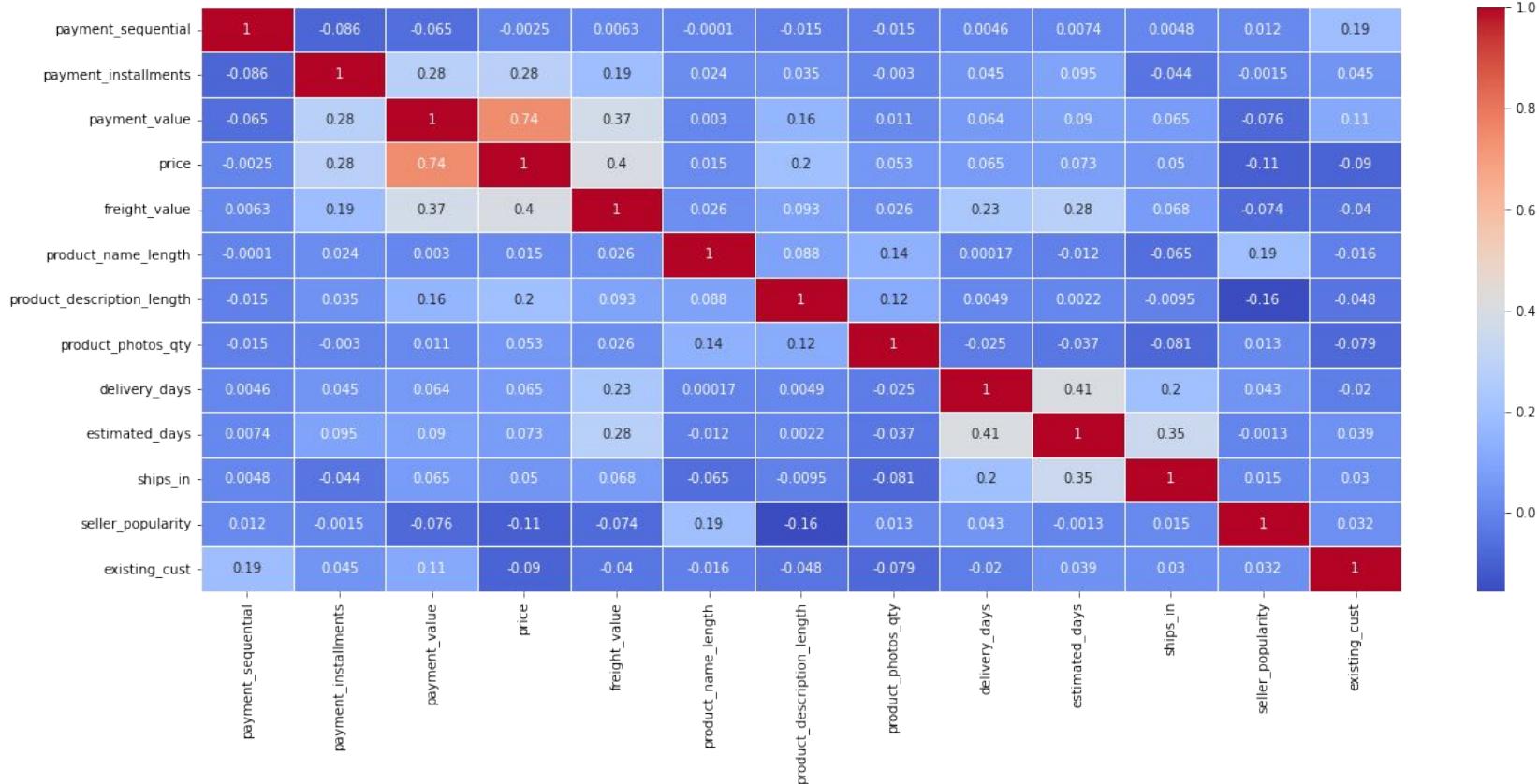
# Zinan's Mega Features



# Features for the highest Recall



## Features for the highest AUC and F1



Threshold	Metric	Logistic Regression	Naive Bayes	DecisionTree	Random Forest	GBoost	XGBoost
4.0	F1_score	0.85	0.78	0.77	0.81	0.86	0.85
	AUC_score	0.68	0.67	0.68	0.68	0.71	0.72
Dropped Payment Value, 4.5	F1_new	0.74	0.68	0.68	0.70	0.74	0.73
	AUC_new	0.63	0.62	0.63	0.63	0.65	0.66
SMOTE w/ Tomek Link, 4.0	F1_smt	0.77	0.78	0.79	0.82	0.83	0.84
	AUC_smt	0.68	0.67	0.67	0.68	0.70	0.68

	0	Naive Bayes	Logistic Regression	LogisticRegression (balanced)	DecisionTree	Random Forest	GBoost
0	AUC_score	0.653	0.663	0.663	0.669	0.674	0.696
1	F1_score	0.71	0.75	0.71	0.70	0.74	0.76

→

	0	Naive Bayes	Logistic Regression	LogisticRegression (balanced)	DecisionTree	Random Forest	GBoost	XGBoost
0	Recall (0)	0.57	0.22	0.57	0.51	0.55	0.23	0.22
1	F1_score	0.75	0.88	0.78	0.81	0.80	0.88	0.88
2	AUC_score	0.68	0.69	0.68	0.69	0.69	0.7	0.70

# Possible actions to improve satisfaction

1. Coupon
2. Refund
3. New item or exchange
4. Improve other aspects for future customers



Thanks!

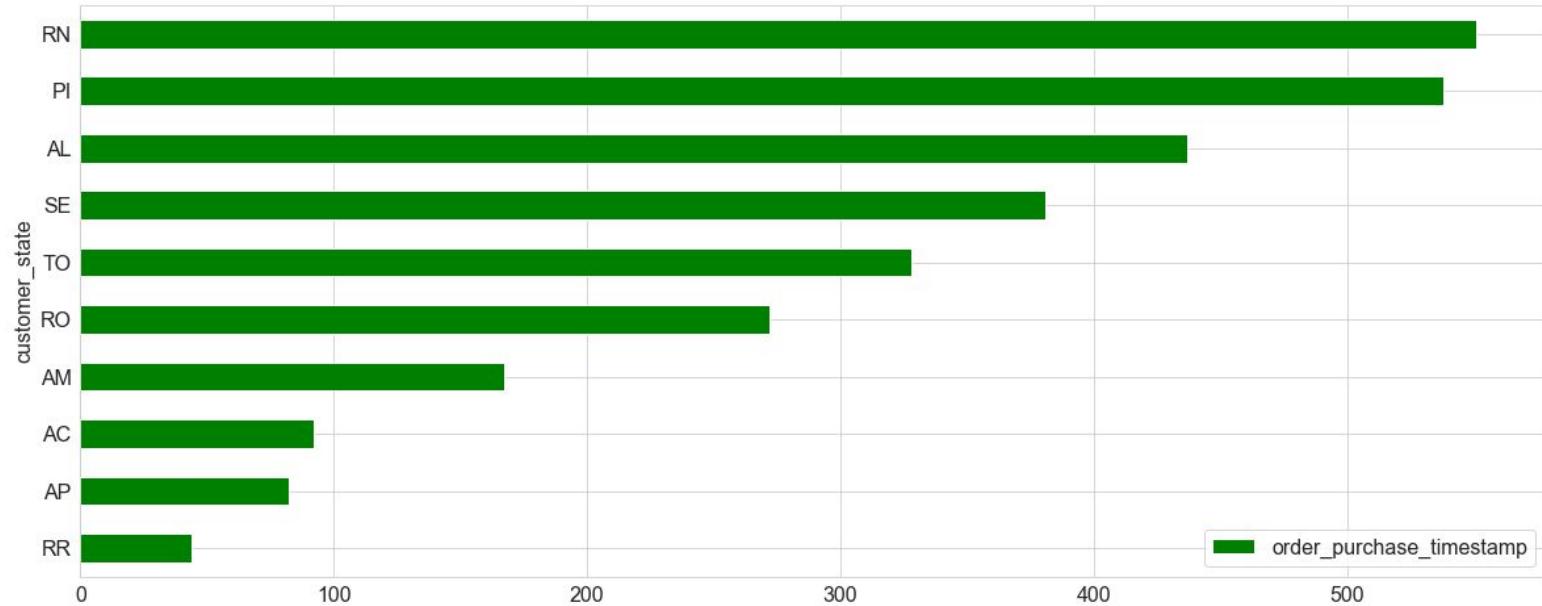


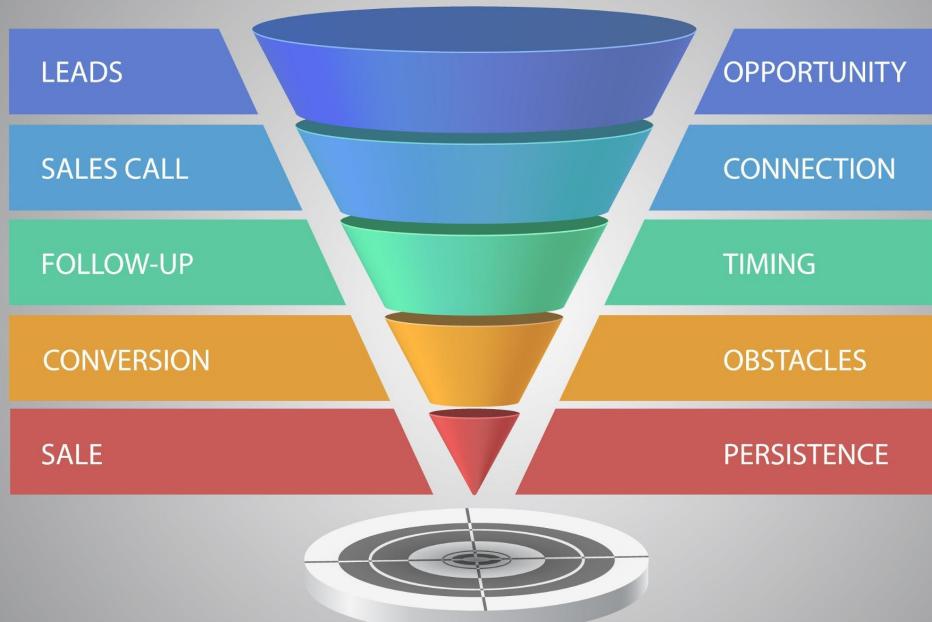




# backup

# Order time by state





<https://www.crazyegg.com/blog/ecommerce-conversion-funnel/>

MARKETING AND SALES FUNNEL