

eCommerce Data Analysis - Profit Improvement Suggestion

Introduction

Background

- This data analysis project was introduced during the coding bootcamp.
- It was a half-guided module and the problem and the preconditions are presented.
- The insights, the visualizations and the suggestions are my original.

Dataset

- Brazilian E-Commerce Public Dataset by Olist
- kaggle.com/olistbr/brazilian-ecommerce
- Olist offers Logistic and Inventory Management Service to sellers
- Information about ~100k orders made between 2016 and 2018
- 8 csv files (orders, customers, reviews, sellers, products...etc)

Problem statement

How should Olist improve its profit margin, given that the revenue and the cost are calculated as the following condition:

Revenue

- Olist takes a 10% cut on the product price (excl. freight) of each order delivered.
- Olist charges 80 BRL by month per seller.

Cost

- Estimated cost occurred by bad review per order

review_score	cost (BRL)
1 star	100
2 stars	50
3 stars	40
4 stars	0
5 stars	0

- IT costs

Olist's IT costs are estimated to be proportional to the square-root of the total cumulated number of orders approved.

The IT department also told you that since the birth of the marketplace, cumulated IT costs have amounted to 500,000 BRL.

Suggestion

After the data analysis on seller data, I would give the following 2 suggestions to improve the profits by at least 10 %.

- By removing the worst 15 sellers (0.5 % of the sellers) who make negative profits, Olist improves the profits by 10 %
- By charging 10 % of the review cost directly to the sellers, Olist improves the profits by 13.9 %

EDA of seller dataset

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import matplotlib.patches as mpatches
import seaborn as sns
import math
from scipy import stats
import plotly.express as px

Original data is imported, cleaned and processed then according to 2,3, revenues and profits for each seller are calculated and the code is saved in seller_updated.py

In [2]: from olist.seller_updated import Seller as SellerUpdated
sellers = SellerUpdated().get_training_data()

In [3]: sellers.shape

Out[3]: (2967, 18)

In [4]: sellers.head()

Out[4]:
```

	seller_id	seller_city	seller_state	delay_to_carrier	wait_time	date_first_sale	date_last_sale	months_...
0	34478959a846ae7ee197c632cb2df15	campinas	SP	0.000000	13.018588	2017-05-05 16:25:18	2017-08-30 12:50:19	
1	d1b65fcd7ebc3361ea8b65f14c68d262	moji guacu	SJ	0.000000	9.065716	2017-03-29 02:10:34	2018-06-06 20:15:21	
2	ce3a9d9e960102a0677a8f1f5d0bb702d	rio de janeiro	RJ	0.000000	4.042292	2018-07-30 12:44:49	2018-07-30 12:44:49	
3	c0f3eeae2e1455b6f6aea3dd58c1b1c3	sao paulo	SP	0.000000	5.667187	2018-08-03 00:44:08	2018-08-03 00:44:08	
4	51a04a8a6bdc232decc082b0c80742cf	braganca paulista	SP	3.353727	35.314861	2017-11-14 12:15:25	2017-11-14 12:15:25	

```
In [5]: sellers.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 2967 entries, 0 to 2966
Data columns (total 18 columns):
# Column Non-Null Count Dtype
0 seller_id 2967 non-null object
1 seller_city 2967 non-null object
2 seller_state 2967 non-null object
3 delay_to_carrier 2967 non-null float64
4 wait_time 2967 non-null float64
5 date_first_sale 2967 non-null datetime64[ns]
6 date_last_sale 2967 non-null datetime64[ns]
7 months_on_olist 2967 non-null int64
8 share_of_one_stars 2967 non-null float64
9 share_of_five_stars 2967 non-null float64
10 review_score 2967 non-null float64
11 review_costs 2967 non-null int64
12 n_orders 2967 non-null int64
13 quantity 2967 non-null int64
14 quantity_per_order 2967 non-null float64
15 sales 2967 non-null float64
16 revenues 2967 non-null float64
17 profits 2967 non-null float64
dtypes: datetime64[ns](2), float64(10), int64(3), object(3)
memory usage: 440.4+ KB

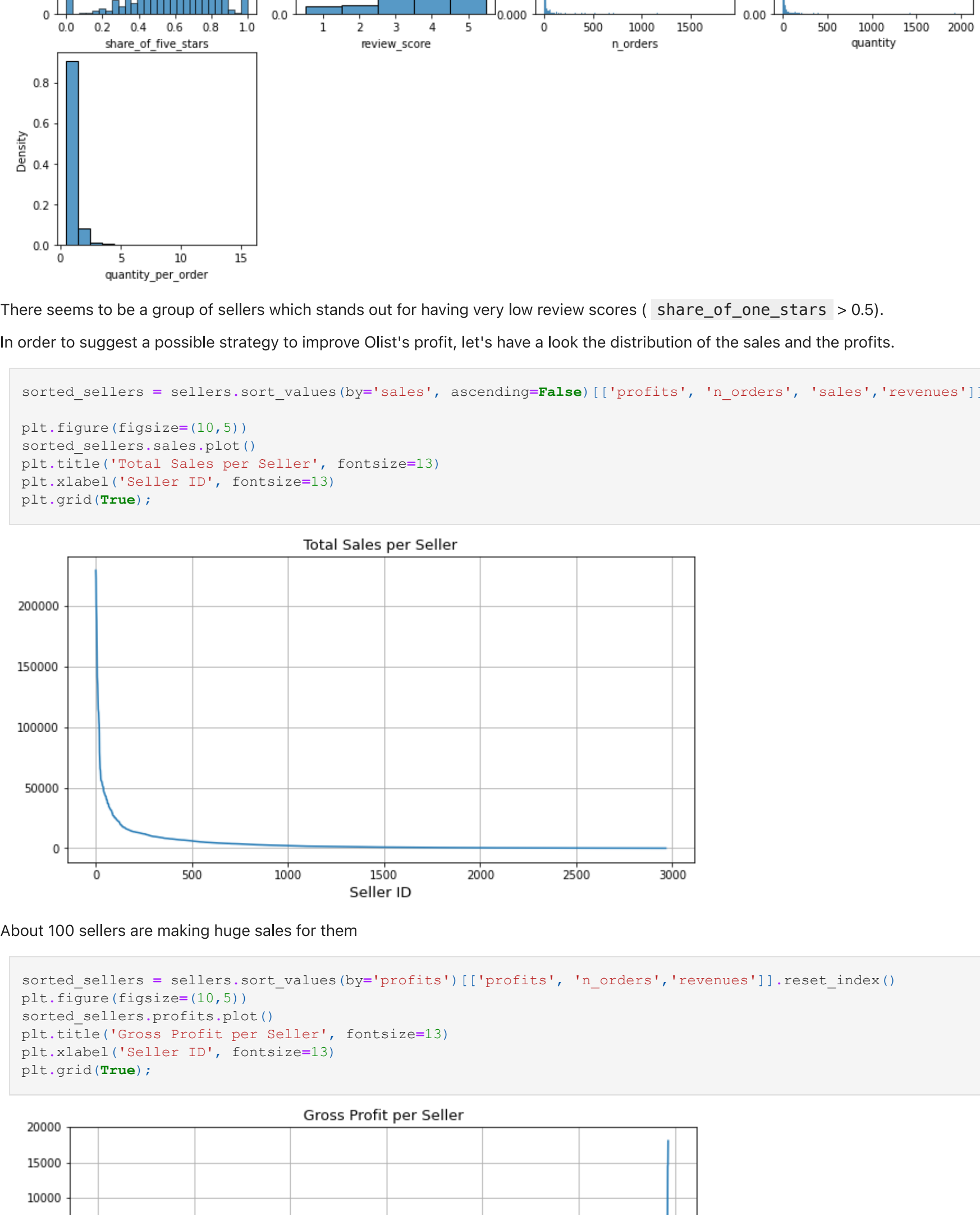
In [6]: sellers.describe()

Out[6]:
```

	delay_to_carrier	wait_time	months_on_olist	share_of_one_stars	share_of_five_stars	review_score	review_costs	n_o...
count	2967.000000	2967.000000	2967.000000	2967.000000	2967.000000	2967.000000	2967.000000	2967.00...
mean	0.385636	12.139417	6.025615	0.120886	0.595077	4.101513	545.402764	3.65...
std	2.259512	7.069233	5.994211	0.188857	0.278763	0.801517	1889.147815	1071...
min	0.000000	1.214718	0.000000	0.000000	0.000000	1.000000	0.000000	1.00...
25%	0.000000	8.287658	1.000000	0.000000	0.493022	3.846154	0.000000	2.00...
50%	0.000000	11.115143	4.000000	0.058824	0.600000	4.210526	100.000000	73.00...
75%	0.000000	14.231984	10.000000	0.159043	0.750000	4.626453	365.000000	230.00...
max	45.434039	189.863160	23.000000	1.000000	1.000000	5.000000	39400.000000	1854.00...

```
In [7]: sellers_eda = sellers.drop(['review_costs', 'sales', 'revenues', 'profits'], axis=1)
plt.figure(figsize=(15,11))
for (i, col) in enumerate(sellers_eda.describe().columns): #["wait_time", "delay_to_carrier", "avg_review_score"]
    plt.subplot(3,4,i+1)
    sns.histplot(sellers[col], kde=False, stat='density', discrete=True, None)[col in ['share_of_one_stars', 'share_of_five_stars']]

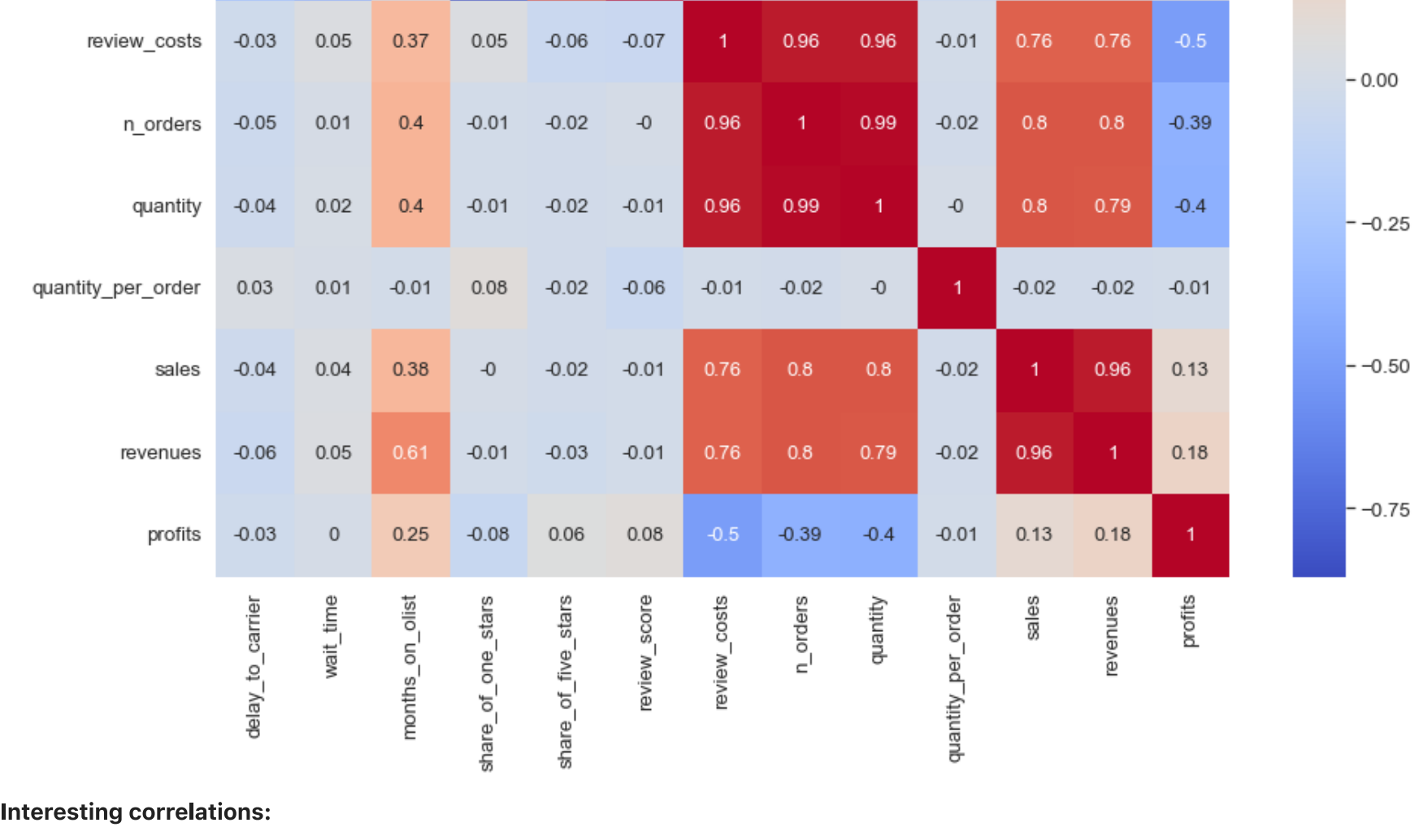
There seems to be a group of sellers which stands out for having very low review scores ( share_of_one_stars > 0.5).
In order to suggest a possible strategy to improve Olist's profit, let's have a look the distribution of the sales and the profits.
```



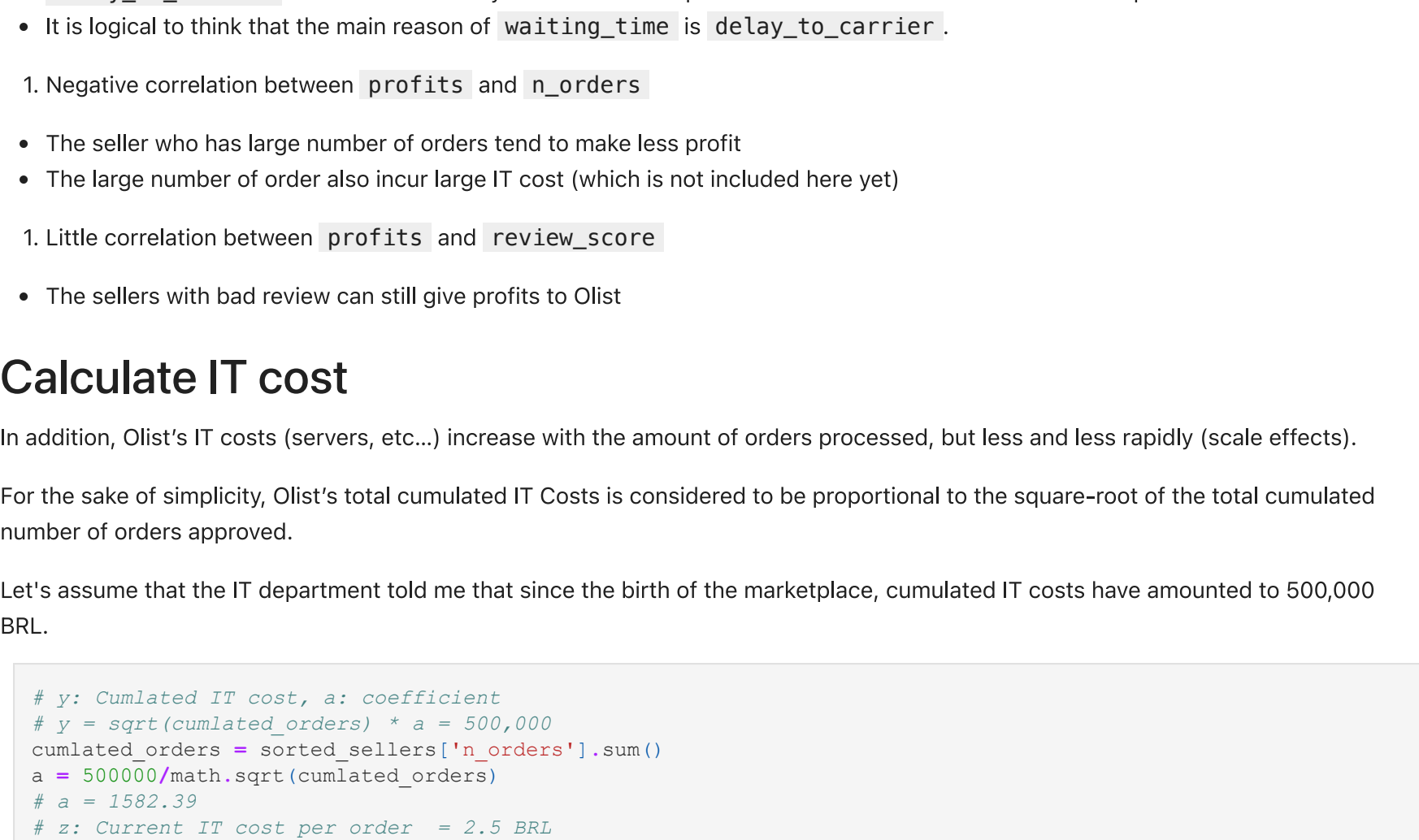
There seems to be a group of sellers which stands out for having very low review scores (share_of_one_stars > 0.5). In order to suggest a possible strategy to improve Olist's profit, let's have a look the distribution of the sales and the profits.



About 100 sellers are making huge sales for them



Only a few percentage of sellers make large negative/positive profits to Olist



Interesting correlations:

- Positive correlation between **waiting_time** and **delay_to_carrier**
 - delay_to_carrier** is the number of days when the seller provide the item to the carrier later than expected
 - It is logical to think that the main reason of **waiting_time** is **delay_to_carrier**.
- Negative correlation between **profits** and **n_orders**
 - The seller who has large number of orders tend to make less profit
 - The large number of order also incur large IT cost (which is not included here yet)
- Little correlation between **profits** and **review_score**
 - The sellers with bad review can still give profits to Olist

Calculate IT cost

In addition, Olist's IT costs (servers, etc...) increase with the amount of orders processed, but less and less rapidly (scale effects).

For the sake of simplicity, Olist's total cumulated IT Costs is considered to be proportional to the square-root of the total cumulated number of orders approved.

Let's assume that the IT department told me that since the birth of the marketplace, cumulated IT costs have amounted to 500,000 BRL.

```
In [11]: # y: Cumulated IT cost, a: coefficient
y = sqrt(cumulated_orders) * a = 500,000
cumulated_orders = sorted_sellers['n_orders'].sum()
a = 500000/math.sqrt(cumulated_orders)
# a = 1582.39
# z: Current IT cost per order = 2.5 BRL
z = math.sqrt(cumulated_orders + 1) * a = 500000

Out[11]: 2.503975060360972

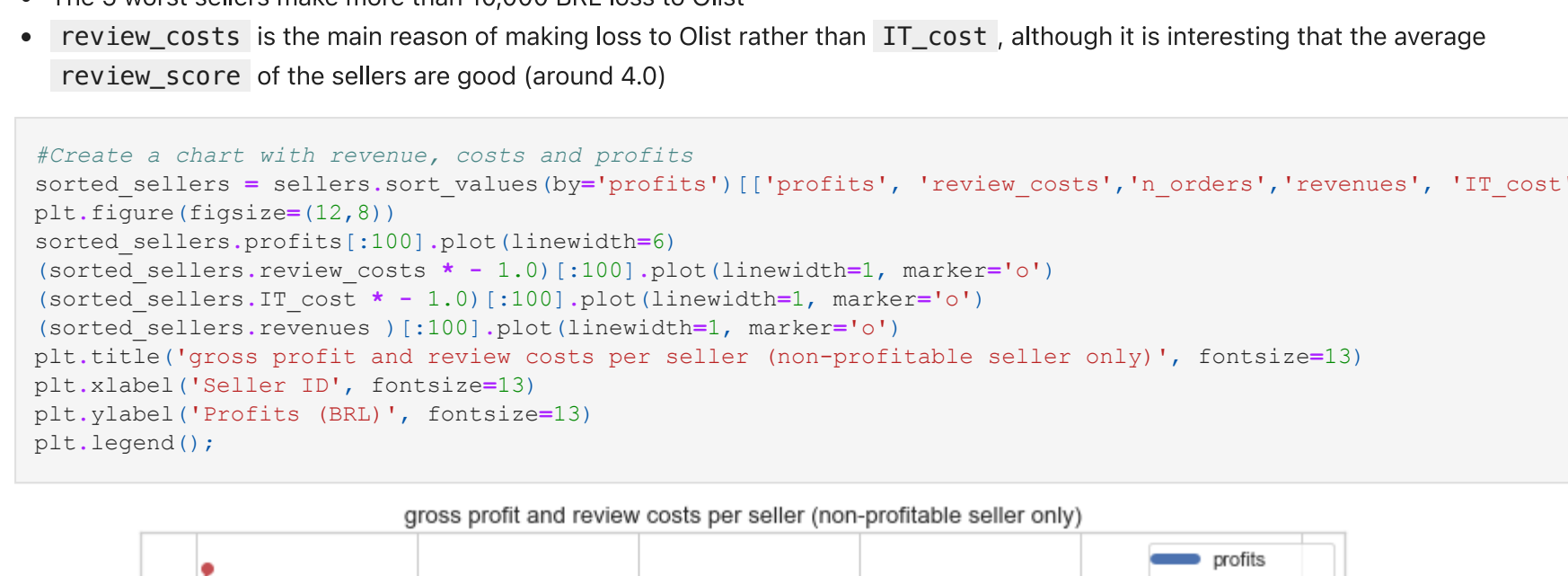
In [12]: # Calculate IT cost for each seller
sellers['IT_cost'] = sellers['n_orders'] * z

In [13]: # Distribution of IT cost and profits
sns.set_theme(style='whitegrid')
fig, ax = plt.subplots(figsize=(12,6))
sorted_sellers.profits.plot()
sns.set_seller_review_costs * 1.0 - 1.0)[:100].plot(linewidth=1, marker='o')
(sorted_sellers.revenues * 1.0)[:100].plot(linewidth=1, marker='o')
plt.title('Gross profit and review costs per seller (non-profitable seller only)', fontsize=13)
plt.xlabel('Seller ID', fontsize=13)
plt.ylabel('Profits (BRL)', fontsize=13)
plt.legend();
```



Seller who is responsible for large **IT_cost** tends to make negative profits.

Let's add **revenue** and **sales** information on the graph.



High revenue sellers do not always give positive profits to Olist

What makes seller unprofitable?

Does **IT_cost** make seller unprofitable (for Olist)?

```
In [15]: worst_sellers = np.round(sellers[['profits', 'revenues', 'review_costs', 'IT_cost', 'review_score']].sort_values(worst_sellers.head(7))

Out[15]:
```

	profits	revenues	review_costs	IT_cost	review_score
0	-21519.52	13770.48	35290	4642.37	3.94
1	-17752.71	21647.29	39400	4522.18	3.83
2	-17581.16	11948.84	29530	4271.78	4.08
3	-15542.25	4677.75	20220	2869.56	4.00
4	-12886.43	2653.57	15440	2361.25	4.07
5	-9926.08	12053.92	21980	3515.58	4.13
6	-9017.61	3562.39	12580	1312.08	3.75

- The 5 worst sellers make more than 10,000 BRL loss to Olist
- The **review_costs** is the main reason of making loss to Olist rather than **IT_cost**, although it is interesting that the average **review_score** of the sellers are good (around 4.0)

```
In [16]: # Create a chart with revenue, costs and profits
charge_df = pd.DataFrame()
charge_df['charge_rate'] = np.linspace(0, 0.2, num=5)
charge_df['charge_rate_pct'] = (charge_df.charge_rate * 100).astype('int64')
charge_df['review_charge'] = (sellers.review_costs.sum() * 100) / (sorted_sellers['profits'].sum())
charge_df['new_profits'] = sellers.profits.sum() * charge_df.review_charge
charge_df['profits_imp_pct'] = (charge_df.new_profits / sellers.profits.sum() - 1) * 100
charge_df['current_profits'] = sellers.profits.sum()
a2.set_ylabel('Bad review charge (%)', fontsize=13)
sns.barplot(x=charge_rate_pct, y=new_profits_imp_pct, data=charge_df, color='darkgreen', label='Removed sellers')
ax1.set_ylabel('Total profit improvement (%)', fontsize=13)
ax2.set_ylabel('Bad review charge (%)', fontsize=13)
ax3.set_ylabel('Total profit improvement (%)', fontsize=13)
ax4.set_ylabel('Bad review charge (%)', fontsize=13)
ax4.set_title('Threshold to remove non-profitable sellers', fontsize=13);

# The graph on the right
ax2 = ax1.twin()
line2 = sns.lineplot(x=rem_orders_pct, y=rel_increase, data=df, color='b', legend=False)
ax2.set_ylabel('Removed orders (%)', color='b', fontsize=13)
ax2.tick_params(axis='x', color='b')
bottom_line = mpatches.Patch(color='b', label='Removed orders')
plt.legend(handles=[top_line, bottom_line])

ax3 = ax1
ax3.annotate('Target profit increase = 10 % (example)', color='r', xy=(5, 10), xytext=(2, 5), fontsize=13,
            arrowprops=dict(arrowstyle='->', color='r'))
ax3.set_title('Profit increase and order volume', fontsize=13)
ax3.plot([0, 9], [10, 10], color='r', label='Removed sellers')

# The graph on the left
ax4 = fig.add_subplot(121)
line3 = sns.lineplot(x=threshold, y=rel_increase, data=df)
ax4.set_ylabel('Total profit improvement (%)', fontsize=13)
ax4.set_title('Threshold to remove non-profitable sellers', fontsize=13);
```



>> Removing the worst 0.5 % of the sellers (15 sellers) increases 10% of Olist's profits

As an alternative, it would also make sense to introduce the new fee to sellers when they receive bad reviews, since the large negative profits were due to **review_costs**.

Introducing bad review charge

What if we charge a partial cost for bad reviews to the sellers, depending on the cost to Olist by percentage?

An example of **bad review charge** per order with different rates. The numbers of the unit is BRL.

review_score	cost to Olist	5 % charge	10 % charge	15 % charge	20 % charge
1 star	100	5	10	15	20
2 stars	50	2.5	5	7.5	10
3 stars	40	2	4	6	8
4 stars	0	0	0	0	0
5 stars	0	0	0	0	0

What is the improvement on the profits based on the different charging rate (0 - 20 %)?

```
In [19]: # Current total profits
sellers.profits.sum()

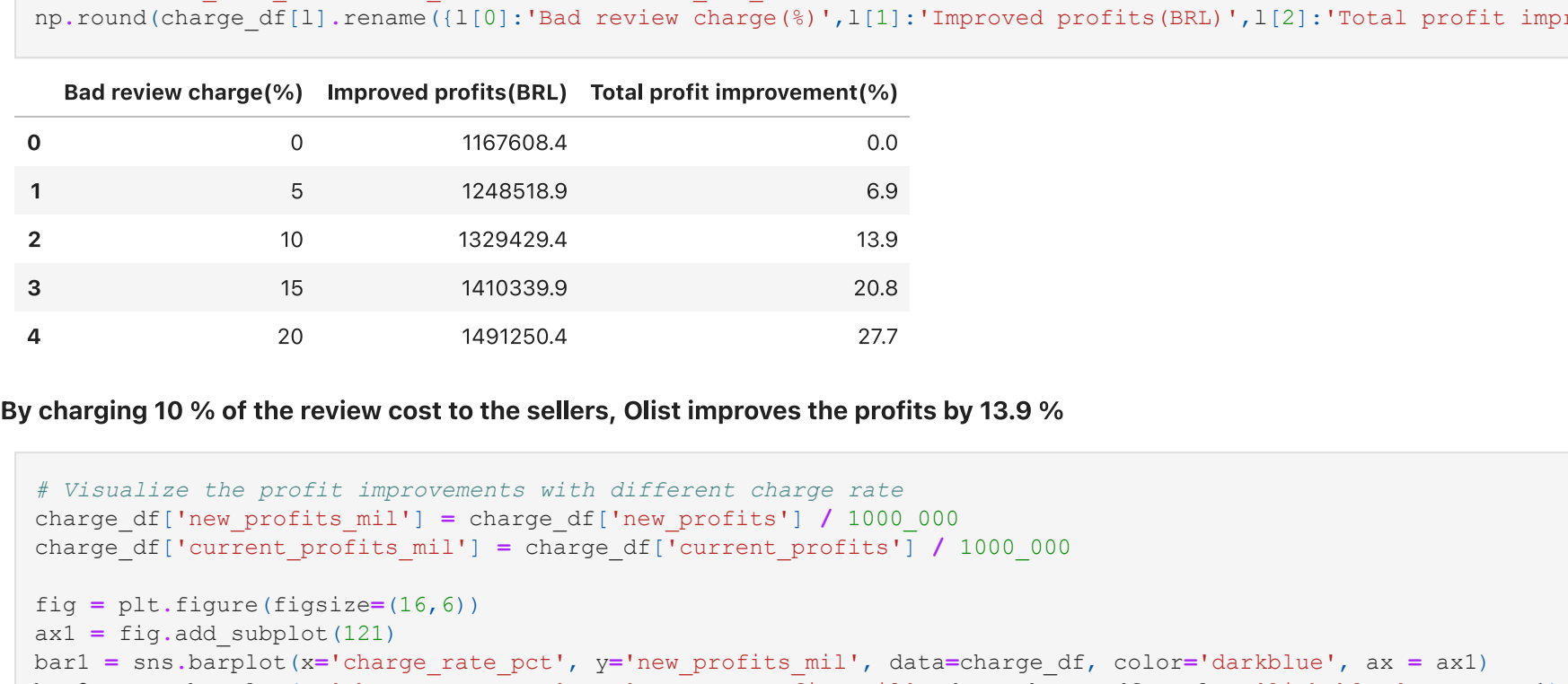
Out[19]: 1167608.4440000001

In [20]: # Create DataFrame for plot
charge_df = pd.DataFrame()
charge_df['charge_rate'] = np.linspace(0, 0.2, num=5)
charge_df['charge_rate_pct'] = (charge_df.charge_rate * 100).astype('int64')
charge_df['review_charge'] = (sellers.review_costs.sum() * 100) / (sorted_sellers['profits'].sum())
charge_df['new_profits'] = sellers.profits.sum() * charge_df.review_charge
charge_df['profits_imp_pct'] = (charge_df.new_profits / sellers.profits.sum() - 1) * 100
charge_df['current_profits'] = sellers.profits.sum()
a2.set_ylabel('Bad review charge (%)', fontsize=13)
sns.barplot(x=charge_rate_pct, y=new_profits_imp_pct, data=charge_df, color='darkgreen', label='Removed sellers')
ax1.set_ylabel('Total profit improvement (%)', fontsize=13)
ax2.set_ylabel('Bad review charge (%)', fontsize=13)
ax2.set_title('Threshold to remove non-profitable sellers', fontsize=13);

# The graph on the right
ax2 = ax1.twin()
line2 = sns.lineplot(x=rem_orders_pct, y=rel_increase, data=df, color='b', legend=False)
ax2.set_ylabel('Removed orders (%)', color='b', fontsize=13)
ax2.tick_params(axis='x', color='b')
bottom_line = mpatches.Patch(color='b', label='Removed orders')
plt.legend(handles=[top_line, bottom_line])

ax3 = ax1
ax3.annotate('Target profit increase = 10 % (example)', color='r', xy=(5, 10), xytext=(2, 5), fontsize=13,
            arrowprops=dict(arrowstyle='->', color='r'))
ax3.set_title('Profit increase and order volume', fontsize=13)
ax3.plot([0, 9], [10, 10], color='r', label='Removed sellers')

# The graph on the left
ax4 = fig.add_subplot(121)
line3 = sns.lineplot(x=threshold, y=rel_increase, data=df)
ax4.set_ylabel('Total profit improvement (%)', fontsize=13)
ax4.set_title('Threshold to remove non-profitable sellers', fontsize=13);
```



Appendix

Describing of features created on seller.py:

feature_name	type	description
seller_id	str	the id of the seller UNIQUE
seller_city	str	the city where seller is located
seller_state	str	the state where seller is located
delay_to_carrier	float	returns 0 if the order is delivered before the shipping_limit_date, otherwise the value of the delay
wait_time	float	average wait-time (duration of deliveries) per seller
date_first_sale	datetime	date of the first sale on Olist
date_last_sale	datetime	date of the last sale on Olist
months_on_olist	float	round number of months on Olist
share_of_one_stars	float	share of one-star reviews for orders in which the seller was involved
share_of_five_stars	float	share of five-star reviews for orders in which the seller was involved
review_score	float	average review score for orders in which the seller was involved
n_orders	int	number of unique orders the seller was involved with
quantity	int	total number of items sold by this seller
quantity_per_order	float	average number of items per order for this seller
sales	float	total sales associated with this seller (excluding freight value) in BRL