

Shopify Fall 2022 Data Science Intern Challenge

Upon examination of the data, I have concluded that the suspect average order value (AOV) of \$3,415.13 is a result of outlier data in both the total_items and order_amount columns. I have gathered that there were shoes with a cost of \$25.25 for a single pair as well as orders whose total items were 2000. Both of these findings when retained lead to a right skewed data distribution and subsequently an inaccurate AOV.

After scrutiny of the data, I chose to work with a subset of the original dataset. The subset I designated optimal contains 4928 rows or 98.56% of the total dataset. This subset of data eliminated all rows containing the outlier price of \$25.725 for a pair of shoes. Additionally, all rows where the total_items entry is greater than or equal to 6 were eliminated as these orders account for only 0.541% of total orders and therefore deemed atypical. Lastly, in order to have a more accurate understanding of each order, this subset contains a new column entitled **cost_one_pair** (As each store only sells one pair of shoes, the cost of one pair is exact).

```
In [1]: import pandas as pd
from matplotlib import pyplot as plt
%matplotlib inline
```

Subset of original dataset, hereafter referred to as dataset

The AOV for this dataset is \$301.47 with the average cost for one pair of shoes at \$151.80.

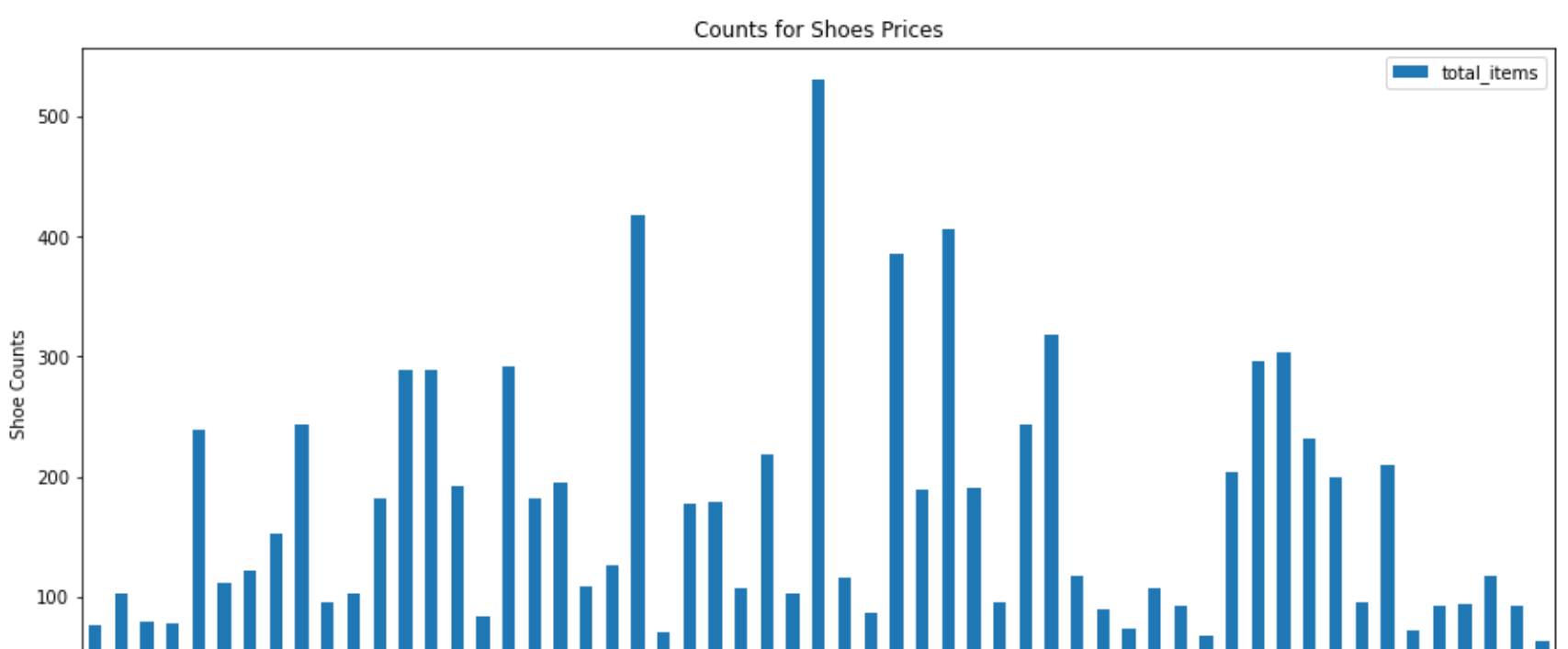
```
In [28]: new_df = df.loc[(df.cost_one_pair <= 352) & (df.total_items < 6)]
new_df.describe()
```

```
Out[28]:
```

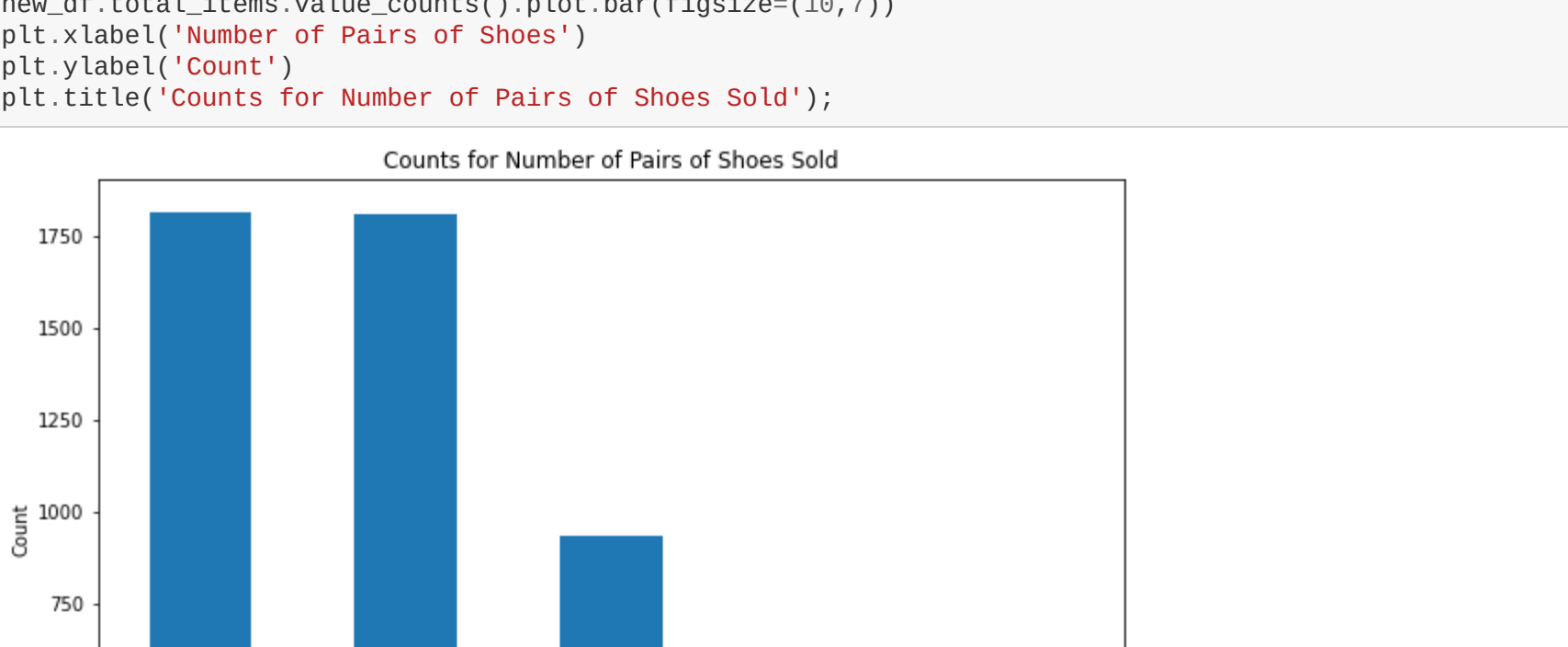
	order_id	shop_id	user_id	order_amount	total_items	cost_one_pair
count	4928.000000	4928.000000	4928.000000	4928.000000	4928.000000	
mean	2497.955560	49.822646	849.810268	301.474432	1.987013	151.798295
std	1444.154569	29.065769	86.812494	158.775640	0.966579	29.050049
min	1.000000	1.000000	700.000000	90.000000	1.000000	90.000000
25%	1245.750000	24.000000	776.000000	163.000000	1.000000	132.000000
50%	2494.500000	50.000000	850.000000	284.000000	2.000000	153.000000
75%	3749.250000	74.000000	925.000000	387.000000	3.000000	166.000000
max	5000.000000	100.000000	999.000000	1760.000000	5.000000	352.000000

The dataset has acceptable mean and 50th percentile (median) values for both order_amount and cost_one_pair. While those values for both columns have a slighter higher mean than median, the difference is considered negligible. Overall, the data is symmetrical.

```
In [29]: new_df[['cost_one_pair', 'total_items']].groupby(['cost_one_pair']).sum().plot.bar(figsize=(15,7));
plt.xlabel('Shoe Price');
plt.ylabel('Shoe Counts');
plt.title('Counts for Shoes Prices');
```



```
In [30]: new_df.total_items.value_counts().plot.bar(figsize=(10,7))
plt.xlabel('Number of Pairs of Shoes')
plt.ylabel('Shoe Price')
plt.title('Counts for Number of Pairs of Shoes Sold');
```



Recommendations

For current as well as any future stores, it is recommended that a record of the shoe price be kept per store. Additionally, as 93.4% of orders were sales of either one, two, or three pairs of shoes (goes down to 74.3% for one or two pairs), an examination of the price range for these shoes would be informative and an avenue for higher profits for Shopify stores. Finally, an exploration of the user ids for stores that had orders with the (outlier) total items of 2000 is suggested as an avenue to procure - or at minimum advertise to - similar customers as well as to ensure continuous communication with those users.

Exploratory Data Analysis

```
In [2]: df = pd.read_excel("Downloads/2019 Winter Data Science Intern Challenge Data Set.xlsx")
```

```
In [3]: df
```

```
Out[3]:
```

	order_id	shop_id	user_id	order_amount	total_items	payment_method	created_at
0	1	53	746	224	2	cash	2017-03-13 12:36:56.190
1	2	92	925	90	1	cash	2017-03-03 17:38:51.999
2	3	44	861	144	1	cash	2017-03-14 04:23:55.595
3	4	18	935	156	1	credit_card	2017-03-26 12:43:36.649
4	5	18	883	156	1	credit_card	2017-03-01 04:35:10.773
...
4995	4996	73	993	330	2	debit	2017-03-30 13:47:16.597
4996	4997	48	789	234	2	cash	2017-03-16 20:36:16.389
4997	4998	56	867	351	3	cash	2017-03-19 05:42:42.228
4998	4999	60	825	354	2	credit_card	2017-03-16 14:51:18.188
4999	5000	44	734	288	2	debit	2017-03-18 15:48:18.205

5000 rows x 7 columns

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 7 columns):
 # Column Non-Null Count Dtype
---  ---
 0 order_id 5000 non-null int64
 1 shop_id 5000 non-null int64
 2 user_id 5000 non-null int64
 3 order_amount 5000 non-null int64
 4 total_items 5000 non-null int64
 5 payment_method 5000 non-null object
 6 created_at 5000 non-null datetime64[ns]
dtypes: datetime64[ns](1), int64(5), object(1)
memory usage: 273.6+ KB
```

Descriptive statistics for dataset. Note the discrepancy in the mean and 50th percentile (median) for the order_amount column. The mean being so much greater signifies the order_amount distribution is greatly skewed right and thus asymmetrical (i.e. there are outlying data points on the high end of order_amount entries).

```
In [5]: df.describe()
```

```
Out[5]:
```

	order_id	shop_id	user_id	order_amount	total_items
count	5000.000000	5000.000000	5000.000000	5000.000000	5000.000000
mean	2500.500000	50.078800	849.092400	3145.128000	8.787120
std	1443.520003	29.006118	87.798982	41282.539349	116.320302
min	1.000000	1.000000	607.000000	90.000000	1.000000
25%	1250.750000	24.000000	775.000000	163.000000	1.000000
50%	2500.500000	50.000000	849.000000	284.000000	2.000000
75%	3750.250000	75.000000	925.000000	390.000000	3.000000
max	5000.000000	100.000000	999.000000	704000.000000	2000.000000

Average order value when all orders and all shoes are retained is \$3145.13 which is the mean of the order_amount column. This is the suspect AOV from the challenge.

```
In [6]: df.payment_method.unique()
```

```
Out[6]: array(['cash', 'credit_card', 'debit'], dtype=object)
```

Total pairs of shoes sold is 43,936

```
In [7]: items_total = df.total_items.sum()
print(items_total)
```

43936

Total amount of all sales is \$15,725,640

```
In [8]: amount_total = df.order_amount.sum()
print(amount_total)
```

15725640

Average cost of a pair of shoes when all orders and all shoes are retained is \$357.92

```
In [9]: print(amount_total/items_total)
```

357.92152221412965

Examining unique amounts in the order_amount column to look for outliers.

```
In [10]: df.order_amount.unique()
```

```
Out[10]: array([ 224,    90,   144,   156,   138,   149,   292,   266,
    161,    11,   352,   176,   164,   129,   136,   182,
    128,   165,   177,   145,   178,   155,   147,   122,
    158,   153,   148,   142,   195,   163,   173,   181,
    187,   138,   168,   184,   116,   127,   118,   132,
    117,   171,   193,   201,   131,   140,   94,   25725,
    114,   172,   169,   134,   196,   101,   168,   190,
    166,   154])
```

Adding new column **cost_one_pair** to the dataframe to ascertain the average price of one pair of shoes for each order

```
In [11]: df[['cost_one_pair']] = df.apply(lambda row: row.order_amount/row.total_items, axis=1)
df
```

```
Out[11]:
```

	order_id	shop_id	user_id	order_amount	total_items	payment_method	created_at	cost_one_pair
0	1	53	746	224	2	cash	2017-03-13 12:36:56.190	112.0
1	2	92	925	90	1	cash	2017-03-03 17:38:51.999	90.0
2	3	44	861	144	1	cash	2017-03-14 04:23:55.595	144.0
3	4	18	935	156	1	credit_card	2017-03-26 12:43:36.649	156.0
4	5	18	883	156	1	credit_card	2017-03-01 04:35:10.773	156.0
...
4995	4996	73	993	330	2	debit	2017-03-30 13:47:16.597	165.0
4996	4997	48	789	234	2	cash	2017-03-16 20:36:16.389	117.0
4997	4998	56	867	351	3	cash	2017-03-19 05:42:42.228	117.0
4998	4999	60	825	354	2	credit_card	2017-03-16 14:51:18.188	177.0
4999	5000	44	734	288	2	debit	2017-03-18 15:48:18.205	144.0

5000 rows x 8 columns

Examine unique amounts in new column **av_cost_one_pair** to look for outliers. The highest prices for a pair of shoes are \$352 and \$25.725.

```
In [12]: df.cost_one_pair.unique()
```

```
Out[12]: array([ 112.,    90.,   144.,   156.,   138.,   149.,   146.,   133.,
    161.,   11.,   352.,   176.,   164.,   129.,   136.,   182.,
    128.,   165.,   177.,   145.,   178.,   155.,   147.,   122.,
    158.,   153.,   148.,   142.,   195.,   163.,   173.,   181.,
    187.,   138.,   168.,   184.,   116.,   127.,   118.,   132.,
    117.,   171.,   193.,   201.,   131.,   140.,   94.,   25725,
    114.,   172.,   169.,   134.,   196.,   101.,   168.,   190.,
    166.,   154.] )
```

Examining subset of dataset which contains only the two highest prices for a pair of shoes. This constraint yields a subset of 97 rows or 1.94% of total dataset.

```
In [13]: df.loc[df.cost_one_pair > 201]
```

```
Out[13]:
```

	order_id	shop_id	user_id	order_amount	total_items	payment_method	created_at	cost_one_pair
15	16	42	607	704000	2000	credit_card	2017-03-07 04:00:00.000	352.0
40	41	42	793	352	1	credit_card	2017-03-24 14:15:40.649	352.0
60	61	42	607	704000	2000	credit_card	2017-03-04 04:00:00.000	352.0
160	161	78	990	25725	1	credit_card	2017-03-12 05:56:8.834	25725.0
308	309	42	770	352	1	credit_card	2017-03-11 18:16:38.774	352.0
...
4745	4746	42	872	352	1	debit	2017-03-24 00:57:24.130	352.0
4767	4768	42	720	704	2	credit_card	2017-03-16 10:26:08.027	352.0
4868	4869	42	607	704000	2000	credit_card	2017-03-22 04:00:00.000	352.0
4882	4883	42	607	704000	2000	credit_card	2017-03-25 04:00:00.000	352.0
4918	4919	78	823	25725	1	cash	2017-03-15 13:26:46.262	25725.0

97 rows x 8 columns

Examining subset of dataset which contains only the highest price (\$25.725) for a pair of shoes. This constraint yields a subset of 46 rows or 0.92% of total dataset. Note, with one pair of shoes being so exorbitant, the order_amount for these entries greatly skews the amount for total sales of all orders.

```
In [14]: df.loc[df.cost_one_pair > 352]
```

```
Out[14]:
```

	order_id	shop_id	user_id	order_amount	total_items	payment_method	created_at	cost_one_pair
160	161	78	990	25725	1	credit_card	2017-03-12 05:56:8.190	25725.0
490	491	78	986	51450	2	debit	2017-03-26 17:08:18.911	25725.0
493	494	78	993	51450	2	cash	2017-03-16 21:39:35.400	25725.0
511	512	78	967	51450	2	cash	2017-03-09 07:23:13.640	25725.0
617	618	78	760	51450	2	cash	2017-03-11 18:41:8.488	25725.0
691	692	78	878	514350	6	debit	2017-03-27 15:41:43.203	25725.0
1056	1057	78	800	25725	1	debit	2017-03-15 10:16:48.830	25725.0
1193	1194	78	944	25725	1	debit	2017-03-16 16:38:25.551	25725.0
1204	1205	78	970	25725	1	credit_card	2017-03-27 11:04:04.863	25725.0
1259	1260	78	775	77175	3	credit_card	2017-03-27 09:27:21.943	25725.0
1384	1385	78	867	25725	1	cash	2017-03-17 16:38:06.279	25725.0
1419	1420	78	912	25725	1	cash	2017-03-20 12:23:42.551	25725.0
1452	1453	78	812	25725	1	credit_card	2017-03-17 18:09:54.089	25725.0
1529	1530	78	810	51450	2	cash	2017-03-19 04:06:46.666	25725.0
2270	2271	78	855	25725	1	credit_card	2017-03-24 17:58:21.635	25725.0
2452	2453	78	709	51450	2	cash	2017-03-27 11:04:04.863	25725.0
2492	2493	78	834	102900	4	debit	2017-03-04 04:37:33.848	25725.0
2495	2496	78	707	51450	2	cash	2017-03-16 04:38:52.497	25725.0
2512	2513	78	935	51450	2	debit	2017-03-28 04:38:57.142	25725.0
2548	2549	78	861	25725	1	cash	2017-03-17 19:35:59.663	25725.0
2564	2565	78	915	77175	3	debit	2017-03-25 01:19:35.410	25725.0
2690	2691	78	962	77175	3	debit	2017-03-27 07:33:25.104	25725.0
2773	2774	78	890	25725	1	cash	2017-03-26 10:36:43.445	25725.0
2818	2819	78	869	51450	2	debit	2017-03-17 06:25:50.921	25725.0
2821	2822	78	814	51450	2	cash	2017-03-02 17:13:25.571	25725.0
2906	2907	78	817	77175	3	debit	2017-03-16 03:45:46.089	25725.0
2922	2923	78	740	25725	1	debit	2017-03-12 20:10:58.008	25725.0
3085	3086	78	910	25725	1	cash	2017-03-26 01:59:26.748	25725.0
3101	3102	78	855	51450	2	credit_card	2017-03-21 05:10:34.147	25725.0
3151	3152	78	745	25725	1	credit_card	2017-03-18 13:10:07.198	25725.0
3167	3168	78	927	51450	2	cash	2017-03-12 12:23:07.516	25725.0
3403	3404	78	928	77175	3	debit	2017-03-16 09:45:04.544	25725.0
3440	3441	78	982	25725	1	debit	2017-03-19 19:02:53.732	25725.0
3705	3706	78	828	51450	2	credit_card	2017-03-16 14:13:25.868	25725.0
3724	3725	78	869	77175	3	credit_card	2017-03-14 14:13:25.868	25725.0
3780	3781	78	786	25725	1	cash	2017-03-11 21:14:49.542	25725.0
4040	4040	78	852	25725	1	cash	2017-03-02 14:31:1	