

Data Storytelling

This step is to explore the dataset to discover certain trends.

The goal of the exploratory data analysis was to find out whether a certain demographic is more likely to get approved compared to others.

The dataset was truncated to create a dataset centered around gender and education. In order to compare the demographics perfectly, different histograms were used to explore these newly-formed dataset.

Plotting graphs for different features includes repetitive line of codes, so to overcome this scenario different plotting functions were defined. These functions were called as per the need of visualization.

★ Starting with the graph that depicts the dataset :-

```
#Percentages of Fraudulent and Non-Fraudulent transactions in data
## Here f-string is used to make string interpolation simple and it is faster than str.format
print(f'Percent of Non-Fraudulent Transactions = {round(data["Class"].value_counts()[0]/len(data) * 100,3)}%')
print(f'Percent of Non-Fraudulent Transactions = {round(data["Class"].value_counts()[1]/len(data) * 100,3)}%')

Percent of Non-Fraudulent Transactions = 99.827%
Percent of Non-Fraudulent Transactions = 0.173%
```

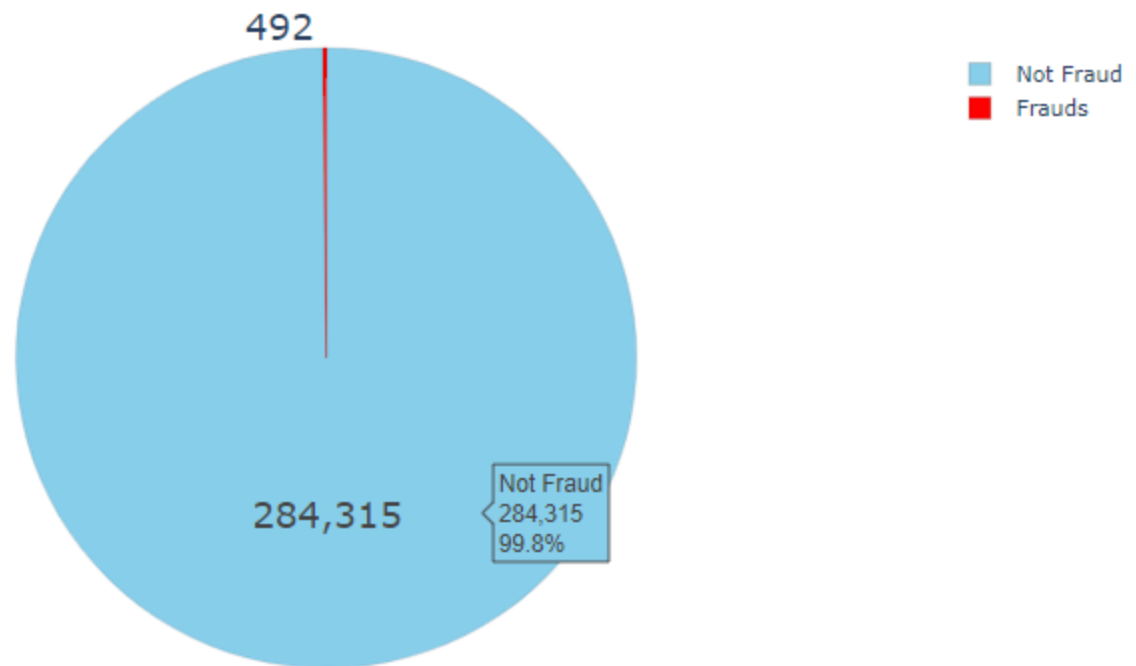
```
#plotting a pie chart for fraud and non-fraud transactions

fraud_or_not = data["Class"].value_counts().tolist()

labels = ['Not Fraud', 'Frauds']
values = [fraud_or_not[0], fraud_or_not[1]]
colors = ['skyblue', 'red']

trace = go.Pie(labels=labels, values=values, textinfo='value',
               textfont=dict(size=20),
               marker=dict(colors=colors, line=dict(color='#000000', width=0.1)))

plotly.offline.iplot([trace], filename='styled_pie_chart')
```

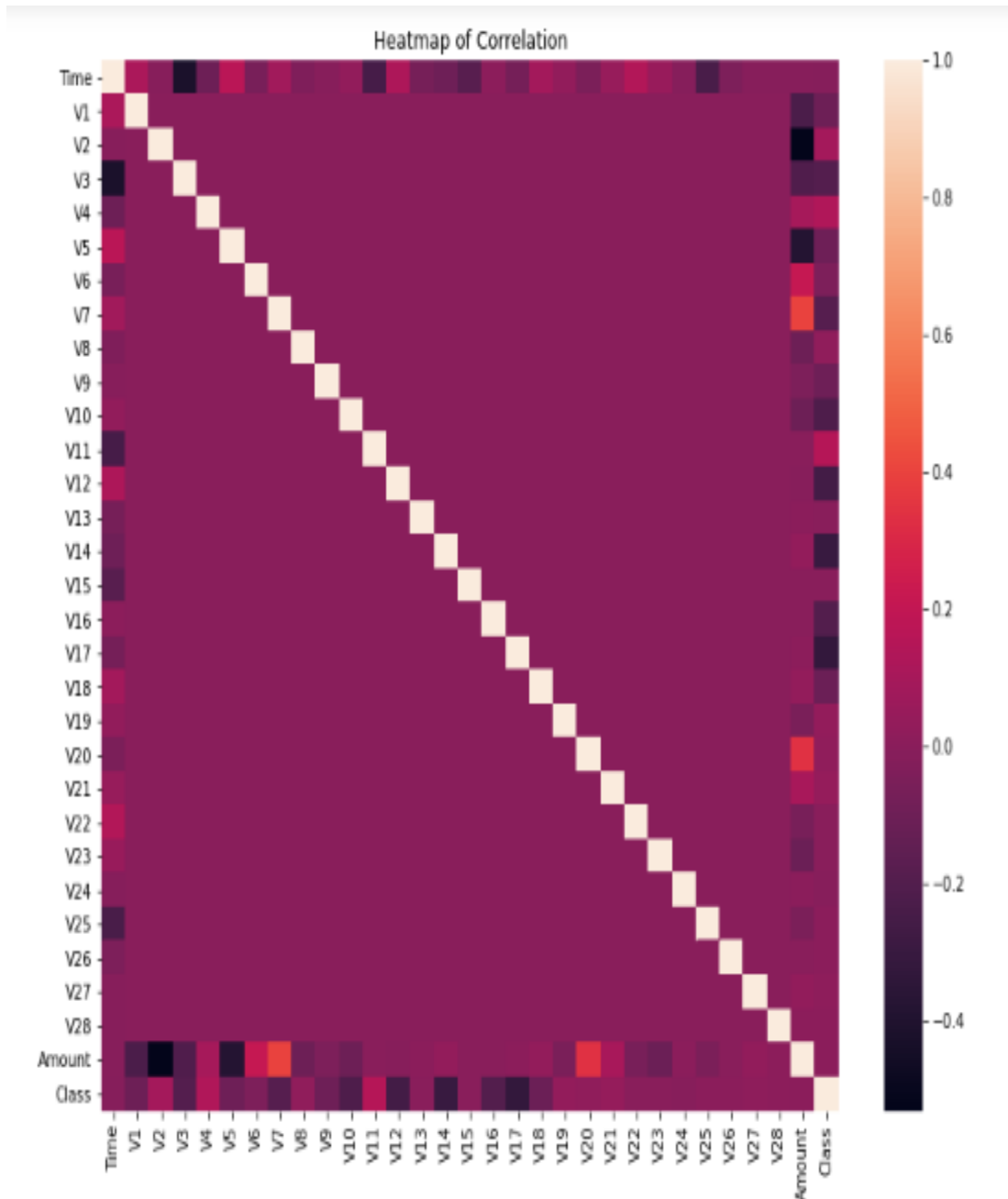


★ This dataset has 492 frauds out of 284,315 transactions. Thus, the dataset is highly unbalanced, the positive class (frauds) account for 0.173% of all transactions.

★ **Heat Map Correlation:-**

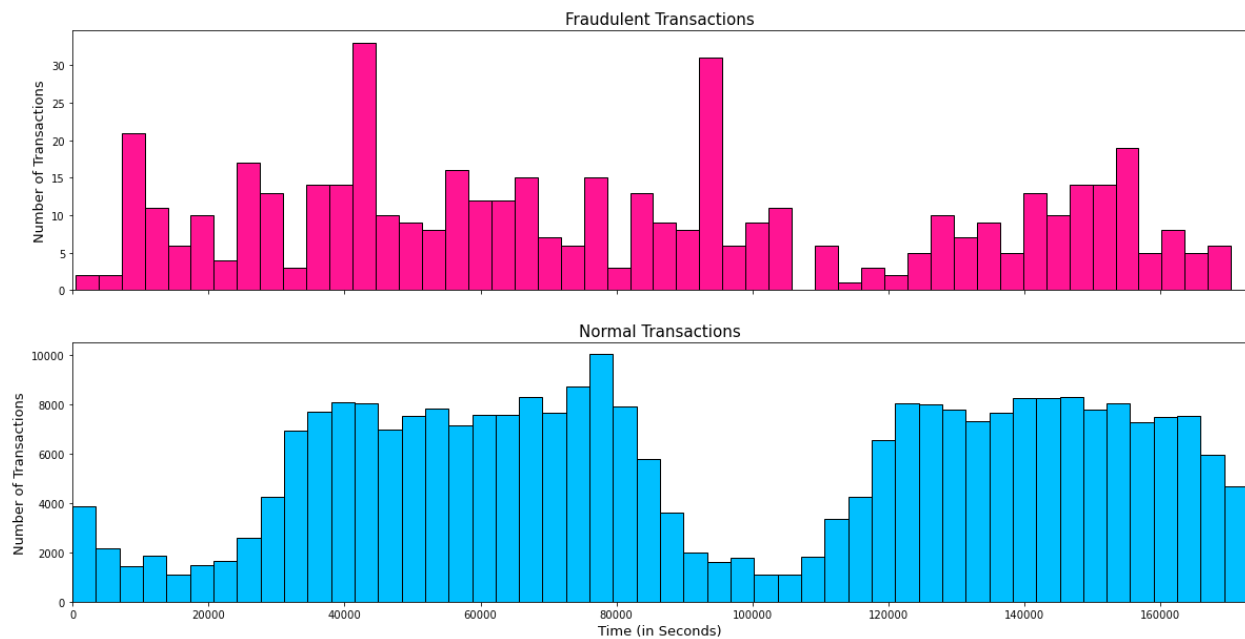
```
#heatmap
corr = data.corr()
plt.figure(figsize=(12,10))
heat = sns.heatmap(data=corr)
plt.title('Heatmap of Correlation')
```

```
Text(0.5, 1.0, 'Heatmap of Correlation')
```



- In the HeatMap we can clearly see that most of the features do not correlate to other features but there are some features that either has a positive or a negative correlation with each other.
- For example “V2” and “V5” are highly negatively correlated with the feature called “Amount”. We also see some correlation with “V20” and “Amount”. This gives us a deeper understanding of the Data available to us.

★ Analysis of Time Column:-



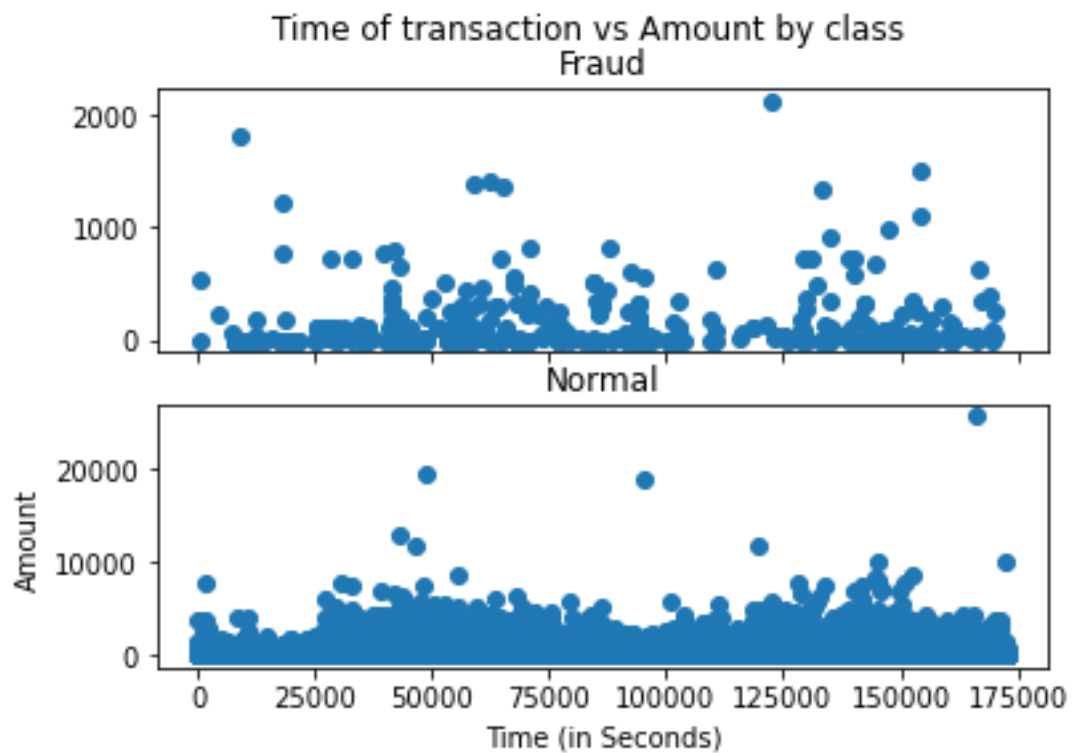
FRAUDULENT

- There are two peaks between 40000 seconds and 100000 seconds which were the maximum number of fraudulent transaction at any time.

NORMAL

- Normal transactions have not much to uncover except the fact that there were less transactions somewhere around 20000 seconds and 100000 seconds which is not very useful.

★ **Class and Amount vs Time:-**



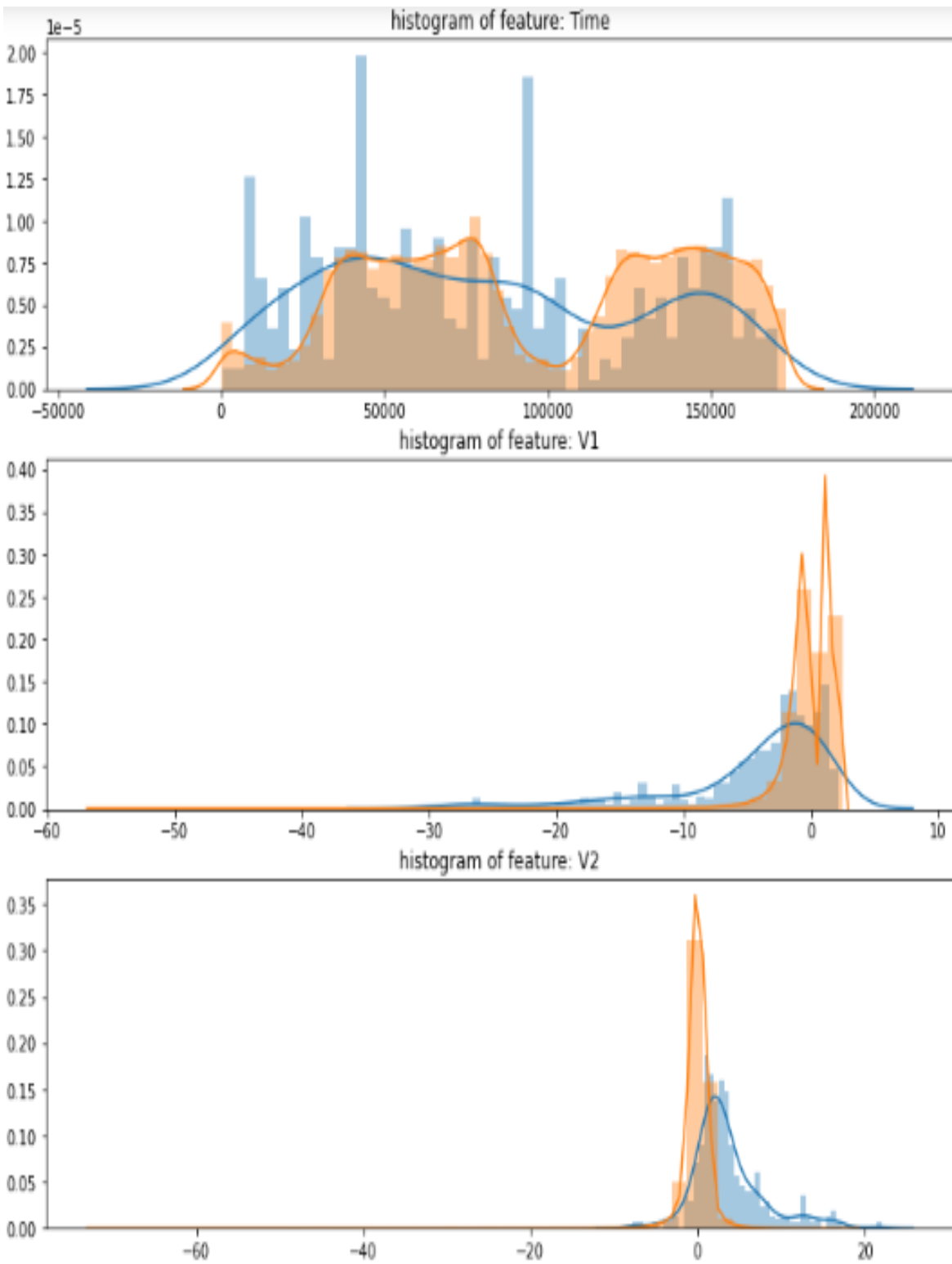
FRAUDULENT

- ▶ There are much more outliers as compared to normal transactions.
- ▶ The plot seems to not have any inherent pattern.

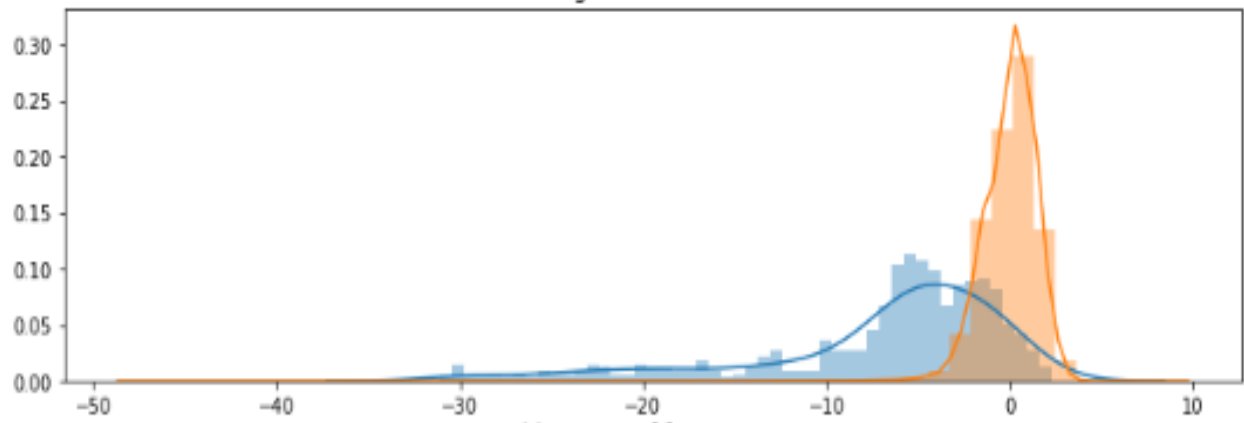
NORMAL

- ▶ There are a less number of outliers as compared to fraudulent transactions.
- ▶ There are a lot of transactions with amount less than 5000.

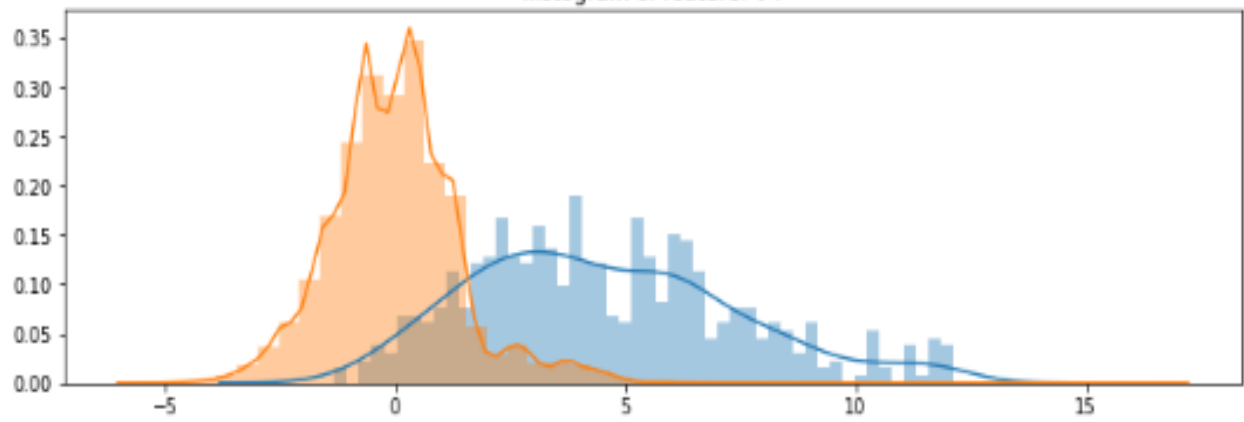
★ Distribution of anomalous features:-



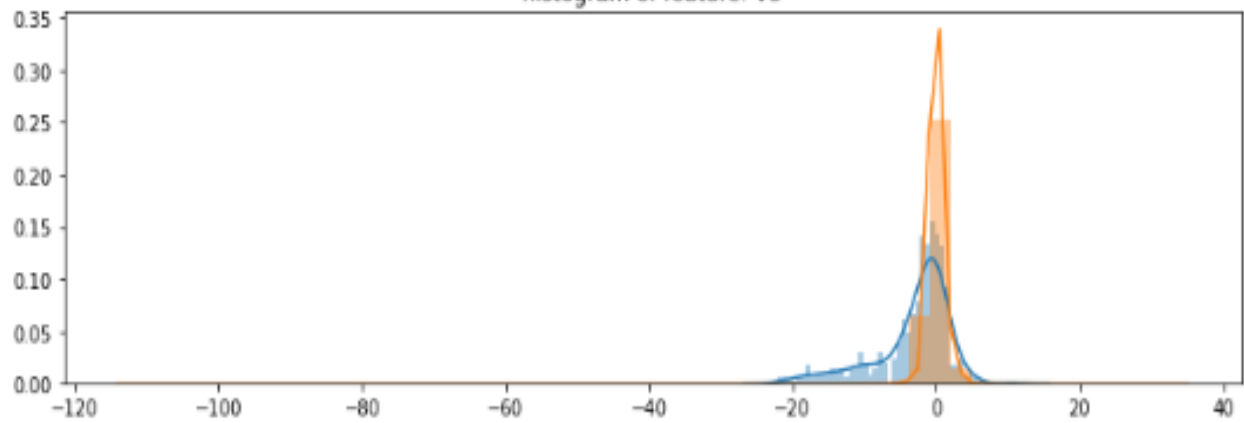
histogram of feature: V3



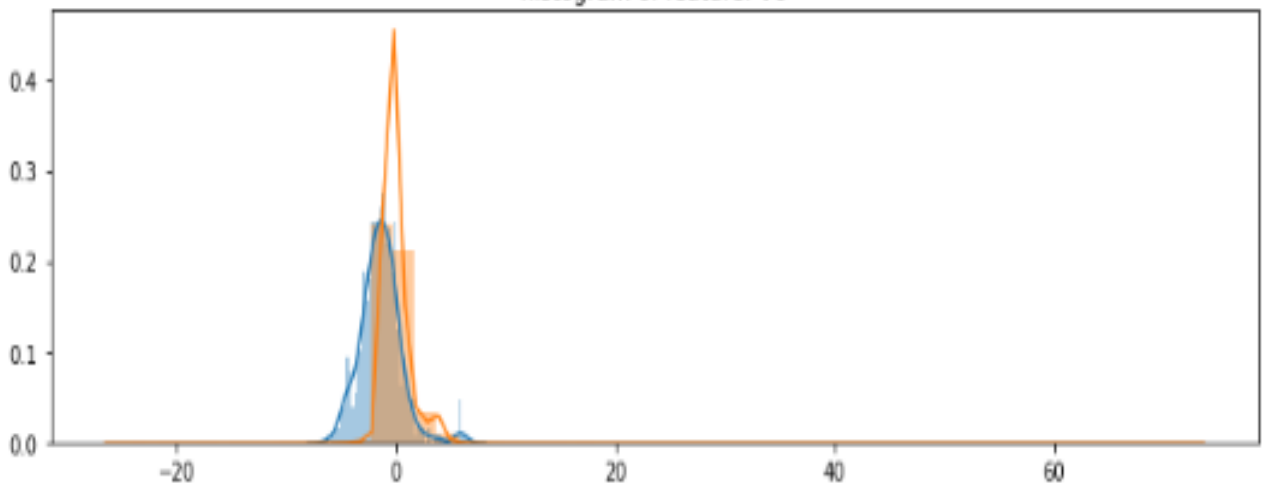
histogram of feature: V4



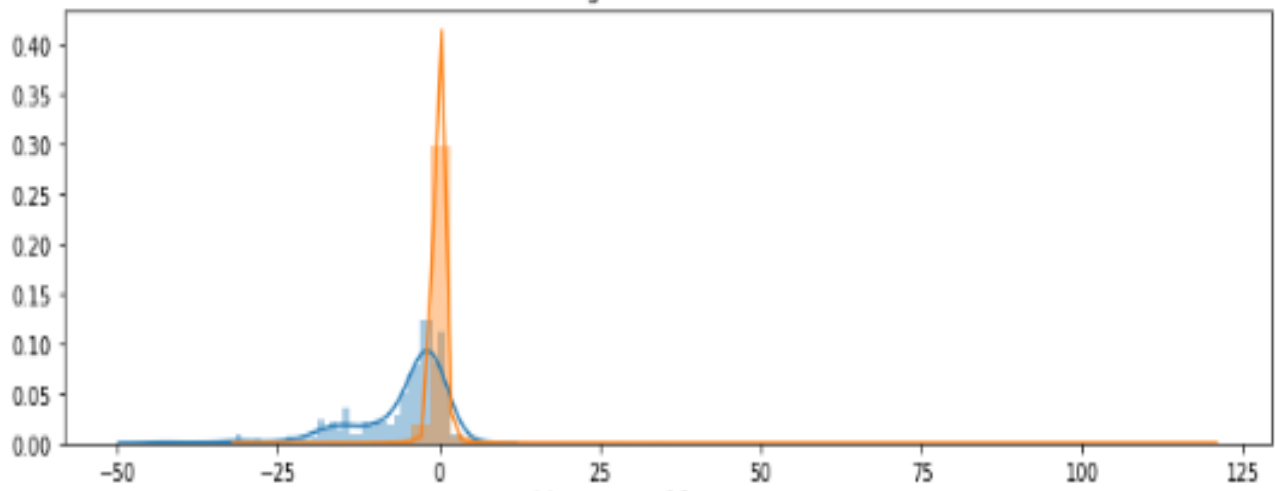
histogram of feature: V5



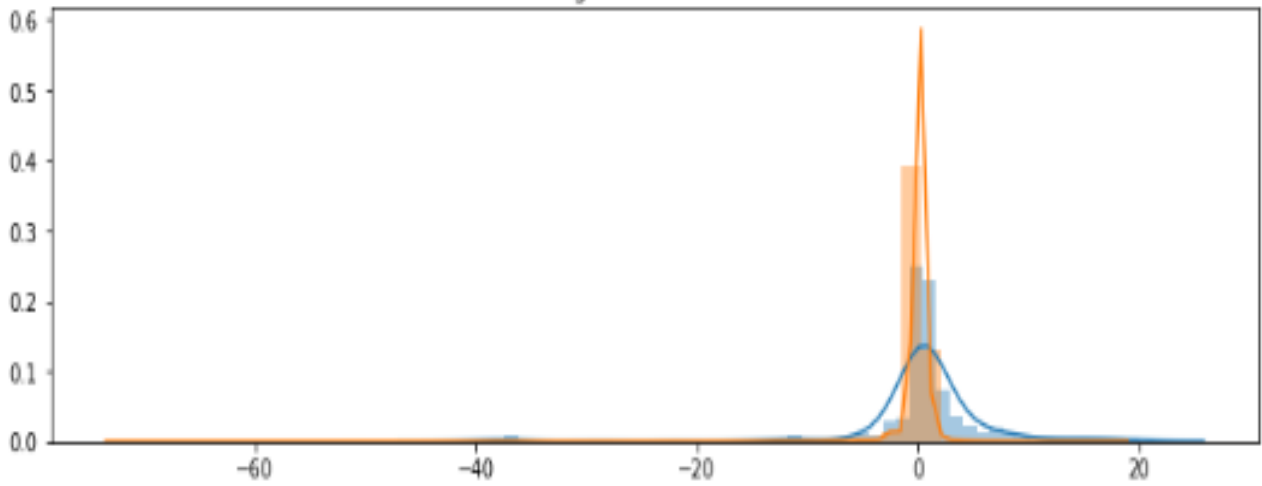
histogram of feature: V6

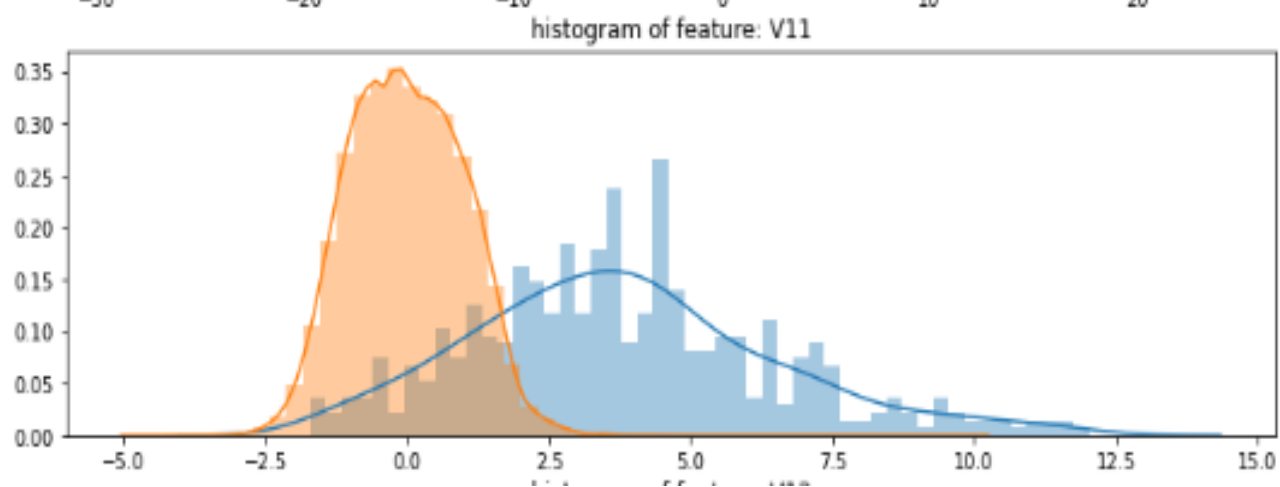
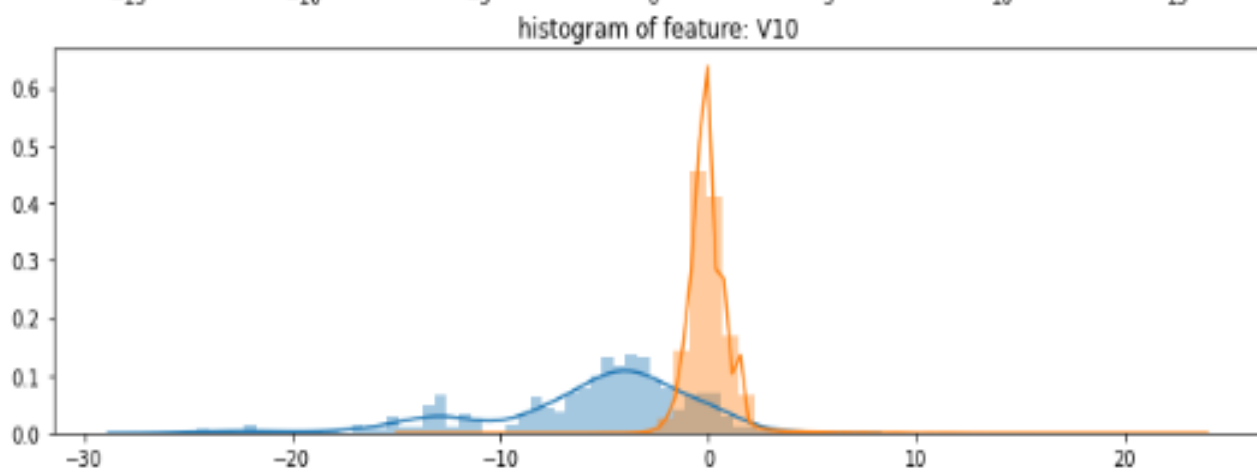
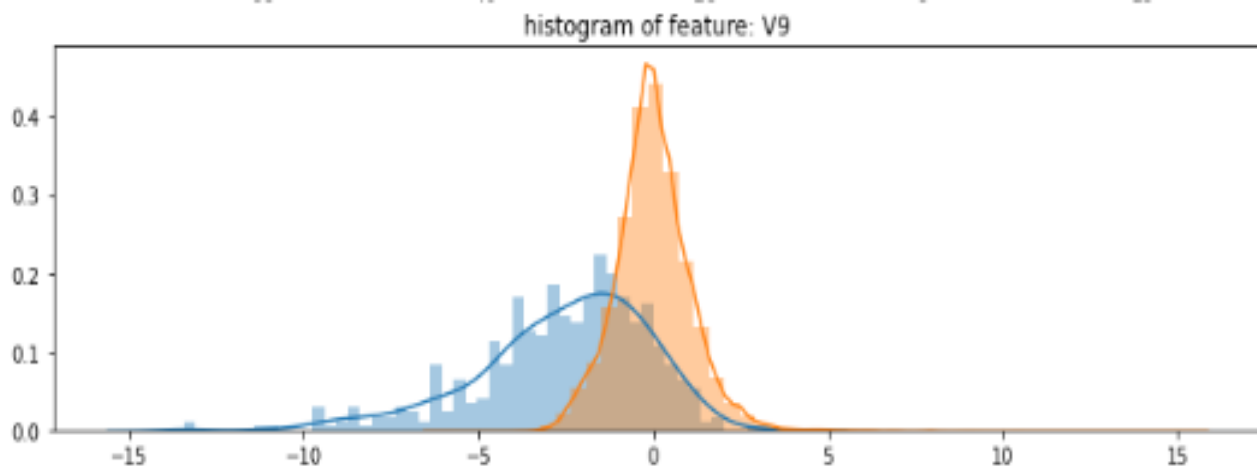


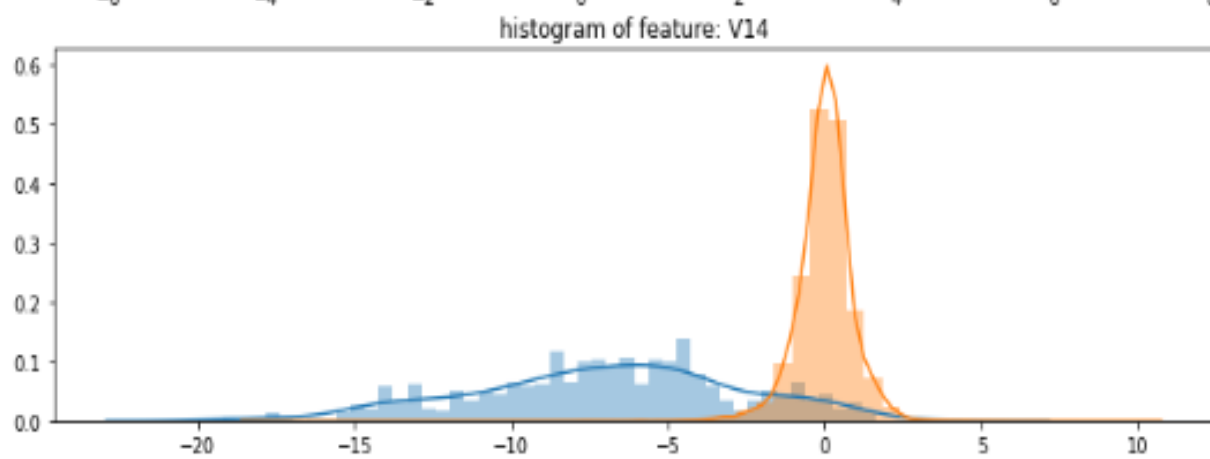
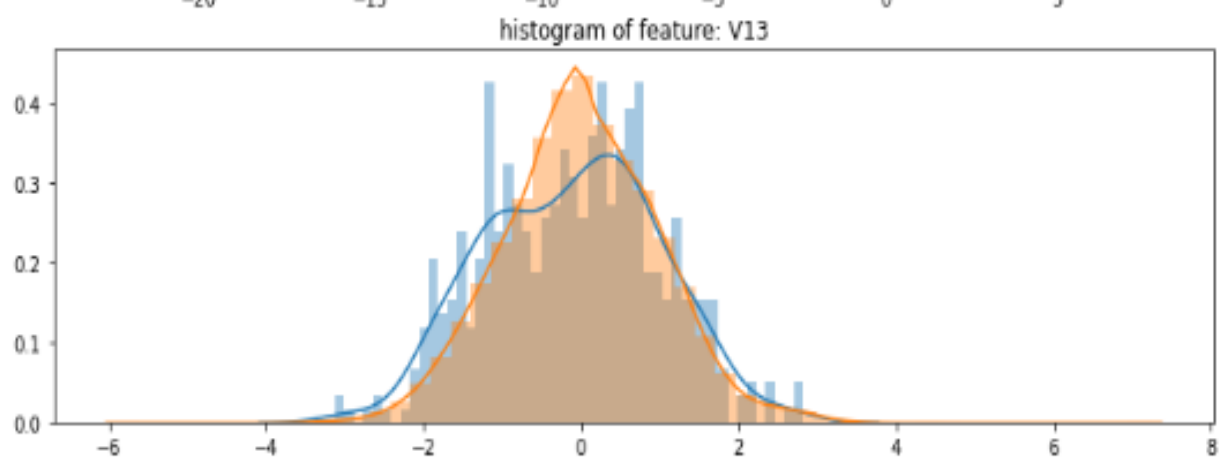
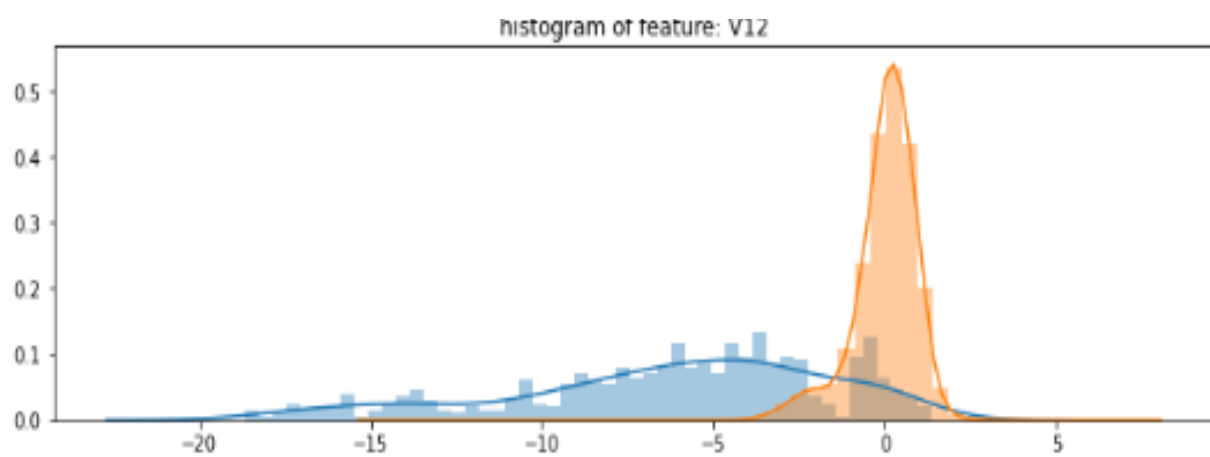
histogram of feature: V7

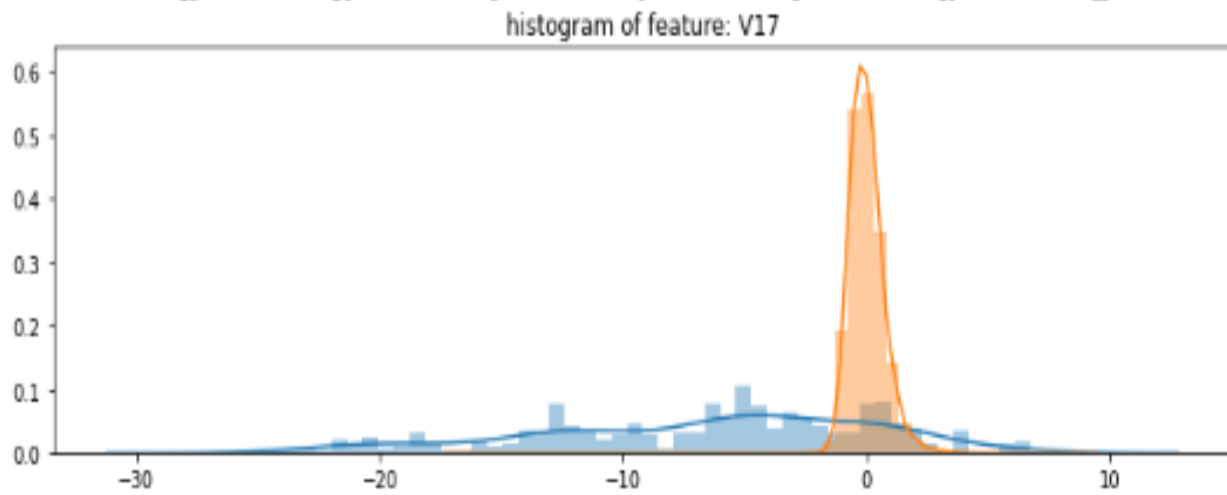
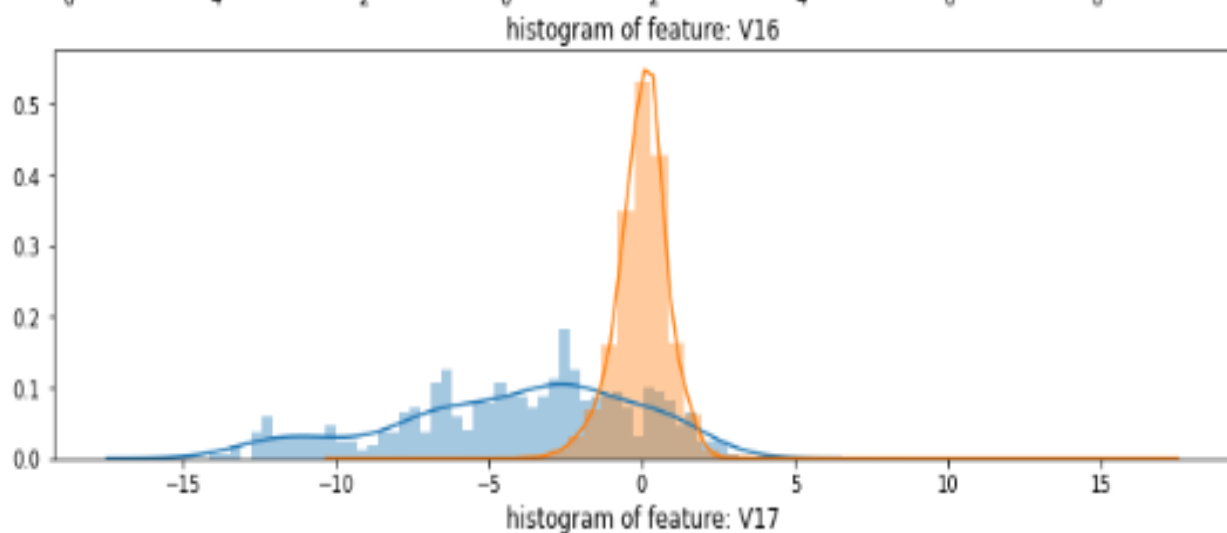
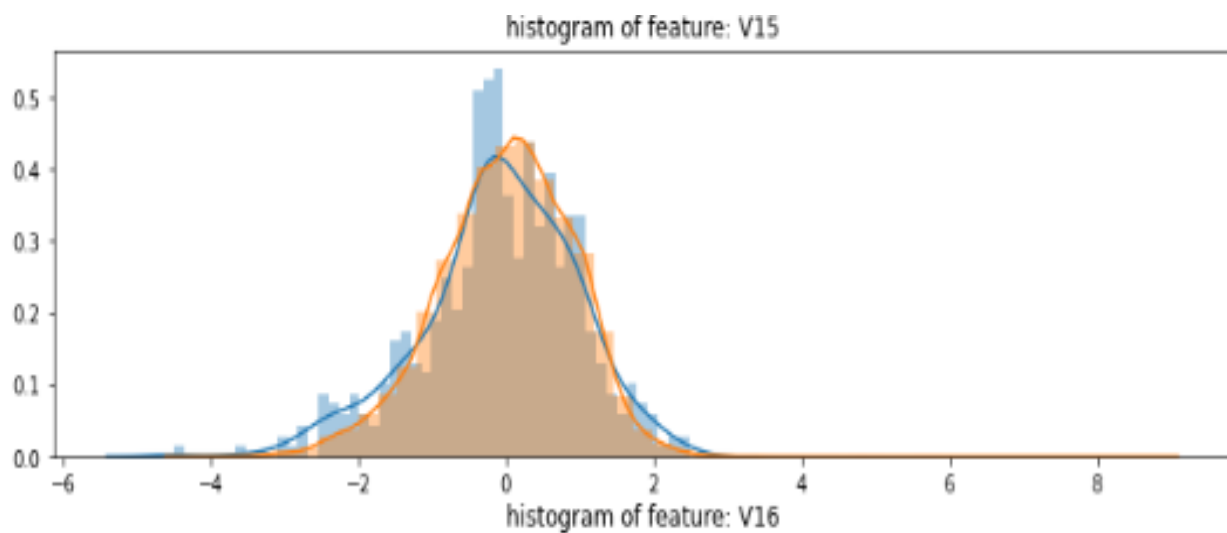


histogram of feature: V8

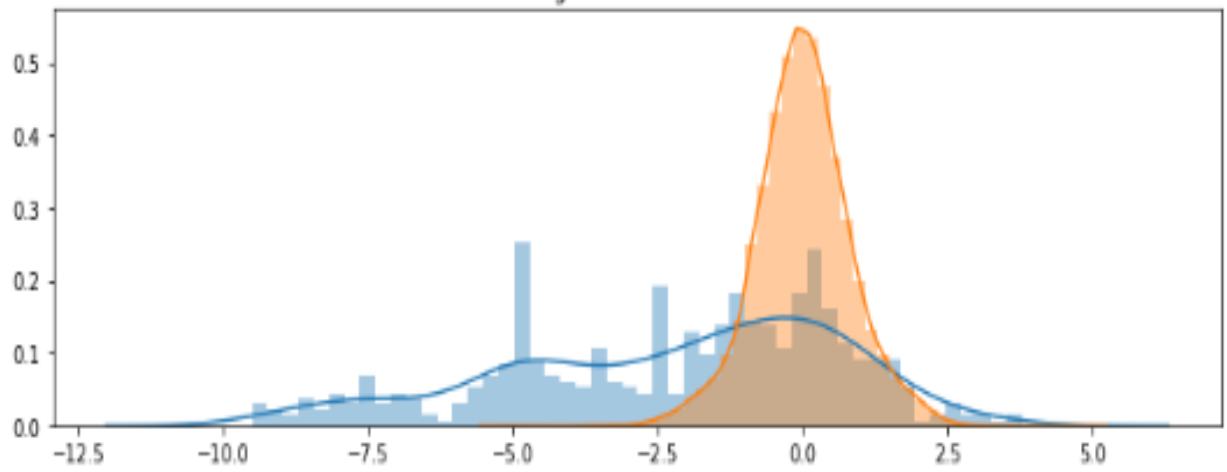




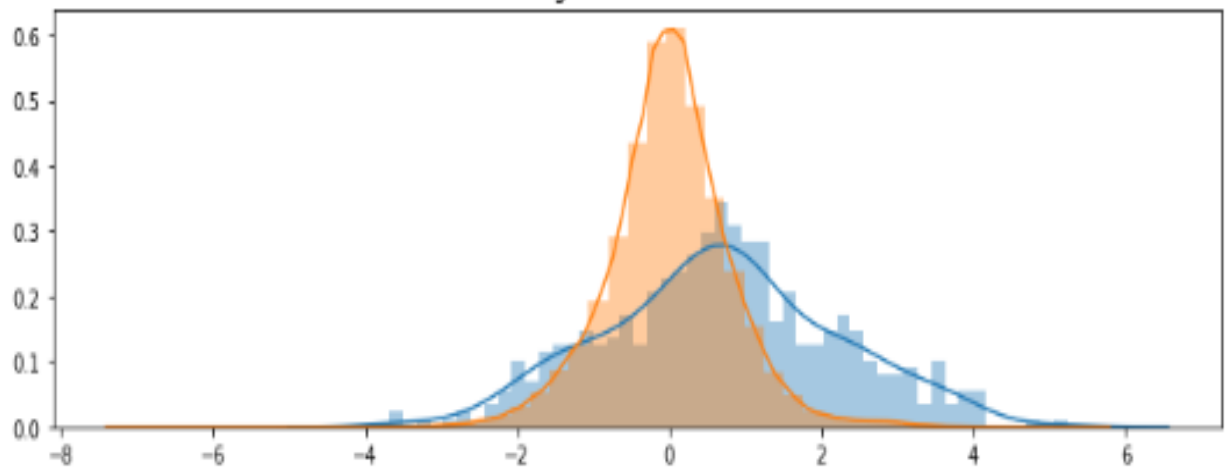




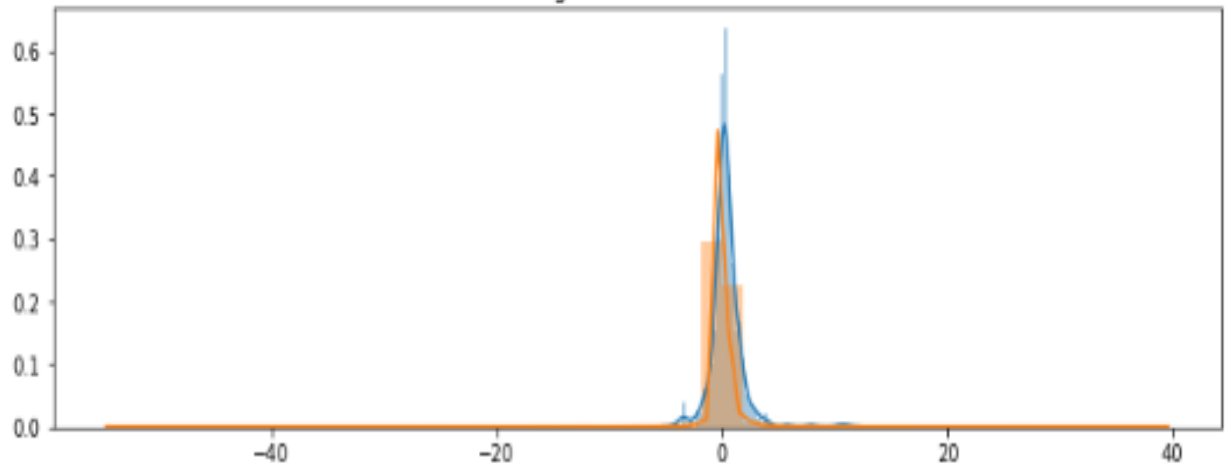
histogram of feature: V18



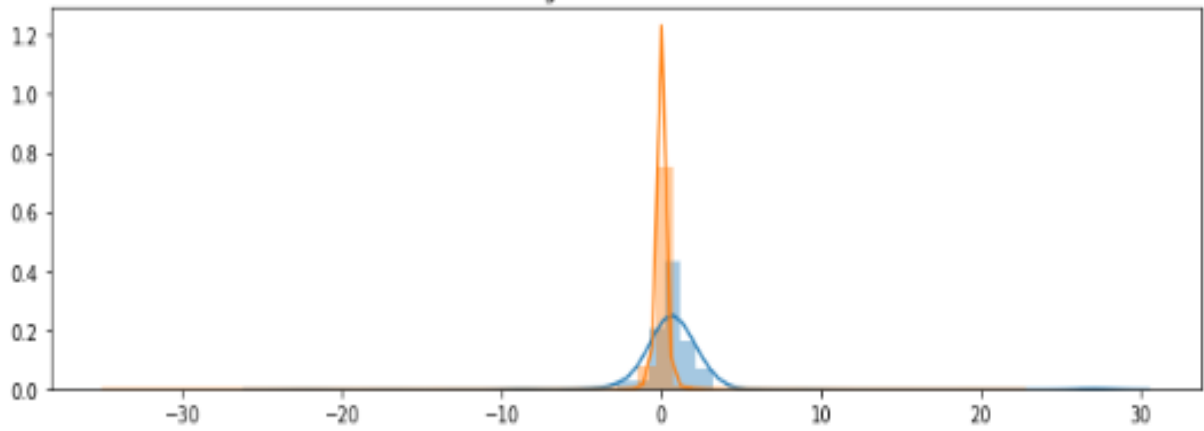
histogram of feature: V19



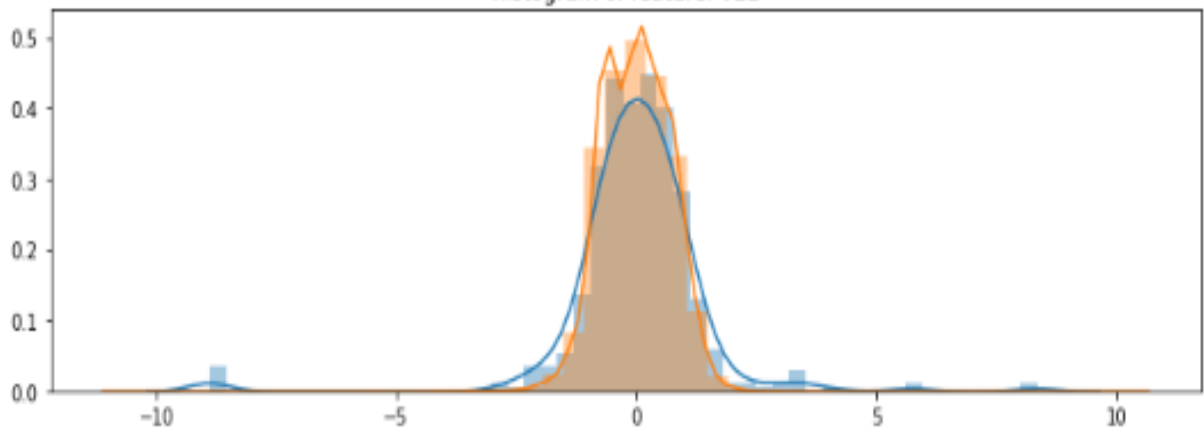
histogram of feature: V20



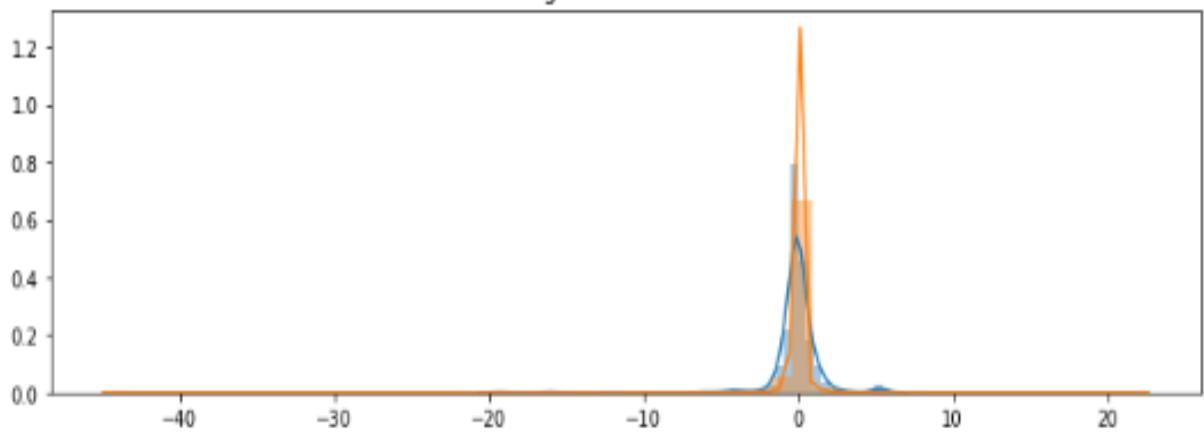
histogram of feature: V21

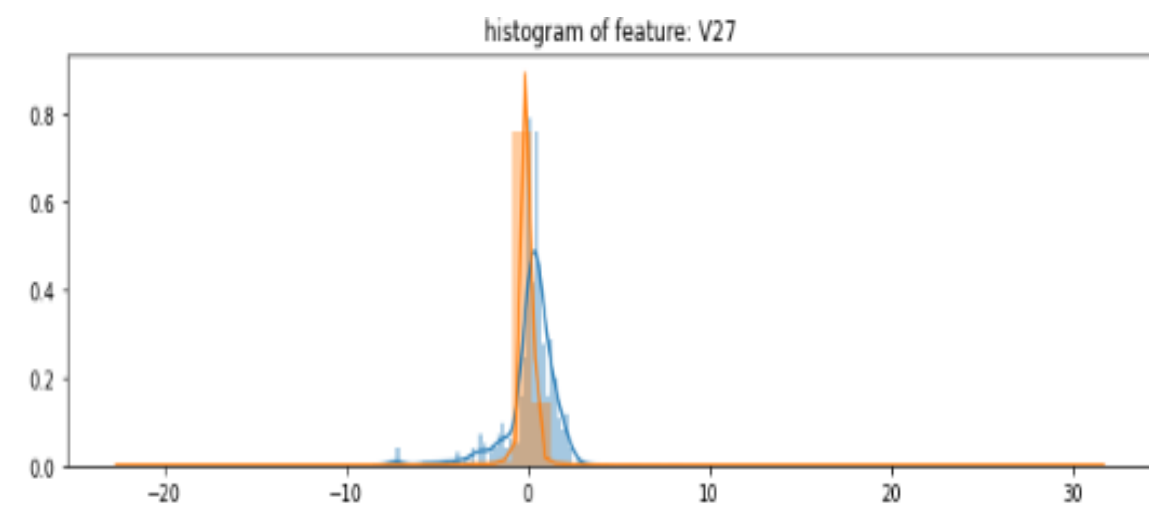
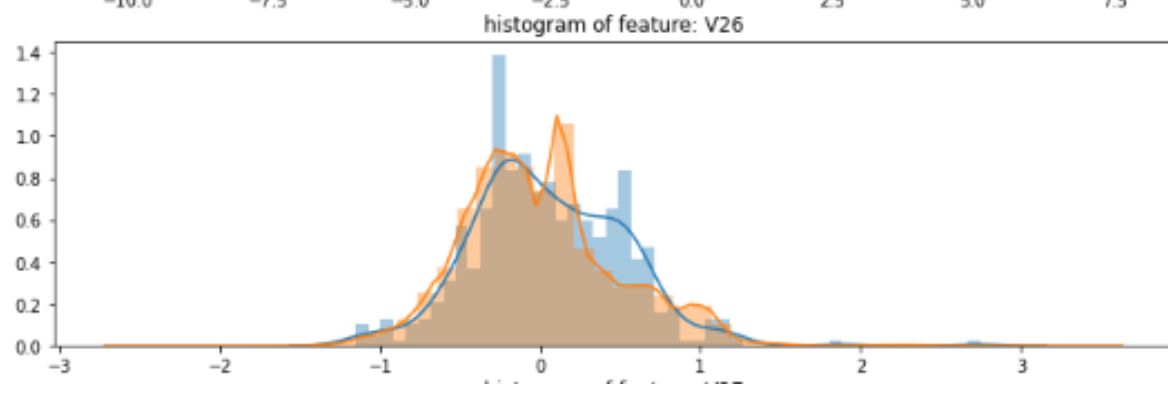
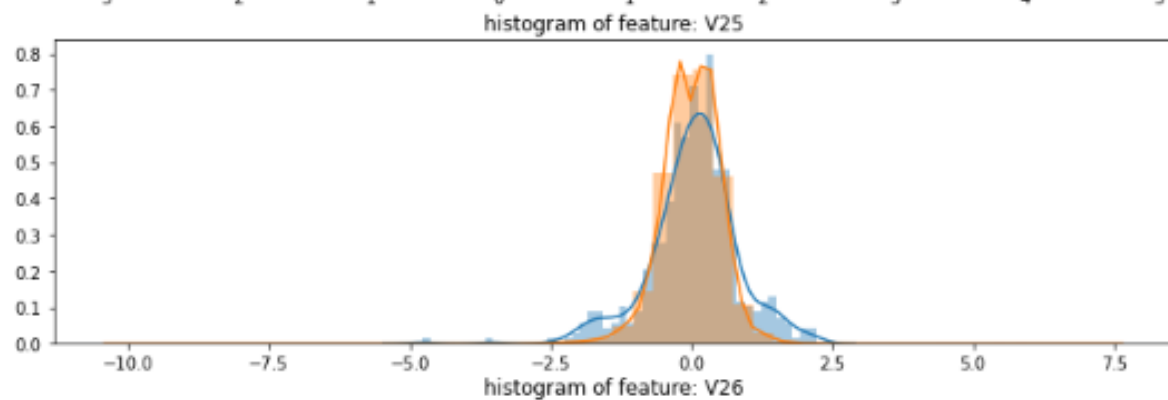
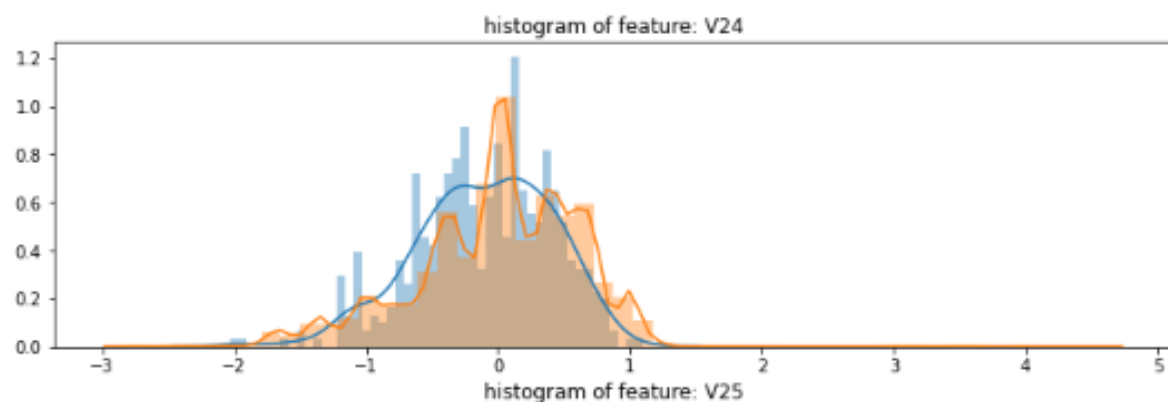


histogram of feature: V22

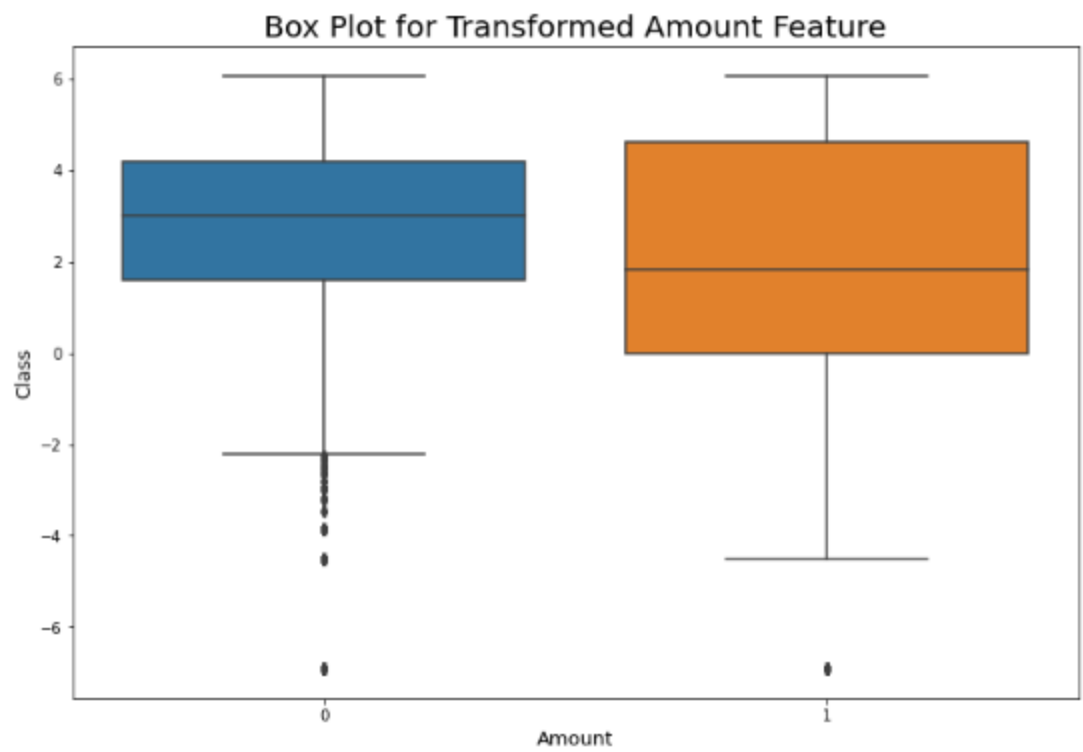


histogram of feature: V23





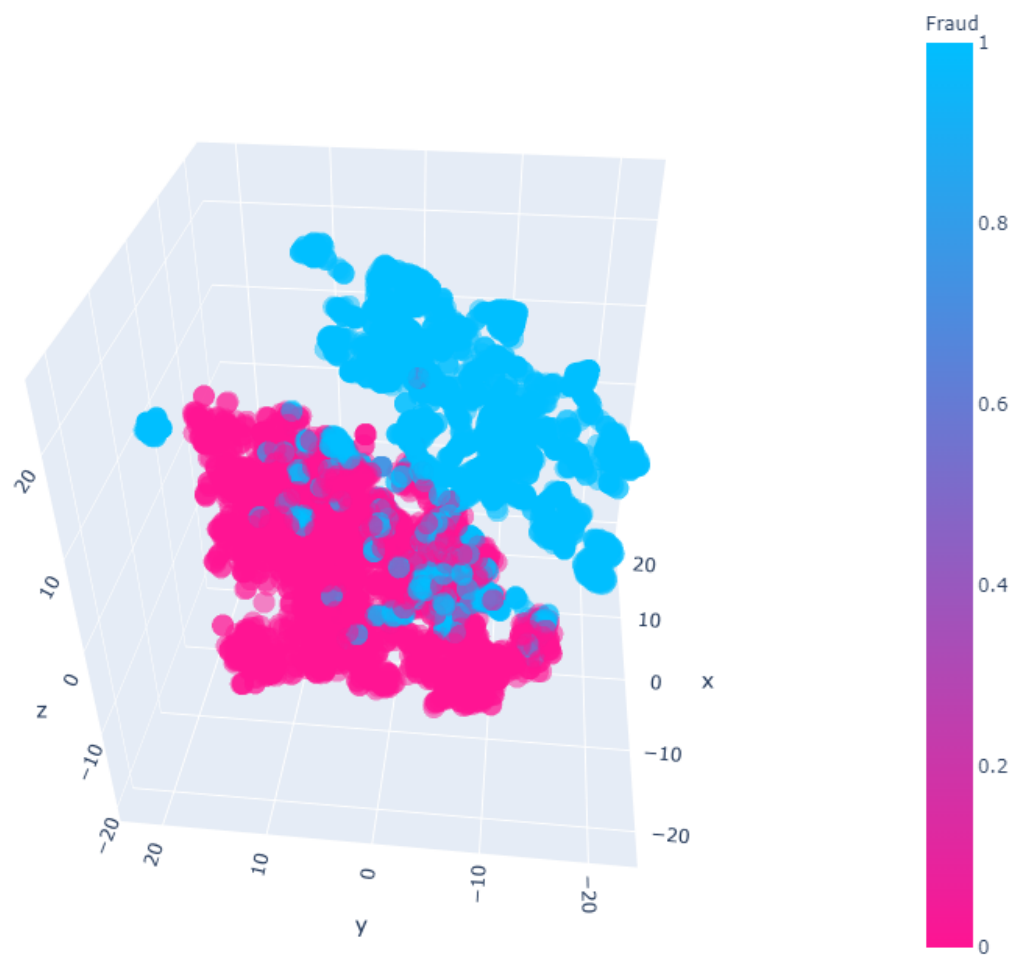
★ BOX PLOT FOR TRANSFORMED AMOUNT FEATURE:-



► After transforming a highly skewed amount feature.

★ Dimensionality Reduction And Clustering:-

TSNE (T-Distributed Stochastic Neighbourhood Embedding)



► As expected, t-SNE is performing well, separating both the classes. The clusters are not perfectly separated but it is able to separate almost 80% of the sample data which is cool.