

Statistical learning: Second assignment

Ali Zamani(96123035)

December 14, 2018

0.1 Visualize dataset

In this project I use creditcard dataset The dataset contains transactions made by credit cards in September 2013 by European cardholders over a two day period. There are 492 frauds out of a total 284,807 examples. Thus, the dataset is highly unbalanced, with the positive class (frauds) accounting for only 0.172% of all transactions. You can imagine that any such dataset would be highly unbalanced, as expected fraud or anomalous cases would only make up for a small percentage of the total transactions. Let's have look at our dataset.

I used seaborn and matplotlib to visualize dataset.

0.1.1 Import packages and dataset

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Fri Dec 14 13:06:00 2018
5
6 @author: ali(zamanilai1995@gmail.com)
7 """
8 %%Import packages
9 import numpy as np # linear algebra
10 import seaborn as sns
11 sns.set(style='whitegrid')
12 import pandas as pd # data processing, CSV file I/O (e.g. pd.
    read_csv)
13 import matplotlib.pyplot as plt
14 %%Check datasret
15 import os
16 print(os.listdir("../dataSet"))
17 %%Read the data
18 print('Loading the dataset....')
19 credit_card = pd.read_csv('../dataSet/creditcard.csv')
20 print('Dataset shape: ',credit_card.shape)
21 print('Dataset was loaded!!!')
```

Output:

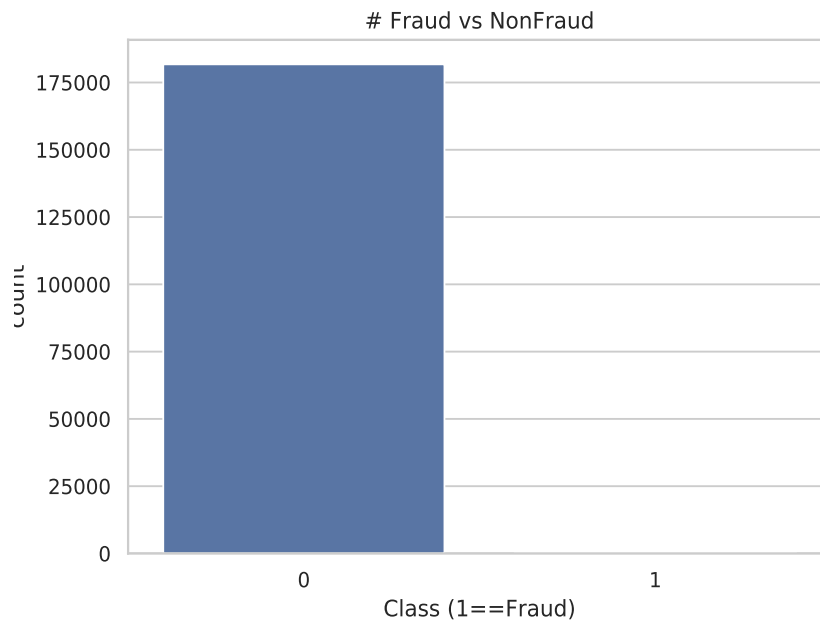
```
['creditcardfraud.zip', 'creditcard.csv', 'creditcard1.csv']
Loading the dataset....
Dataset shape: (182131, 31)
Dataset was loaded!!!
```

0.1.2 Balance of Data Visualization

Let's get a visual confirmation of the unbalanced data in this fraud dataset.

```
1 %%Plot fraud vs nonfraud
2 f, ax = plt.subplots(figsize=(7, 5))
3 sns.countplot(x='Class', data=credit_card)
4 _ = plt.title('# Fraud vs NonFraud')
5 _ = plt.xlabel('Class (1==Fraud)')
6 plt.savefig("fraudvsnonfraud"+" .pdf")
7 plt.show()
```

Output:

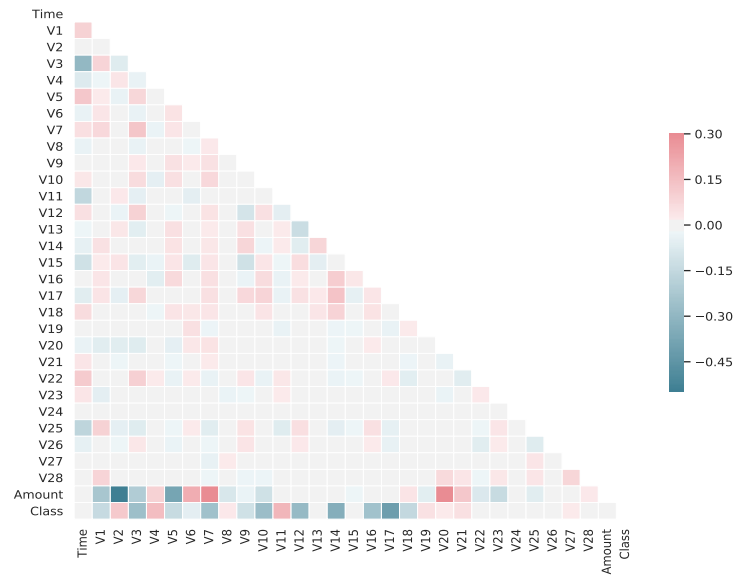


As you can see, the non-fraud cases strongly outweigh the fraud cases.

0.1.3 Heatmap

```
1 #%%Heatmap
2 corr=credit_card.corr()
3 mask = np.zeros_like(corr, dtype=np.bool)
4 mask[np.triu_indices_from(mask)] = True
5 cmap = sns.diverging_palette(220, 10, as_cmap=True)
6 # Draw the heatmap with the mask and correct aspect ratio
7 f, ax = plt.subplots(figsize=(11, 9))
8 sns.heatmap(corr, mask=mask, cmap=cmap, vmax=.3, center=0,
9             square=True, linewidths=.5, cbar_kws={"shrink": .5})
10 plt.savefig("heatmap"+" .pdf")
11 plt.show()
```

Output:



0.1.4 Fraud and non-fraud data describe

We will cut up the dataset into two data frames, one for non-fraud transactions and the other for fraud.

```
1 %%Fraud and non-fraud data describe
2 non_fraud = credit_card[credit_card.Class == 0]
3 fraud = credit_card[credit_card.Class == 1]
```

Let's look at some summary statistics and see if there are obvious differences between fraud and non-fraud transactions.

```
1 non_fraud.Amount.describe()
```

Output:

```
count    181766.000000
mean         88.435107
std        247.579620
min           0.000000
25%          5.780000
50%         22.400000
75%         78.000000
max        19656.530000
Name: Amount, dtype: float64
```

```
1 fraud.Amount.describe()
```

Output:

```
count      365.000000
mean       116.533205
std        249.276178
min         0.000000
25%         1.000000
50%        11.400000
75%        104.030000
max       2125.870000
Name: Amount, dtype: float64
```

Although the mean is a little higher in the fraud transactions, it is certainly within a standard deviation and so is unlikely to be easy to discriminate in a highly precise manner between the classes with pure statistical methods. I could run statistical tests (e.g. t-test) to support the claim that the two samples likely come from populations with similar means and deviations. However, such statistical methods are not the focus of this article on autoencoders.

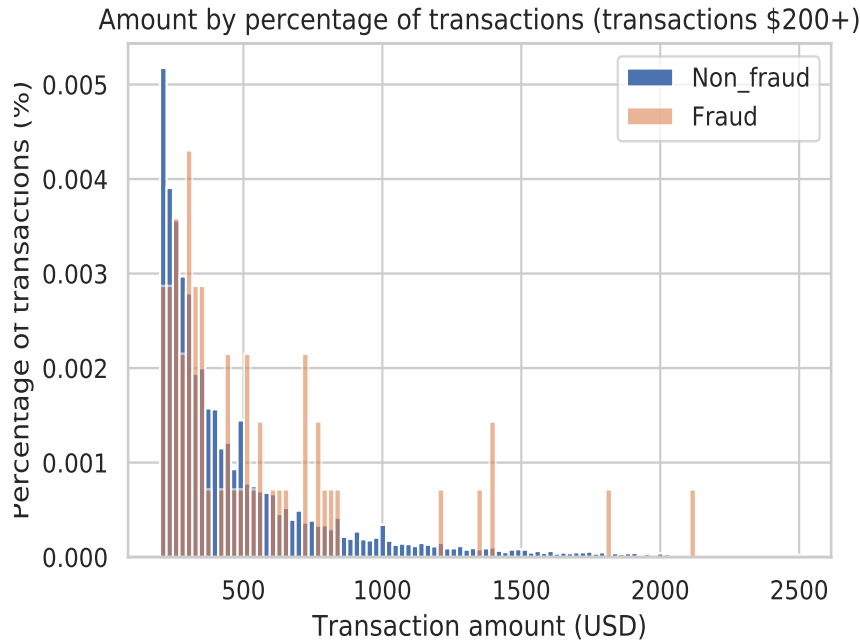
0.1.5 Visual Exploration of the Transaction Amount Data

We are going to get more familiar with the data and try some basic visuals. In anomaly detection datasets it is common to have the areas of interest "washed out" by abundant data. The most common method is to simply 'slice and dice' the data in a couple different ways until something interesting is found. Although this practice is common, it is not a scientifically sound way to explore data. There are always non-meaningful quirks to real data, so just looking until you "find something interesting" is likely going to result in you finding false positives. In other words, you find a random pattern in the current data set that will never be seen again. As a famous economist wrote, "If you torture the data long enough, it will confess."

In this dataset, I expect a lot of low-value transactions that will be generally uninteresting (buying cups of coffee, lunches, etc). This abundant data is likely to wash out the rest of the data, so I decided to look at the data in a number different \$100 and \$1,000 intervals. Since it would be tedious to show reader these graphs, I will only show the final graph that only visualizes the transactions above \$200.

```
1  #%%plot of high value transactions
2  bins = np.linspace(200, 2500, 100)
3  plt.hist(non_fraud.Amount, bins, alpha=1, normed=True, label='
   Non_fraud')
4  plt.hist(fraud.Amount, bins, alpha=0.6, normed=True, label='Fraud')
5  plt.legend(loc='upper right')
6  plt.title("Amount by percentage of transactions (transactions \ $200
   +)")
7  plt.xlabel("Transaction amount (USD)")
8  plt.ylabel("Percentage of transactions (%)");
9  plt.savefig("Amountbypercentageoftransactions"+" .pdf")
10 plt.show()
```

Output:



Since the fraud cases are relatively few in number compared to bin size, we see the data looks predictably more variable. In the long tail, especially, we are likely observing only a single fraud transaction. It would be hard to differentiate fraud from non-fraud transactions by transaction amount alone.

0.1.6 Visual Exploration of the Data by Hour

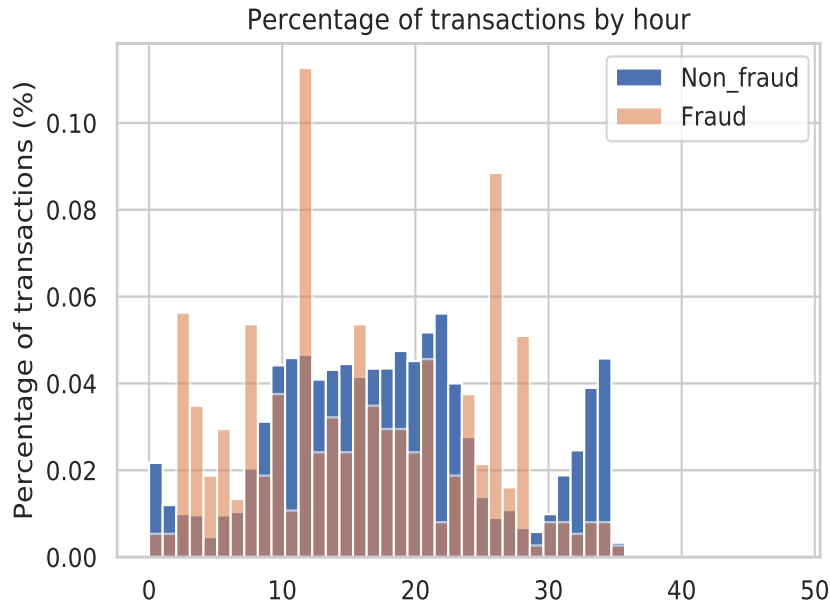
With a few exceptions, the transaction amount does not look very informative. Let's look at the time of day next.

```

1  ##Visual Exploration of the Data by Hour
2  bins = np.linspace(0, 48, 48) #48 hours
3  plt.hist((non_fraud.Time/(60*60)), bins, alpha=1, normed=True,
4           label='Non_fraud')
5  plt.hist((fraud.Time/(60*60)), bins, alpha=0.6, normed=True, label=
6           'Fraud')
7  plt.legend(loc='upper right')
8  plt.title("Percentage of transactions by hour")
9  plt.xlabel("Transaction time as measured from first transaction in
10             the dataset (hours)")
11 plt.ylabel("Percentage of transactions (%)");
12 plt.savefig("VisualExplorationoftheDataByHour"+" .pdf")
13 plt.show()

```

Output:



Transaction time as measured from first transaction in the dataset (hour
Hour "zero" corresponds to the hour the first transaction happened and not necessarily 12-1am. Given the heavy decrease in non-fraud transactions from hours 1 to 8 and again roughly at hours 24 to 32, I am assuming those time correspond to nighttime for this dataset. If this is true, fraud tends to occur at higher rates during the night. Statistical tests could be used to give evidence for this fact, but are not in the scope of this article. Again, however, the potential time offset between normal and fraud transactions is not enough to make a simple, precise classifier. Next, we will explore the potential interaction between transaction amount and hour to see if any patterns emerge.

0.1.7 Visual Exploration of Transaction Amount vs. Hour

```
1 #%%Visual Exploration of Transaction Amount vs. Hour
2 plt.scatter((non_fraud.Time/(60*60)), non_fraud.Amount, alpha=0.6,
3             label='non-fraud')
4 plt.scatter((fraud.Time/(60*60)), fraud.Amount, alpha=0.9, label='
5             Fraud')
6 plt.title("Amount of transaction by hour")
7 plt.xlabel("Transaction time as measured from first transaction in
8             the dataset (hours)")
9 plt.ylabel('Amount (USD)')
10 plt.legend(loc='upper right')
11 plt.savefig("VisualExplorationofTransactionAmountvsHour".pdf")
12 plt.show()
```

Output:

Again, this is not enough to make a good classifier. For example, it would be

hard to draw a line that cleanly separates fraud and normal transactions. For the experienced Data Scientists in the readership, I am excluding more advanced techniques such as the kernel trick.

0.2 Logistic model with sklearn

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Created on Thu Dec  6 13:10:34 2018
5
6  @author: ali
7  """
8  #%%
9  import numpy as np
10 import pandas as pd
11 from sklearn.model_selection import train_test_split
12 from sklearn.preprocessing import StandardScaler
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.pipeline import Pipeline
15 from sklearn.metrics import roc_curve, roc_auc_score,
16     classification_report, accuracy_score, confusion_matrix
17 import matplotlib.pyplot as plt
18 import os
19 print(os.listdir("../dataSet"))
20 credit_card = pd.read_csv('../dataSet/creditcard.csv')
21 #%%
22 X = credit_card.drop(columns='Class', axis=1)
23 y = credit_card.Class.values
24 #%%
25 np.random.seed(42)
26 X_train, X_test, y_train, y_test = train_test_split(X, y)
27 #%%
28 scaler = StandardScaler()
29 lr = LogisticRegression()
30 modell = Pipeline([('standardize', scaler),
31     ('log-reg', lr)])
32 modell.fit(X_train, y_train)
33 y_test_hat = modell.predict(X_test)
34 y_test_hat_probs = modell.predict_proba(X_test)[: ,1]
35 test_accuracy = accuracy_score(y_test, y_test_hat)*100
36 test_auc_roc = roc_auc_score(y_test, y_test_hat_probs)*100
37 print('Confusion matrix:\n', confusion_matrix(y_test, y_test_hat))
38 print('Training accuracy: %.4f %%' % test_accuracy)
39 print('Training AUC: %.4f %%' % test_auc_roc)
40 print(classification_report(y_test, y_test_hat, digits=6))
41 fpr, tpr, thresholds = roc_curve(y_test, y_test_hat_probs,
42     drop_intermediate=True)
43 f, ax = plt.subplots(figsize=(9, 6))
44 - = plt.plot(fpr, tpr, [0,1], [0, 1])
45 - = plt.title('AUC ROC')
46 - = plt.xlabel('False positive rate')
47 - = plt.ylabel('True positive rate')
48 plt.style.use('seaborn')
```

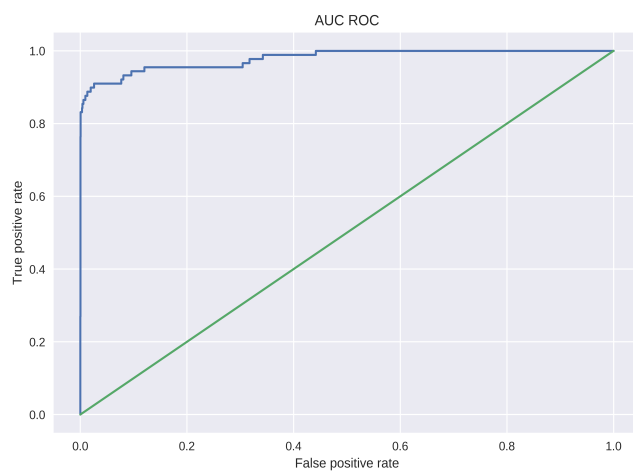


```

48 plt.savefig('auc_roc.png', dpi=600)
49 y_hat_90 = (y_test_hat_probs > 0.90 )*1
50 print('Confusion matrix for 90%:\n', confusion_matrix(y_test ,
    y_hat_90))
51 print('Report for 90%',classification_report(y_test , y_hat_90 ,
    digits=6))
52 y_hat_10 = (y_test_hat_probs > 0.05)*1
53 print('Confusion matrix for 5%:\n', confusion_matrix(y_test ,
    y_hat_10))
54 print('Report for 5%',classification_report(y_test , y_hat_10 ,
    digits=4))

```

Results:



```

Confusion matrix:
[[45433  11]
 [  33  56]]
Training accuracy: 99.9834 %
Training AUC: 97.8998 %
      precision    recall  f1-score   support

 0   0.999274   0.999758   0.999516   45444
 1   0.835821   0.629213   0.717949    89

 micro avg   0.999034   0.999034   0.999034   45533
 macro avg   0.917548   0.814486   0.858732   45533
weighted avg   0.998955   0.999034   0.998966   45533

Confusion matrix for 90%:
[[45436   8]
 [  47  42]]
Report for 90%
      precision    recall  f1-score   support

 0   0.998967   0.998824   0.999395   45444
 1   0.848000   0.471910   0.604317    89

 micro avg   0.998792   0.998792   0.998792   45533
 macro avg   0.919483   0.735867   0.801856   45533
weighted avg   0.998656   0.998792   0.998623   45533

Confusion matrix for 5%:
[[45421  23]
 [  16  73]]
Report for 5%
      precision    recall  f1-score   support

 0   0.9996   0.9995   0.9996   45444
 1   0.7604   0.8202   0.7892    89

 micro avg   0.9991   0.9991   0.9991   45533
 macro avg   0.8800   0.9099   0.8944   45533
weighted avg   0.9992   0.9991   0.9992   45533

```