

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(zamanilail1995@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:

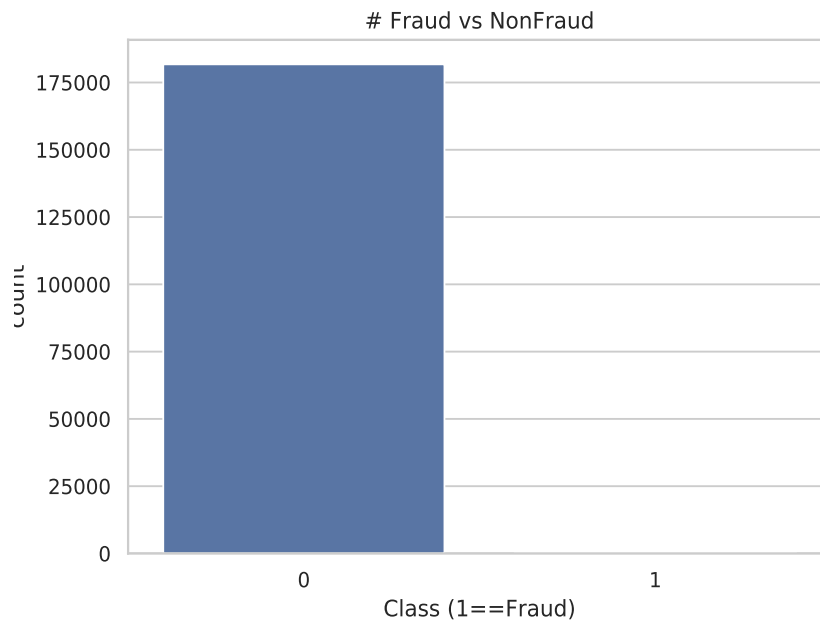
```
['creditcardfraud.zip', 'creditcard.csv', 'creditcard1.csv']
Loading the dataset....
Dataset shape: (182131, 31)
Dataset was loaded!!!
```

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:

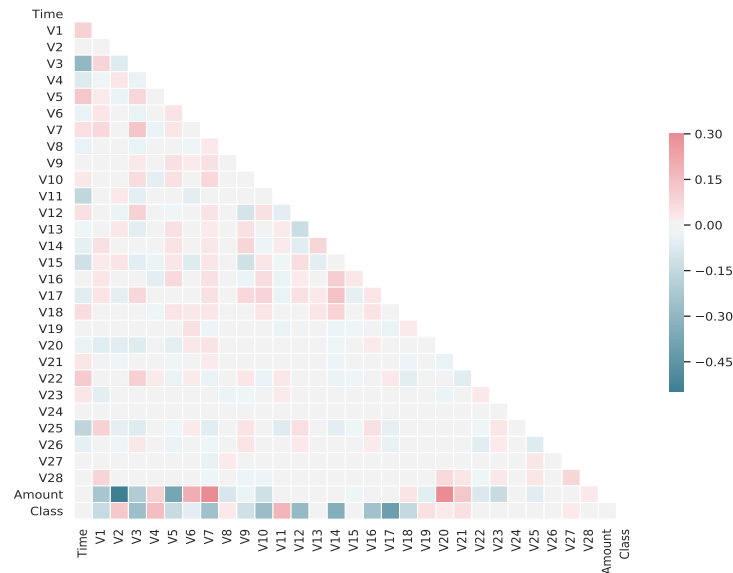


As you can see, the non-fraud cases strongly outweigh the fraud cases.

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.

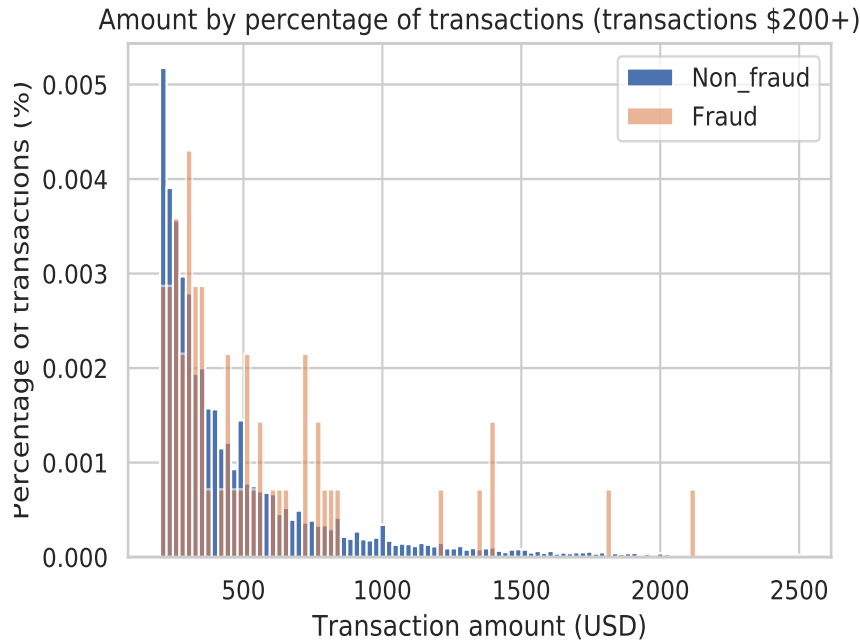
0.1.5 Visual Exploration of the Transaction Amount Data

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

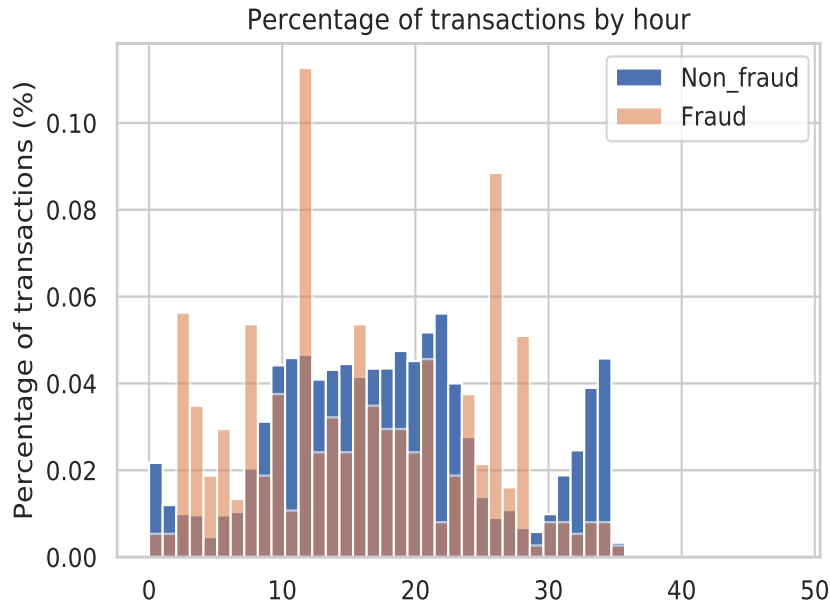
With a few exceptions, the transaction amount does not look very informative. Let's look at the time of day next.

```

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')
```

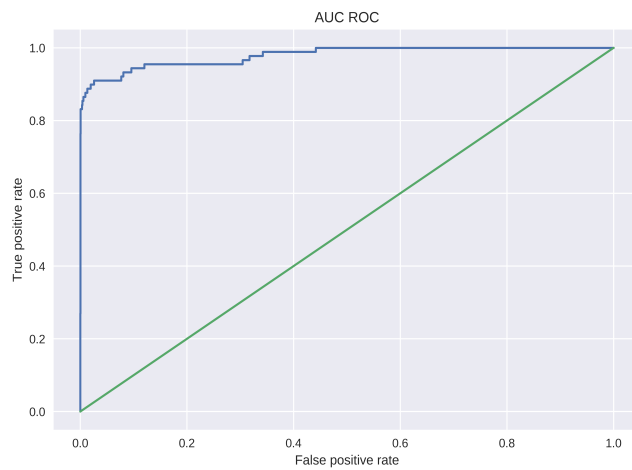


```

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,
51 y_hat_90))
52 print('Report for 90%', classification_report(y_test, y_hat_90,
53 digits=6))
54 y_hat_10 = (y_test_hat_probs > 0.05)*1
55 print('Confusion matrix for 5%:\n', confusion_matrix(y_test,
56 y_hat_10))
57 print('Report for 5%', classification_report(y_test, y_hat_10,
58 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

```

0.3 Logistic model with tensorflow

```

1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-

```

```

3 """
4 Created on Fri Dec 7 05:47:51 2018
5
6 @author: ali
7 """
8 ###Parameters
9 train_set_num=.8
10 seed=5
11 ## Define the learning rate batch_size etc.
12 learning_rate = 0.0003
13 batch_size = 1000
14 epoch_num = 600
15 ##Import packages
16 import numpy as np # linear algebra
17 import seaborn as sns
18 sns.set(style='whitegrid')
19 import pandas as pd # data processing, CSV file I/O (e.g. pd.
    read_csv)
20 import matplotlib.pyplot as plt
21 import tensorflow as tf
22 from sklearn.metrics import roc_curve, roc_auc_score,
    classification_report, accuracy_score, confusion_matrix
23 ## Define the normalized function
24 def min_max_normalized(data):
25     col_max = np.max(data, axis=0)
26     col_min = np.min(data, axis=0)
27     col_mean = np.mean(data, axis=0)
28     return np.divide(data - col_mean, col_max - col_min)
29 ##Check datasret
30 import os
31 print(os.listdir("../dataSet"))
32 ##Read the data
33 print('Loading the dataset....')
34 credit_card = pd.read_csv('../dataSet/creditcard.csv')
35 print('Dataset shape: ',credit_card.shape)
36 print('Dataset was loaded!!!')
37 ##
38 # set replace=False, Avoid double sampling
39 X = credit_card.drop(columns='Class', axis=1).values.reshape(-1,30)
40 y = credit_card.Class.values.reshape(-1,1)
41 train_index = np.random.choice(len(X), round(len(X) * train_set_num
    ),
42                                replace=False)
43 test_index = np.array(list(set(range(len(X))) - set(train_index)))
44 train_X = X[train_index]
45 train_y = y[train_index]
46 test_X = X[test_index]
47 test_y = y[test_index]
48 ## Normalized processing
49 train_X = min_max_normalized(train_X)
50 test_X = min_max_normalized(test_X)
51 ##Build the model framework
52 # Begin building the model framework
53 # Declare the variables that need to be learned and initialization
54 # There are 30 features here, A's dimension is (30, 1)
55 w = tf.Variable(tf.random.normal(shape=[30, 1]))
56 b = tf.Variable(tf.random.normal(shape=[1, 1]))

```

```

57 init = tf.global_variables_initializer()
58 sess = tf.Session()
59 sess.run(init)
60 # Define placeholders
61 x = tf.placeholder(dtype=tf.float32, shape=[None, 30], name="x")
62 y = tf.placeholder(dtype=tf.float32, shape=[None, 1], name="y")
63 # Define logistic Regression
64 logit = tf.matmul(x, w) + b
65 y_predicted = 1.0 / (1.0 + tf.exp(-logit))
66 # Declare loss function
67 loss = -1 * tf.reduce_sum(y * tf.log(y_predicted) +
68                             (1 - y) * tf.log(1 - y_predicted))
69 # Define optimizer: GradientDescent
70 optimizer = tf.train.GradientDescentOptimizer(
71     learning_rate=learning_rate).minimize(loss)
72 # Define the accuracy
73 # The default threshold is 0.5, rounded off directly
74 prediction = tf.round(tf.sigmoid(logit))
75 # Bool into float32 type
76 correct = tf.cast(tf.equal(prediction, y), dtype=tf.float32)
77 # Average
78 accuracy = tf.reduce_mean(correct)
79 # End of the definition of the model framework
80 #label=[tf.count_nonzero(y), tf.subtract(tf.size(y), tf.
81     count_nonzero(y))]
82 #confusion_matrix_tf = tf.confusion_matrix(labels=[10, 100],
83     predictions=[2, 108])
84 #FN=tf.metrics.false_negatives(labels=y, predictions=tf.round(
85     y_predicted))
86 conifution=np.zeros(shape=[2,2])
87 #%%
88 print("Parameters were initialized, Session is runing ...")
89 train_error_list = []
90 train_acc_list = []
91 test_acc_list = []
92 test_error_list=[]
93 with tf.Session() as sess:
94     sess.run(tf.global_variables_initializer())
95     for epoch in range(epoch_num):
96         train_loss = 0
97         for idx in range(len(train_X)//batch_size):
98             input_list = {x: train_X[idx*batch_size:(idx+1)*
99                 batch_size],
100                 y: train_y[idx*batch_size:(idx+1)*
101                     batch_size]}
102             -, train_loss1 = sess.run([optimizer, loss], feed_dict=
103                 input_list)
104             train_loss += train_loss1
105             train_error_list.append(train_loss/len(train_X))
106             train_acc_list.append(sess.run(
107                 accuracy, feed_dict={x: train_X, y:
108                     train_y})*100)
109             test_acc_list.append(sess.run(accuracy
110                 , feed_dict={x: test_X,
111                     y: test_y})*100)
112             test_error_list.append(sess.run(loss,

```

```

108                                     feed_dict={x: test_X,
109                                               y: test_y})/len(
110     test_y))
111     if (epoch + 1) % 50 == 0:
112         print('epoch: {:4d} loss: {:.5f} train_acc: {:.5f}%
113         test_acc: {:.5f}%'
114               .format(epoch + 1, train_loss/len(train_X),
115                       train_acc_list[epoch], test_acc_list[
116 epoch]))
117     w_value, b_value = sess.run([w, b])
118     for i in range(len(train_X)):
119         logit1 = np.matmul(train_X[i], w_value) + b_value
120         if (np.round(1.0 / (1.0 + np.exp(-logit1)))):
121             if (sess.run(tf.round(y_predicted), feed_dict={x: train_X[i
122 ]})):
123                 if train_y[i]:
124                     confuition[1,1]+=1
125                 else:
126                     confuition[0,1]+=1
127             else:
128                 if train_y[i]:
129                     confuition[1,0]+=1
130                 else:
131                     confuition[0,0]+=1
132     print ('Confusion matrix:', confuition)
133     #%%
134     train_y_hat=np.round(1.0/(1.0 + np.exp(-(np.matmul(train_X, w_value)
135 +b_value))))
136     test_y_hat=np.round(1.0/(1.0 + np.exp(-(np.matmul(test_X, w_value)+
137 b_value))))
138     test_accuracy = accuracy_score(test_y, test_y_hat)*100
139     test_auc_roc = roc_auc_score(test_y, test_y_hat)*100
140     print('Confusion matrix for train data:\n', confusion_matrix(test_y
141 ,
142     test_y_hat))
143     print('Confusion matrix for test data:\n', confusion_matrix(train_y
144 ,
145     train_y_hat))
146     print('Training accuracy: ', test_accuracy)
147     print('Training AUC: ', test_auc_roc)
148     print(classification_report(test_y, test_y_hat, digits=6))
149     fpr, tpr, thresholds = roc_curve(test_y, 1.0/(1.0 + np.exp(-(np.
150 matmul(test_X, w_value)+b_value)))
151 , drop_intermediate=True)
152     #select_tereshold=np.zeros_like(thresholds)
153     #recall=tpr/(tpr+fpr)
154     #precision=
155     #select_tereshold=2*(tpr*fpr)/(tpr+fpr)
156     #select_tereshold.append
157     f, ax = plt.subplots(figsize=(9, 6))
158     - = plt.plot(fpr, tpr, [0,1], [0, 1])
159     - = plt.title('AUC ROC')
160     - = plt.xlabel('False positive rate')
161     - = plt.ylabel('True positive rate')
162     plt.style.use('seaborn')

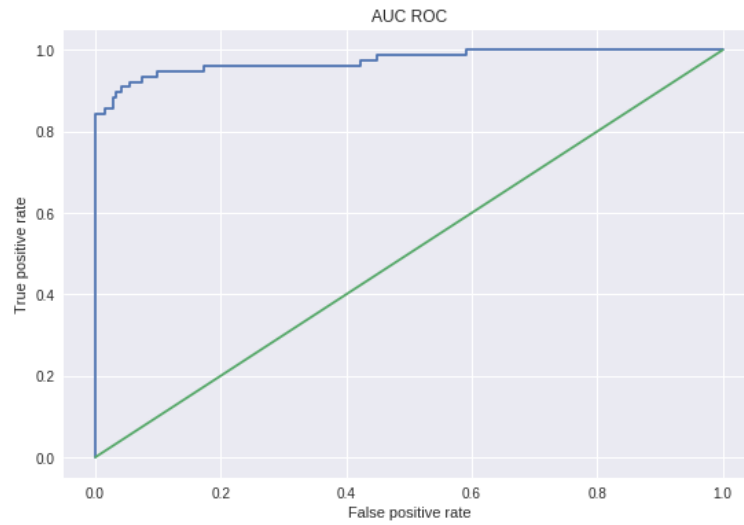
```

```

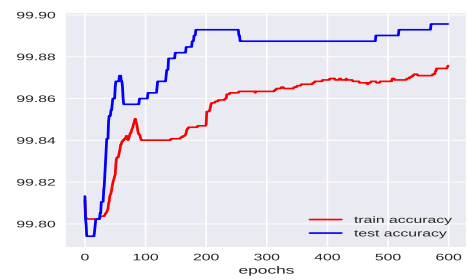
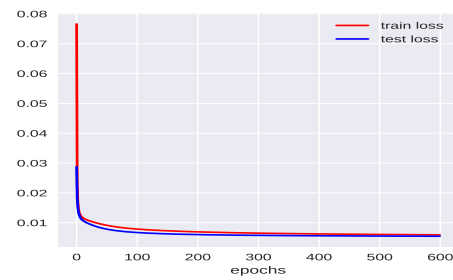
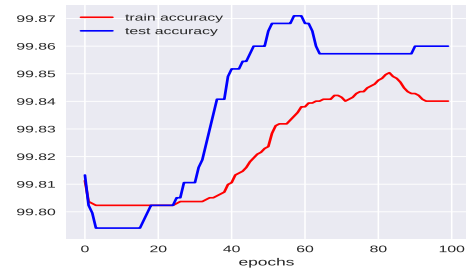
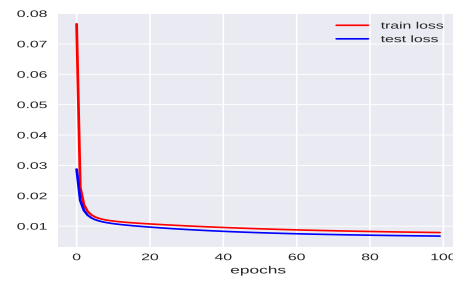
154 plt.savefig('auc-roc.png', dpi=600)
155 test_y_hat_10 = (1.0/(1.0 + np.exp(-(np.matmul(test_X, w_value)+
    b_value))) > 0.05 )*1
156 train_y_hat_10 = (1.0/(1.0 + np.exp(-(np.matmul(train_X, w_value)+
    b_value))) > 0.05 )*1
157 test_accuracy_10 = accuracy_score(test_y, test_y_hat_10)*100
158 test_auc_roc_10 = roc_auc_score(test_y, test_y_hat_10)*100
159 print('Confusion matrix for train data (0.05):\n', confusion_matrix(
    test_y,
160
    test_y_hat_10))
161 print('Confusion matrix for test data(0.05):\n', confusion_matrix(
    train_y,
162
    train_y_hat_10))
163 print('Training accuracy (0.05): ', test_accuracy_10)
164 print('Training AUC (0.05): ', test_auc_roc_10)
165 print(classification_report(test_y, test_y_hat_10, digits=6))
166 ##%
167 fig, ax = plt.subplots(2, 2, figsize=(10, 10))
168 fig.suptitle("test accuracy = " + str(test_acc_list[epoch]))
169 for a in ax.reshape(-1,1):
170     a[0].set_xlabel("epochs")
171 ax[0][0].plot(train_error_list[:100], color='red', label='train
    loss')
172 ax[0][0].plot(test_error_list[:100], color='blue', label='test loss
    ')
173 ax[0][0].legend()
174 ax[1][0].plot(train_error_list, color='red', label='train loss')
175 ax[1][0].plot(test_error_list, color='blue', label='test loss')
176 ax[1][0].legend()
177 ax[0][1].plot(train_acc_list[:100], color='red', label='train
    accuracy')
178 ax[0][1].plot(test_acc_list[:100], color='blue', label='test
    accuracy')
179 ax[0][1].legend()
180 ax[1][1].plot(train_acc_list, color='red', label='train accuracy')
181 ax[1][1].plot(test_acc_list, color='blue', label='test accuracy')
182 ax[1][1].legend()
183 plt.savefig("trainandtest"+" .pdf")
184 #End main program

```

Results:



test accuracy = 99.89567995071411



```
[creditcardfraud.zip, 'creditcard.csv', 'creditcard1.csv']
loading the dataset...
Dataset shape: (182131, 31)
Dataset was loaded!!!
Parameters were initialized, Session is running...
epoch: 50 loss: 0.002173 train_acc: 99.82232% test_acc: 99.859989%
epoch: 100 loss: 0.007893 train_acc: 99.840807% test_acc: 99.859989%
epoch: 150 loss: 0.007308 train_acc: 99.84077% test_acc: 99.881953%
epoch: 200 loss: 0.006951 train_acc: 99.846953% test_acc: 99.892932%
epoch: 250 loss: 0.006701 train_acc: 99.86276% test_acc: 99.892932%
epoch: 300 loss: 0.006516 train_acc: 99.863422% test_acc: 99.887443%
epoch: 350 loss: 0.006372 train_acc: 99.863408% test_acc: 99.887443%
epoch: 400 loss: 0.006256 train_acc: 99.868912% test_acc: 99.887443%
epoch: 450 loss: 0.006168 train_acc: 99.868226% test_acc: 99.887443%
epoch: 500 loss: 0.006088 train_acc: 99.868226% test_acc: 99.898198%
epoch: 550 loss: 0.006018 train_acc: 99.878974% test_acc: 99.892932%
epoch: 600 loss: 0.005958 train_acc: 99.873778% test_acc: 99.892688%
```

```

Confusion matrix: [[1.45384e+05 3.36800e+01]
 [1.48000e+02 1.40000e+02]]
Confusion matrix for train data:
[[36342  77]
 [ 31  46]]
Confusion matrix for test data:
[[145384  33]
 [ 148  140]]
Training accuracy: 99.89567891066821
Training AUC: 79.86050099450743
      precision    recall  f1-score   support

0     0.999148    0.999807    0.999477    36349
1     0.867925    0.597403    0.707692      77

micro avg     0.998957    0.998957    0.998957    36426
macro avg     0.933536    0.798605    0.853583    36426
weighted avg     0.998870    0.998957    0.998861    36426

Confusion matrix for train data (0.05):
[[36338  11]
 [ 12  65]]
Confusion matrix for test data(0.05):
[[145374  43]
 [  91 197]]
Training accuracy (0.05): 99.93685828803602
Training AUC (0.05): 92.1926611752836
      precision    recall  f1-score   support

0     0.999670    0.999627    0.999644    36349
1     0.853263    0.841356    0.849673      77

micro avg     0.999369    0.999369    0.999369    36426
macro avg     0.927467    0.921927    0.924678    36426
weighted avg     0.999365    0.999369    0.999367    36426

```

0.4 compare sklearn and tensorflow results for logistic model

Accuracy model of sklearn and tensorflow are 99.9% and 99.89% but we know because of unbalance data, accuracy is not suitable criteria hence we compare recall(percentage of one data is fraud and our model predict it fraud) and precision(percentage of one data is non-fraud and our model predicts it non-fraud) and f1-score(recall and precision) of two packages.

Recall:

sklearn:82%

tensorflow:84%

Precision:

sklearn:76%

tensorflow:85%

F1-score:

sklearn:79%

tensorflow:85%

We see that tensorflow predict better.

0.5 SVM with sklearn

```

1  ###Parameters
2  train_set_num=.8
3  seed=5
4  ###Import packages
5  import numpy as np # linear algebra
6  import pandas as pd # data processing, CSV file I/O (e.g. pd.
   read_csv)
7  import matplotlib.pyplot as plt
8  from sklearn import svm
9  from sklearn.svm import SVC
10 from sklearn.metrics import classification_report
11 from sklearn.metrics import confusion_matrix
12 from sklearn.metrics import f1_score
13 from sklearn.metrics import accuracy_score
14 from sklearn.metrics import recall_score
15 from sklearn.metrics import precision_score, precision_recall_curve

```

```

16 from sklearn.metrics import roc_auc_score, roc_curve, auc,
    average_precision_score
17 from sklearn.preprocessing import StandardScaler
18 from sklearn.model_selection import train_test_split
19 from mlxtend.plotting import plot_confusion_matrix
20 import seaborn as sns
21 import warnings
22 warnings.filterwarnings('ignore')
23
24 ### Define the normalized function
25 def min_max_normalized(data):
26     col_max = np.max(data, axis=0)
27     col_min = np.min(data, axis=0)
28     col_mean = np.mean(data, axis=0)
29     return np.divide(data - col_min, col_max - col_min)
30
31 ### Check dataset
32 import os
33 print(os.listdir("../dataSet"))
34 ### Read the data
35 print('Loading the dataset....')
36 credit_card = pd.read_csv('../dataSet/creditcard.csv')
37 print('Dataset shape: ', credit_card.shape)
38 print('Dataset was loaded!!!')
39
40 # set replace=False, Avoid double sampling
41 X = credit_card.drop(columns='Class', axis=1).values.reshape(-1,30)
42 y = credit_card.Class.values.reshape(-1,1)
43 train_index = np.random.choice(len(X), round(len(X) * train_set_num
    ),
44                                replace=False)
45 test_index = np.array(list(set(range(len(X))) - set(train_index)))
46 train_X = X[train_index]
47 train_y = y[train_index]
48 test_X = X[test_index]
49 test_y = y[test_index]
50
51 ### Normalized processing
52 train_X = min_max_normalized(train_X)
53 test_X = min_max_normalized(test_X)
54
55 # Applying SVM Algorithm
56 print("_____")
57 print("                                Support Vector Machine
    ")
58 print("_____")
59
60 #Using the rbf kernel to build the initail model.
61 classifier= svm.SVC(C= 1, kernel= 'linear', random_state= 0)
62
63 #Fit into Model
64 classifier.fit(train_X, train_y)
65
66 #Predict the class using X_test
67 y_pred = classifier.predict(test_X)
68
69 con_mat = confusion_matrix(test_y, y_pred)
70 average_precision = average_precision_score(test_y, y_pred)
71 cls_report = classification_report(test_y, y_pred)

```

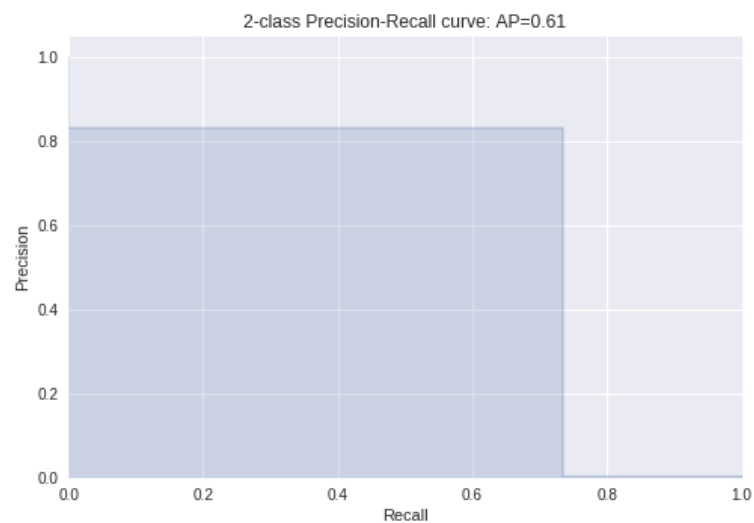
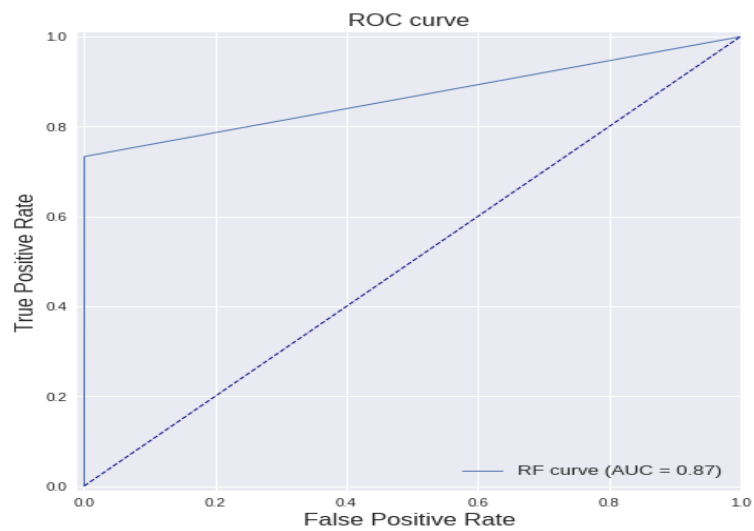


```

70
71 print("
    *****
    ")
72 print("Area under the curve : %f" % (roc_auc_score(test_y, y_pred))
    )
73 print("Average precision-recall score RF: {}".format(
    average_precision))
74 print(con_mat)
75 print(cls_report)
76 print("
    *****
    ")
77 #%%
78 precision, recall, _ = precision_recall_curve(test_y, y_pred)
79 plt.step(recall, precision, color='b', alpha=0.2, where='post')
80 plt.fill_between(recall, precision, step='post', alpha=0.2, color='
    b')
81
82 plt.xlabel('Recall')
83 plt.ylabel('Precision')
84 plt.ylim([0.0, 1.05])
85 plt.xlim([0.0, 1.0])
86 plt.title('2-class Precision-Recall curve: AP={0:0.2f}'.format(
    average_precision))
87
88 fpr_rf, tpr_rf, _ = roc_curve(test_y, y_pred)
89 roc_auc_rf = auc(fpr_rf, tpr_rf)
90 plt.figure(figsize=(8,8))
91 plt.xlim([-0.01, 1.00])
92 plt.ylim([-0.01, 1.01])
93 plt.plot(fpr_rf, tpr_rf, lw=1, label='{} curve (AUC = {:.0.2f})'.
    format('RF', roc_auc_rf))
94
95 plt.xlabel('False Positive Rate', fontsize=16)
96 plt.ylabel('True Positive Rate', fontsize=16)
97 plt.title('ROC curve', fontsize=16)
98 plt.legend(loc='lower right', fontsize=13)
99 plt.plot([0, 1], [0, 1], color='navy', lw=1, linestyle='--')
100 plt.axes().set_aspect('equal')
101 plt.show()

```

Results:



```

Support Vector Machine
-----
Area under the curve : 0.86515
Average precision-recall score RF: 0.6116091094769152
[[36340 11]
 [ 20 55]]
precision recall f1-score support
0 1.00 1.00 1.00 36351
1 0.83 0.73 0.78 75
micro avg 1.00 1.00 1.00 36426
macro avg 0.92 0.87 0.89 36426
weighted avg 1.00 1.00 1.00 36426

```

0.6 SVM with tensorflow

0.7 Compare sklearn and tensorflow results for SVM