

Credit card fraud detection

Table of contents

1. Defining the problem statement
2. Imports
3. EDA
4. Model prediction

1. Defining the problem statement

- The dataset contains transactions made by credit cards in September 2013 by European cardholders. This dataset presents transactions that occurred in two days, where we have 492 frauds out of 284,807 transactions. The dataset is highly unbalanced, the positive class (frauds) account for 0.172% of all transactions.

□



2. Imports

In [3]:

```
import numpy as np
import pandas as pd
import sklearn
import scipy
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import classification_report, accuracy_score
from sklearn.ensemble import IsolationForest
from sklearn.neighbors import LocalOutlierFactor
from sklearn.svm import OneClassSVM
from pylab import rcParams
rcParams['figure.figsize'] = 14, 8
RANDOM_SEED = 42
LABELS = ["Normal", "Fraud"]
```

2.1 Collecting the data

In [4]:

```
data = pd.read_csv('creditcard.csv')
data.head()
```

Out[4]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23
0	0.0	1.359807	0.072781	2.536347	1.378155	0.338321	0.462388	0.239599	0.098698	0.363787	...	0.018307	0.277838	0.110474
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	0.082361	0.078803	0.085102	0.255425	...	0.225775	0.638672	0.101288
2	1.0	1.358354	1.340163	1.773209	0.379780	0.503198	1.800499	0.791461	0.247676	1.514654	...	0.247998	0.771679	0.909412
3	1.0	0.966272	0.185226	1.792993	0.863291	0.010309	1.247203	0.237609	0.377436	1.387024	...	0.108300	0.005274	0.190321
4	2.0	1.158233	0.877737	1.548718	0.403034	0.407193	0.095921	0.592941	0.270533	0.817739	...	0.009431	0.798278	0.137458

5 rows × 31 columns



In [3]:

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 284807 entries, 0 to 284806
Data columns (total 31 columns):
Time                284807 non-null float64
V1                  284807 non-null float64
V2                  284807 non-null float64
V3                  284807 non-null float64
V4                  284807 non-null float64
V5                  284807 non-null float64
V6                  284807 non-null float64
V7                  284807 non-null float64
V8                  284807 non-null float64
V9                  284807 non-null float64
V10                 284807 non-null float64
V11                 284807 non-null float64
V12                 284807 non-null float64
V13                 284807 non-null float64
V14                 284807 non-null float64
V15                 284807 non-null float64
V16                 284807 non-null float64
V17                 284807 non-null float64
V18                 284807 non-null float64
V19                 284807 non-null float64
V20                 284807 non-null float64
V21                 284807 non-null float64
V22                 284807 non-null float64
V23                 284807 non-null float64
V24                 284807 non-null float64
V25                 284807 non-null float64
V26                 284807 non-null float64
V27                 284807 non-null float64
V28                 284807 non-null float64
Amount              284807 non-null float64
Class               284807 non-null int64
dtypes: float64(30), int64(1)
memory usage: 67.4 MB
```

In [4]:

```
data.describe()
```

Out[4]:

	Time	V1	V2	V3	V4	V5	V6	V7	V8
count	284807.000000	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	94813.859575	3.919560e-15	5.688174e-16	-8.769071e-16	2.782312e-15	-1.552563e-15	2.010663e-15	-1.694249e-15	-1.927028e-15

	Time	V1	V2	V3	V4	V5	V6	V7	V8
std	47488.145955	1.958696e+00	1.651309e+00	1.516255e+00	1.415869e+00	1.380247e+00	1.332271e+00	1.237094e+00	1.194353e+00
min	0.000000	5.640751e+01	7.271573e+01	4.832559e+01	5.683171e+00	1.137433e+02	2.616051e+01	4.355724e+01	7.321672e+00
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903648e-01	-8.486401e-01	-6.915971e-01	-7.682956e-01	-5.540759e-01	-2.086297e-01
50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-01	-1.984653e-02	-5.433583e-02	-2.741871e-01	4.010308e-02	2.235804e-01
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01	6.119264e-01	3.985649e-01	5.704361e-01	3.273459e-01
max	172792.000000	2.454930e+00	2.205773e+01	9.382558e+00	1.687534e+01	3.480167e+01	7.330163e+01	1.205895e+02	2.000721e+01

8 rows × 31 columns

3.Exploratory Data Analysis

In [5]:

```
#Any missing values
data.isnull().values.any()
```

Out[5]:

False

In [6]:

```
# Determine the number of fraud and normal transactions in the entire dataset using bar graph

count_classes = pd.value_counts(data['Class'], sort = True)

count_classes.plot(kind = 'bar', rot=0)

plt.title("Transaction Class Distribution")

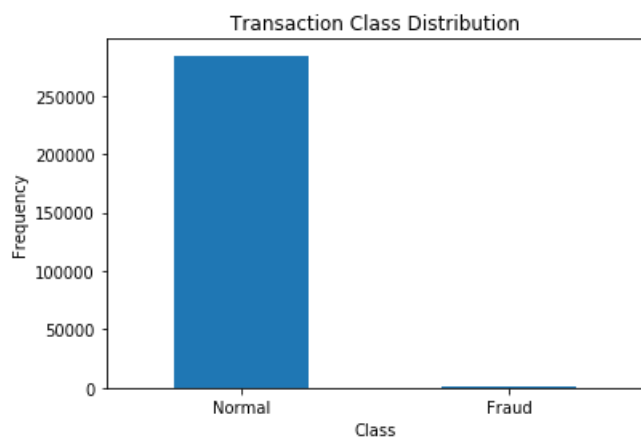
plt.xticks(range(2), LABELS)

plt.xlabel("Class")

plt.ylabel("Frequency")
```

Out[6]:

Text(0, 0.5, 'Frequency')



How many are fraud and how many are not fraud ?

In [8]:

```
class_names = {0:'Not Fraud', 1:'Fraud'}
```

```
print(data.Class.value_counts().rename(index = class_names))
```

```
Not Fraud    284315
Fraud         492
Name: Class, dtype: int64
```

In [9]:

```
fig = plt.figure(figsize = (15, 12))
```

<Figure size 1080x864 with 0 Axes>

In [10]:

```
## Get the Fraud and the normal dataset
```

```
fraud = data[data['Class']==1]
```

```
normal = data[data['Class']==0]
```

In [11]:

```
print(fraud.shape,normal.shape)
```

```
(492, 31) (284315, 31)
```

In [12]:

```
## We need to analyze more amount of information from the transaction data
#How different are the amount of money used in different transaction classes?
fraud.Amount.describe()
```

Out[12]:

```
count    492.000000
mean      122.211321
std       256.683288
min         0.000000
25%         1.000000
50%         9.250000
75%        105.890000
max       2125.870000
Name: Amount, dtype: float64
```

In [13]:

```
normal.Amount.describe()
```

Out[13]:

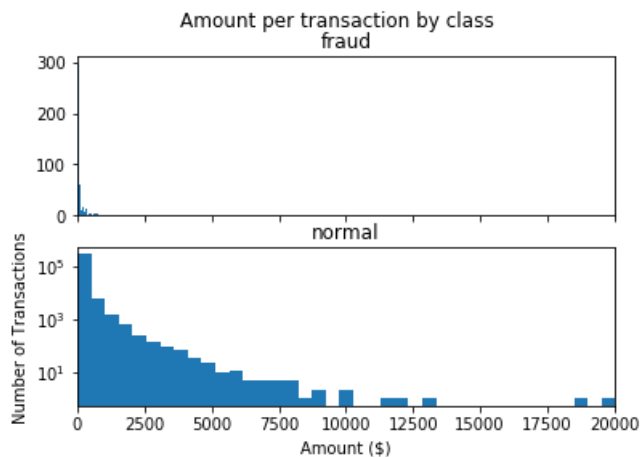
```
count    284315.000000
mean         88.291022
std       250.105092
min         0.000000
25%         5.650000
50%        22.000000
75%        77.050000
max       25691.160000
Name: Amount, dtype: float64
```

In [14]:

```
#Graphical representation of the data

f, (ax1, ax2) = plt.subplots(2, 1, sharex=True)
f.suptitle('Amount per transaction by class')
bins = 50
ax1.hist(fraud.Amount, bins = bins)
```

```
ax1.set_title('fraud')
ax2.hist(normal.Amount, bins = bins)
ax2.set_title('normal')
plt.xlabel('Amount ($)')
plt.ylabel('Number of Transactions')
plt.xlim((0, 20000))
plt.yscale('log')
plt.show();
```



In [15]:

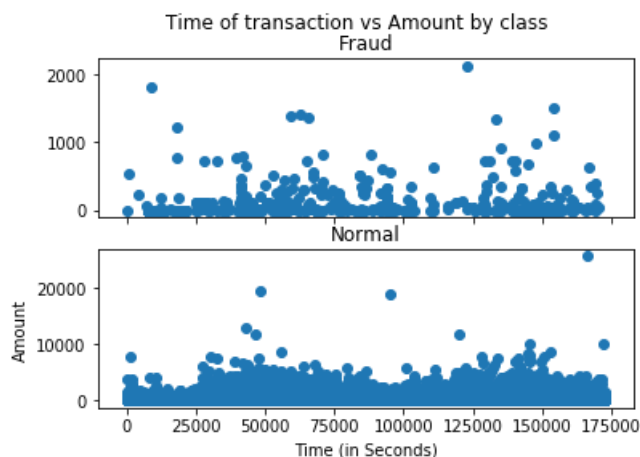
```
Normal=data[data['Class']==0]
```

In [17]:

```
Fraud = data[data['Class']==1]
```

In [18]:

```
# Visual representation showing whether fraudulent transactions occur more often during certain ti
me frame
f, (ax1, ax2) = plt.subplots(2, 1, sharex=True)
f.suptitle('Time of transaction vs Amount by class')
ax1.scatter(Fraud.Time, Fraud.Amount)
ax1.set_title('Fraud')
ax2.scatter(Normal.Time, Normal.Amount)
ax2.set_title('Normal')
plt.xlabel('Time (in Seconds)')
plt.ylabel('Amount')
plt.show()
```



- Doesn't seem like the time of transaction really matters here as per above observation

In [19]:

```
## Take some sample of the data

data1= data.sample(frac = 0.1,random_state=1)

data1.shape
```

```
Out[19]:

(28481, 31)
```

```
In [20]:
```

```
data.shape
```

```
Out[20]:

(284807, 31)
```

```
In [21]:
```

```
#Determine the number of fraud and valid transactions in the dataset

Fraud = data1[data1['Class']==1]

Valid = data1[data1['Class']==0]

outlier_fraction = len(Fraud)/float(len(Valid))
```

```
In [22]:
```

```
#Print the outlier fraction and number of Fraud and Valid Transaction cases
print(outlier_fraction)

print("Fraud Cases : {}".format(len(Fraud)))

print("Valid Cases : {}".format(len(Valid)))
```

```
0.0017234102419808666
Fraud Cases : 49
Valid Cases : 28432
```

```
In [23]:
```

```
#Get all the columns from the dataframe#
#Create independent and Dependent Features
columns = data1.columns.tolist()
# Filter the columns to remove data we do not want
columns = [c for c in columns if c not in ["Class"]]
# Store the variable we are predicting
target = "Class"
# Define a random state
state = np.random.RandomState(42)
X = data1[columns]
Y = data1[target]
X_outliers = state.uniform(low=0, high=1, size=(X.shape[0], X.shape[1]))
# Print the shapes of X & Y
print(X.shape)
print(Y.shape)
```

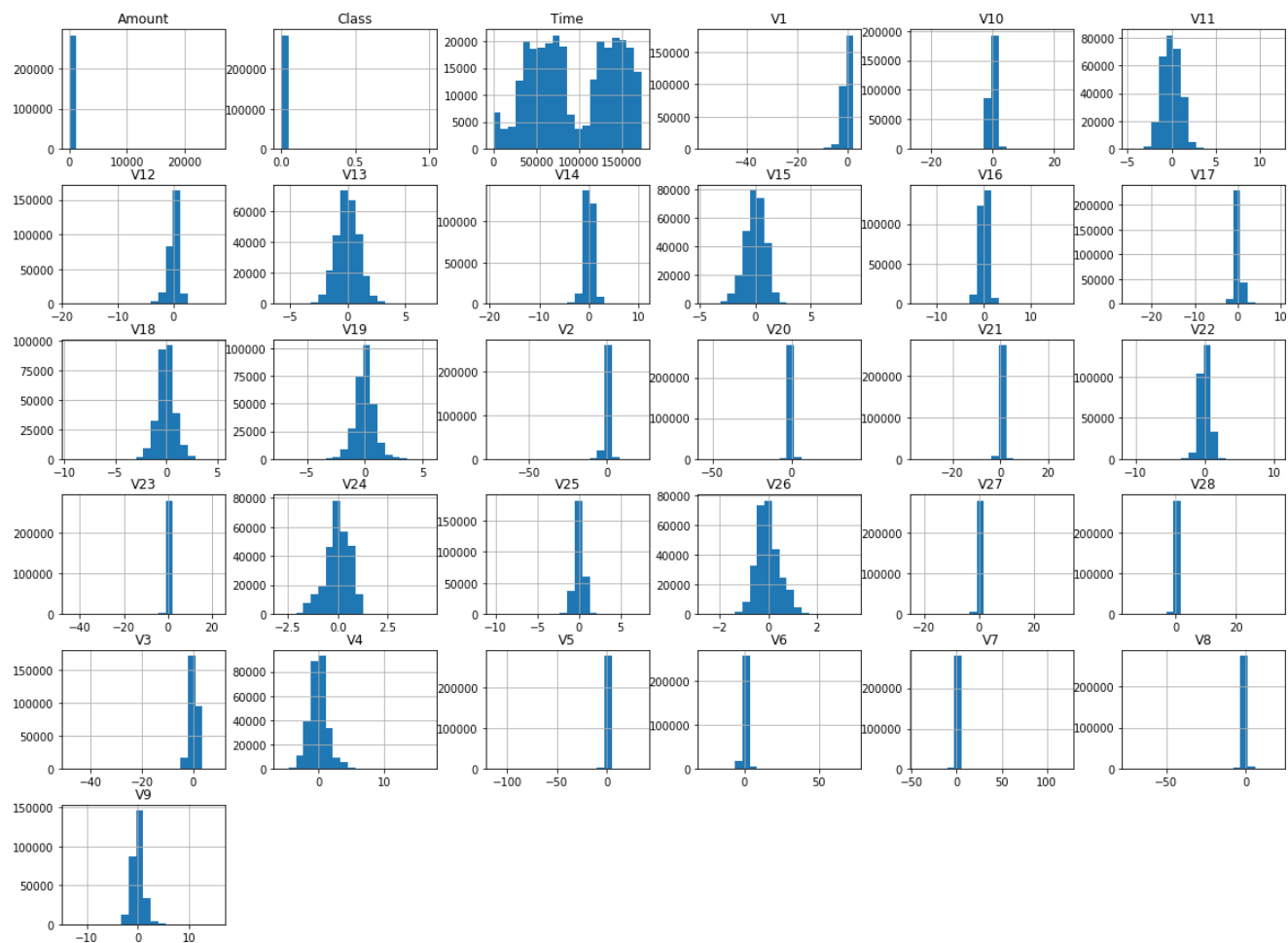
```
(28481, 30)
(28481,)
```

Plot histogram of each parameter

```
In [5]:
```

```
# Histograms of the features
# most of the data has a quasi-normal/gaussian distribution
```

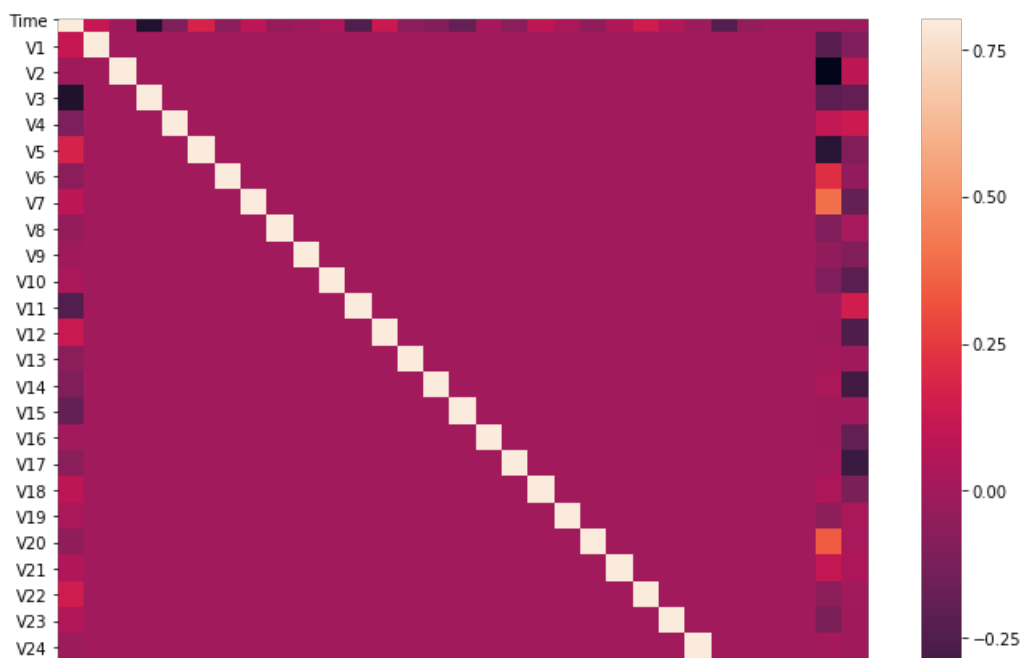
```
data.hist(bins=20, figsize=(20,15))
plt.show()
```

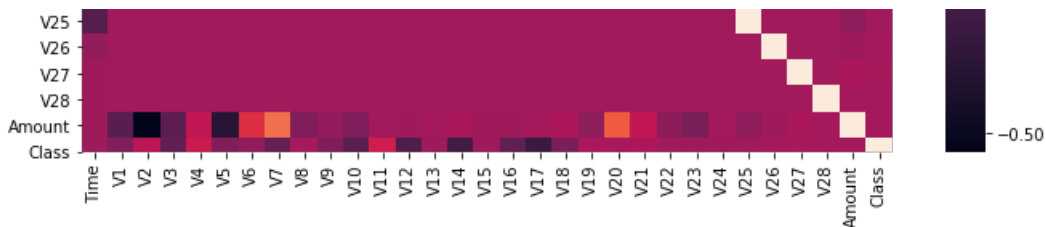


Correlation Matrix

In [24]:

```
corrmat = data.corr()
fig = plt.figure(figsize = (12, 9))
sns.heatmap(corrmat, vmax = .8, square = True)
plt.show()
```





- The above correlation matrix shows that none of the V1 to V28 PCA components have any correlation to each other however if we observe Class has some form positive and negative correlations with the V components but has no correlation with Time and Amount.

4. Model prediction

In [38]:

```
##Define the outlier detection methods

classifiers = {
    "Isolation Forest":IsolationForest(n_estimators=100, max_samples=len(X),
                                       contamination=outlier_fraction, verbose=0),
    "Local Outlier Factor":LocalOutlierFactor(n_neighbors=20, algorithm='auto',
                                             leaf_size=30, metric='minkowski',
                                             p=2, metric_params=None, contamination=outlier_fracti
),
    "Support Vector Machine":OneClassSVM(kernel='rbf', degree=3, gamma=0.1,nu=0.05,
                                       max_iter=-1)
}
```

In [39]:

```
type(classifiers)
```

Out[39]:

```
dict
```

In [40]:

```
n_outliers = len(Fraud)
for i, (clf_name,clf) in enumerate(classifiers.items()):
    #Fit the data and tag outliers
    if clf_name == "Local Outlier Factor":
        y_pred = clf.fit_predict(X)
        scores_prediction = clf.negative_outlier_factor_
    elif clf_name == "Support Vector Machine":
        clf.fit(X)
        y_pred = clf.predict(X)
    else:
        clf.fit(X)
        scores_prediction = clf.decision_function(X)
        y_pred = clf.predict(X)
    #Reshape the prediction values to 0 for Valid transactions , 1 for Fraud transactions
    y_pred[y_pred == 1] = 0
    y_pred[y_pred == -1] = 1
    n_errors = (y_pred != Y).sum()
    # Run Classification Metrics
    print("{}: {}".format(clf_name,n_errors))
    print("Accuracy Score :")
    print(accuracy_score(Y,y_pred))
    print("Classification Report :")
    print(classification_report(Y,y_pred))
```

```
Isolation Forest: 73
```

```
Accuracy Score :
```

```
0.9974368877497279
```

```
Classification Report :
```

```
precision    recall  f1-score   support
```


	precision	recall	f1-score	support
0	1.00	1.00	1.00	28432
1	0.26	0.27	0.26	49
accuracy			1.00	28481
macro avg	0.63	0.63	0.63	28481
weighted avg	1.00	1.00	1.00	28481

Local Outlier Factor: 97

Accuracy Score :

0.9965942207085425

Classification Report :

	precision	recall	f1-score	support
0	1.00	1.00	1.00	28432
1	0.02	0.02	0.02	49
accuracy			1.00	28481
macro avg	0.51	0.51	0.51	28481
weighted avg	1.00	1.00	1.00	28481

Support Vector Machine: 8516

Accuracy Score :

0.7009936448860644

Classification Report :

	precision	recall	f1-score	support
0	1.00	0.70	0.82	28432
1	0.00	0.37	0.00	49
accuracy			0.70	28481
macro avg	0.50	0.53	0.41	28481
weighted avg	1.00	0.70	0.82	28481

Observations :

- Isolation Forest detected 73 errors versus Local Outlier Factor detecting 97 errors vs. SVM detecting 8516 errors
- Isolation Forest has a 99.74% more accurate than LOF of 99.65% and SVM of 70.09
- When comparing error precision & recall for 3 models , the Isolation Forest performed much better than the LOF as we can see that the detection of fraud cases is around 27 % versus LOF detection rate of just 2 % and SVM of 0%.
- So overall Isolation Forest Method performed much better in determining the fraud cases which is around 30%.
- We can also improve on this accuracy by increasing the sample size or use deep learning algorithms however at the cost of computational expense. We can also use complex anomaly detection models to get better accuracy in determining more fraudulent cases

In []:

In []:

In []:

In []:

In []:

In []:

