# Credit Card Kaggle Anamoly Detection ¶

## Context

It is important that credit card companies are able to recognize fraudulent credit card transactions so that customers are not charged for items that they did not purchase.

## Content

The datasets 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.

It contains only numerical input variables which are the result of a PCA transformation. Unfortunately, due to confidentiality issues, we cannot provide the original features and more background information about the data. Features V1, V2, ... V28 are the principal components obtained with PCA, the only features which have not been transformed with PCA are 'Time' and 'Amount'. Feature 'Time' contains the seconds elapsed between each transaction and the first transaction in the dataset. The feature 'Amount' is the transaction Amount, this feature can be used for example-dependant cost-senstive learning. Feature 'Class' is the response variable and it takes value 1 in case of fraud and 0 otherwise.

## Inspiration

Identify fraudulent credit card transactions.

Given the class imbalance ratio, we recommend measuring the accuracy using the Area Under the Precision-Recall Curve (AUPRC). Confusion matrix accuracy is not meaningful for unbalanced classification.

## Acknowledgements

The dataset has been collected and analysed during a research collaboration of Worldline and the Machine Learning Group (http://mlg.ulb.ac.be (http://mlg.ulb.ac.be)) of ULB (Université Libre de Bruxelles) on big data mining and fraud detection. More details on current and past projects on related topics are available on https://www.researchgate.net/project/Fraud-detection-5 (https://www.researchgate.net/project/Fraud-detection-5) and the page of the DefeatFraud project

```
In [2]: 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"]
```

```
In [3]:  data = pd.read_csv('creditcard.csv',sep=',')
         data.head()
```

Out[3]:

| | Time | V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 |
| 1 | 0.0 | 1.191857 | 0.266151 | 0.166480 | 0.448154 | 0.060018 | -0.082361 | -0.078803 | 0.085102 |
| 2 | 1.0 | -1.358354 | -1.340163 | 1.773209 | 0.379780 | -0.503198 | 1.800499 | 0.791461 | 0.247676 |
| 3 | 1.0 | -0.966272 | -0.185226 | 1.792993 | -0.863291 | -0.010309 | 1.247203 | 0.237609 | 0.377436 |
| 4 | 2.0 | -1.158233 | 0.877737 | 1.548718 | 0.403034 | -0.407193 | 0.095921 | 0.592941 | -0.270533 |

5 rows × 31 columns

```
In [4]:  data.info()
```
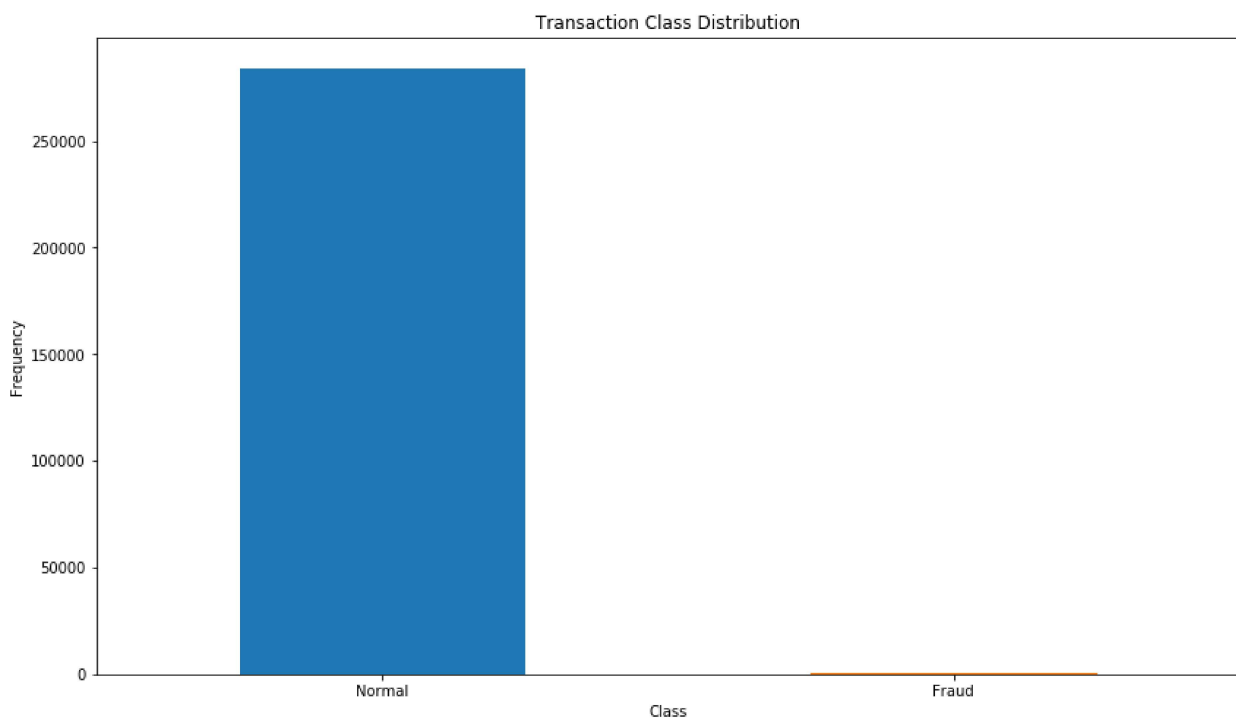
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 284807 entries, 0 to 284806
Data columns (total 31 columns):
Time      284807 non-null float64
V1        284807 non-null float64
V2        284807 non-null float64
V3        284807 non-null float64
V4        284807 non-null float64
V5        284807 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
```

# Exploratory Data Analysis

```
In [5]:  data.isnull().values.any()

Out[5]:  False

In [7]:  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[7]:  Text(0, 0.5, 'Frequency')
```



Transaction Class Distribution

```
In [9]:  ## Get the Fraud and the normal dataset

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

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

In [10]: 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
```
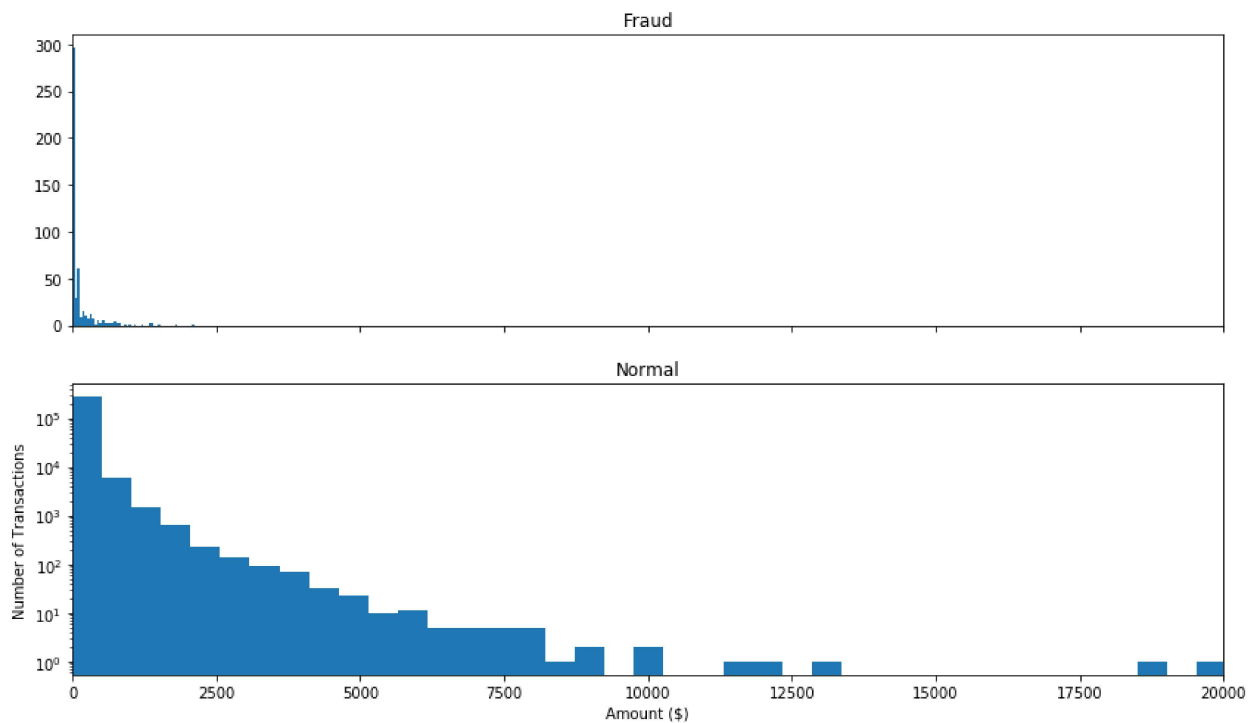
```
In [13]:  normal.Amount.describe()
```

```
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
```
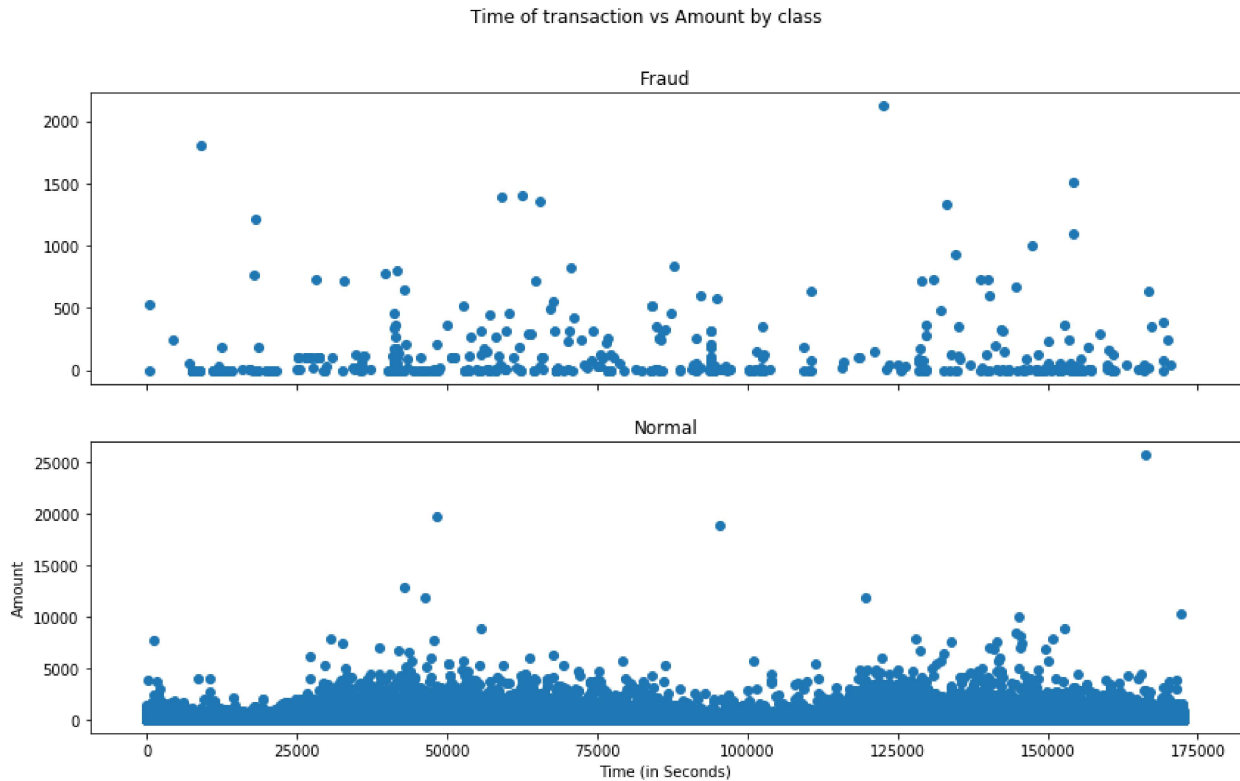
```
In [15]: f, (ax1, ax2) = plt.subplots(2, 1, sharex=True)
         f.suptitle('Amount per transaction by class')
         bins = 50
         ax1.hist(fraud.Amount, bins = bins)
         ax1.set_title('Fraud')
         ax2.hist(normal.Amount, bins = bins)
         ax2.set_title('Normal')
         plt.xlabel('Amount ($)')
         plt.ylabel('Number of Transactions')
         plt.xlim((0, 20000))
         plt.yscale('log')
         plt.show();
```



Amount per transaction by class

```
In [16]:  # We Will check Do fraudulent transactions occur more often during certain time
           frame ? Let us find out with a visual representation.

          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()
```



```
In [17]:  ## Take some sample of the data

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

          data1.shape

Out[17]: (28481, 31)

In [18]:  data.shape

Out[18]: (284807, 31)

In [19]:  #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))
```
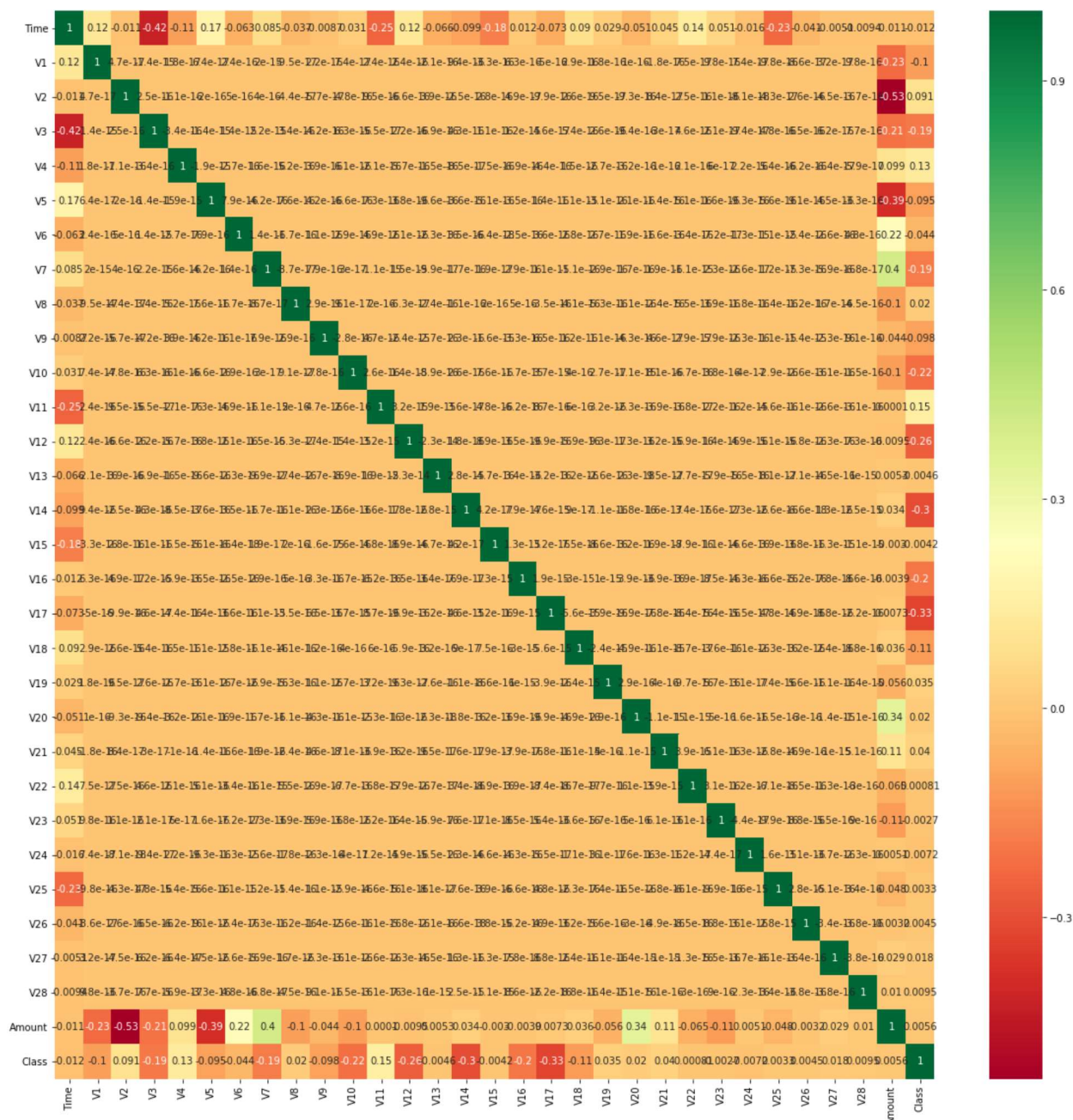
```
In [20]:  print(outlier_fraction)

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

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

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

```
In [22]:  ## Correlation
          import seaborn as sns
          #get correlations of each features in dataset
          corrmat = data1.corr()
          top_corr_features = corrmat.index
          plt.figure(figsize=(20,20))
          #plot heat map
          g=sns.heatmap(data[top_corr_features].corr(),annot=True,cmap="RdYlGn")
```

```
In [23]: #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,)
```

# Model Prediction

Now it is time to start building the model .The types of algorithms we are going to use to try to do anomaly detection on this dataset are as follows

## Isolation Forest Algorithm :

One of the newest techniques to detect anomalies is called Isolation Forests. The algorithm is based on the fact that anomalies are data points that are few and different. As a result of these properties, anomalies are susceptible to a mechanism called isolation.

This method is highly useful and is fundamentally different from all existing methods. It introduces the use of isolation as a more effective and efficient means to detect anomalies than the commonly used basic distance and density measures. Moreover, this method is an algorithm with a low linear time complexity and a small memory requirement. It builds a good performing model with a small number of trees using small sub-samples of fixed size, regardless of the size of a data set.

Typical machine learning methods tend to work better when the patterns they try to learn are balanced, meaning the same amount of good and bad behaviors are present in the dataset.

How Isolation Forests Work The Isolation Forest algorithm isolates observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature. The logic argument goes: isolating anomaly observations is easier because only a few conditions are needed to separate those cases from the normal observations. On the other hand, isolating normal observations require more conditions. Therefore, an anomaly score can be calculated as the number of conditions required to separate a given observation.

The way that the algorithm constructs the separation is by first creating isolation trees, or random decision trees. Then, the score is calculated as the path length to isolate the observation.

## Local Outlier Factor(LOF) Algorithm

The LOF algorithm is an unsupervised outlier detection method which computes the local density deviation of a given data point with respect to its neighbors. It considers as outlier samples that have a substantially lower density than their neighbors.

The number of neighbors considered, (parameter n_neighbors) is typically chosen 1) greater than the minimum number of objects a cluster has to contain, so that other objects can be local outliers relative to this cluster, and 2) smaller than the maximum number of close by objects that can potentially be local outliers. In practice, such informations are generally not available, and taking n_neighbors=20 appears to work well in general.

```
In [24]:   ##Define the outlier detection methods

           classifiers = {
               "Isolation Forest":IsolationForest(n_estimators=100, max_samples=len(X),
                                                   contamination=outlier_fraction,random_sta
           te=state, verbose=0),
               "Local Outlier Factor":LocalOutlierFactor(n_neighbors=20, algorithm='auto',
                                                         leaf_size=30, metric='minkowski',
                                                         p=2, metric_params=None, contamina
           tion=outlier_fraction),
               "Support Vector Machine":OneClassSVM(kernel='rbf', degree=3, gamma=0.1,nu=0.
           05,
                                                    max_iter=-1, random_state=state)

           }
```

```
In [26]:   type(classifiers)
```

Out[26]:   dict

```python
In [27]: 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 tra
         nsactions
             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))
```

```
C:\Users\krish.naik\AppData\Local\Continuum\anaconda3\envs\myenv\lib\site-package
s\sklearn\ensemble\iforest.py:223: FutureWarning: behaviour="old" is deprecated a
nd will be removed in version 0.22. Please use behaviour="new", which makes the d
ecision_function change to match other anomaly detection algorithm API.
  FutureWarning)
C:\Users\krish.naik\AppData\Local\Continuum\anaconda3\envs\myenv\lib\site-package
s\sklearn\ensemble\iforest.py:417: DeprecationWarning: threshold_ attribute is de
precated in 0.20 and will be removed in 0.22.
  " be removed in 0.22.", DeprecationWarning)
Isolation Forest: 73
Accuracy Score :
0.9974368877497279
Classification Report :
              precision    recall  f1-score   support

           0       1.00      1.00      1.00     28432
           1       0.26      0.27      0.26        49

   micro avg       1.00      1.00      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

   micro avg       1.00      1.00      1.00     28481
   macro avg       0.51      0.51      0.51     28481
weighted avg       1.00      1.00      1.00     28481


C:\Users\krish.naik\AppData\Local\Continuum\anaconda3\envs\myenv\lib\site-package
s\sklearn\svm\classes.py:1175: DeprecationWarning: The random_state parameter is
deprecated and will be removed in version 0.22.
  " be removed in version 0.22.", DeprecationWarning)
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

   micro avg       0.70      0.70      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 [ ]: