# CAPSTONE PROJECT INTERIM REPORT

| Batch details | PGP – DSE ONLINE FEB 21 (CHENNAI) |
|---|---|
| Team members | DHANESH K |
| | MANOJ KUMARAN R |
| | ANIRUDH RAMACHANDRA |
| | SUDVEEP KUMAR REDDY |
| | VENKATRAMANA S |
| | NAGARAJU ARVETI |
| Domain of Project | FRAUD DETECTION |
| Proposed project title | CREDIT CARD FRAUD DETECTION |
| Group Number | 4 |
| Team Leader | DHANESH K |
| Mentor Name | MS ANJANA AGRAWAL |

# TABLE OF CONTENTS

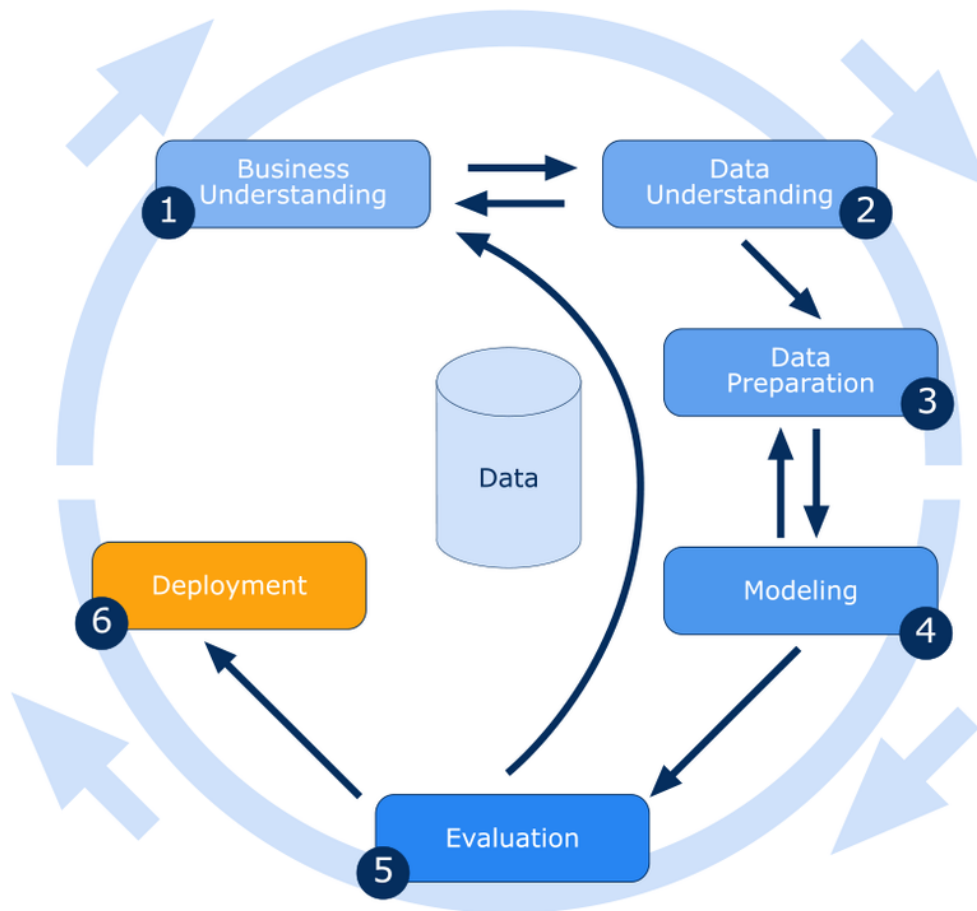# OVERVIEW

In today's world, high dependency on internet technology has enjoyed increased credit card transactions but, credit card fraud has also accelerated as online and offline transaction. As credit card transactions become a widespread mode of payment, focus has been given to recent computational methodologies to handle the credit card fraud problem.

Data mining technique is one notable and popular methods used in solving credit fraud detection problem. It is impossible to be sheer certain about the true intention and rightfulness behind an application or transaction. Credit card transaction datasets are rarely available, highly imbalanced and skewed. Optimal feature (variables) selection for the models, suitable metric is most important part of data mining to evaluate performance of techniques on skewed credit card fraud data.

From the experiments, the result that has been concluded is that Logistic regression has an accuracy of 0.99, KNN with 0.99. The boosting methods like Gradient boosting, AdaBoost and XGBoost shows an accuracy of 1.00.

# METHODOLOGY



## BUSINESS UNDERSTANDING

As online transactions increased over a period of time online fraud also subsequently increased. Although the proportion of fraud to nonfraud transactions is very low the chances of online fraud to grow is increasing gradually.

The goal is to identify and reduce the online fraud transactions without disturbing the flow of legal transactions.

By reducing the number of frauds, this will increase the trust of the customers in approaching the bank for the future.

# DATA UNDERSTANDING





- The dataset contains 284807 rows and 31 columns.

- There is no categorical columns present in the dataset as it contains only numerical columns.
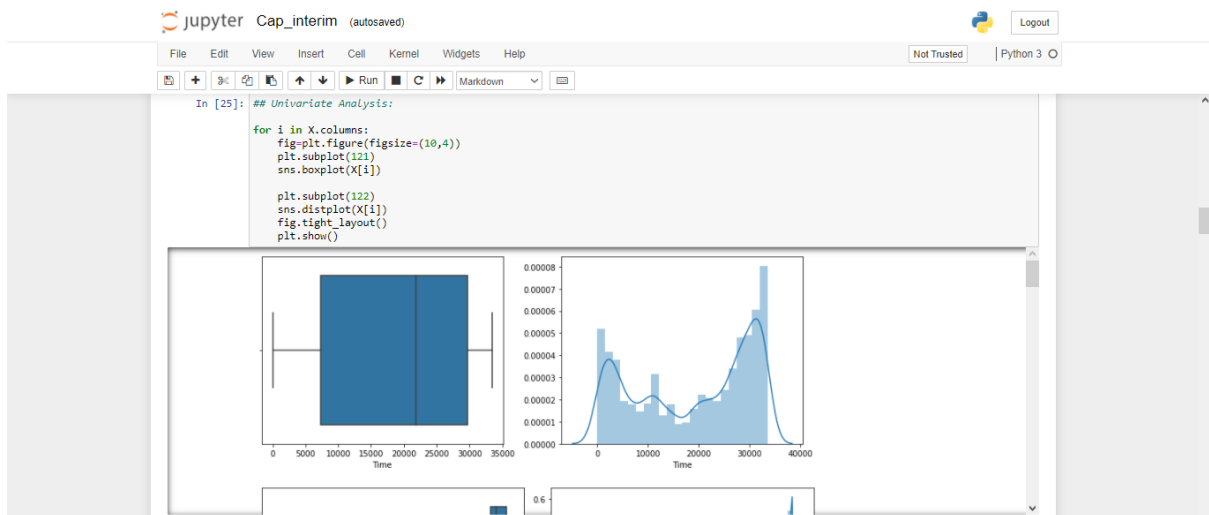
## DATA PREPARATION

| No of rows / columns | 284807 rows / 31 columns | |
|---|---|---|
| Missing values | No missing values | |
| Categorical features | None | |
| Skewness value (numerical data) | V1 = -3.280667 | V10 = 1.187141 |
| | V2 = -4.624866 | V12 = -2.278401 |
| | V3 = -2.240155 | V14 = -1.995176 |
| | V5 = -2.425901 | V16 = -1.100966 |
| | V6 = 1.826581 | V17 = -3.844914 |
| | V7 = 2.553907 | V28 = 11.192091 |
| | V8 = -8.521944 | Amount = 16.977724 |

## MISSING VALUES



```
In [5]:   df.isnull().sum()

Out[5]: Time      0
        V1        0
        V2        0
        V3        0
        V4        0
        V5        0
        V6        0
        V7        0
        V8        0
        V9        0
        V10       0
        V11       0
        V12       0
        V13       0
        V14       0
        V15       0
        V16       0
        V17       0
        V18       0
        V19       0
        V20       0
        V21       0
        V22       0
        V23       0
        V24       0
        V25       0
        V26       0
        V27       0
        V28       0
        Amount    0
        Class     0
```

## SKEWNESS
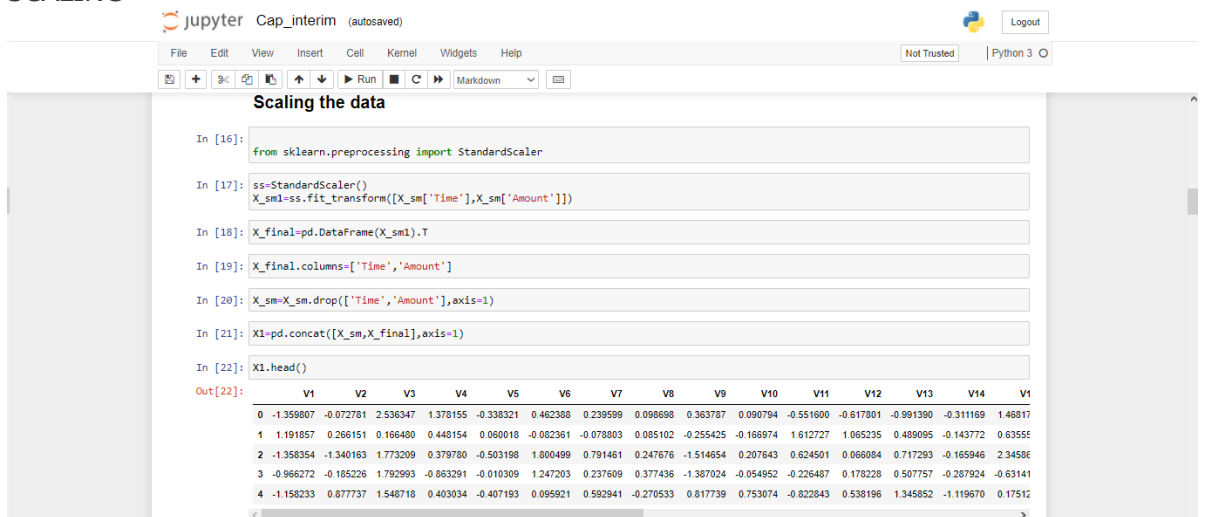
- Column time does not have outliers in it.
- Other than Time all other columns have extreme values or outliers which can't be denoted due to PCA.
- So, a distplot is also used to see the distribution and check skewness.
- Columns V1, V2, V3, V5, V6, V7, V8, V12, V14, V16, V17, V20, V21, V23, V27, V28, Amount all seem to be skewed.
- So, the skewness has to be removed from these columns.

```
In [25]: ## Univariate Analysis:

         for i in X.columns:
             fig=plt.figure(figsize=(10,4))
             plt.subplot(121)
             sns.boxplot(X[i])

             plt.subplot(122)
             sns.distplot(X[i])
             fig.tight_layout()
             plt.show()
```

OUTLIER TREATMENT

- Outlier treatment can't be done as this is a dataset which has undergone Principal Component Analysis.

- Capping method using IQR is not preferred because the information will be lost.

- Transformation techniques like log, sqrt, exponential, box-cox won't work as there are values that are very close to zero and negative values as well.

- So, it is better not to do outlier treatment here.

SCALING



**Scaling the data**

```
In [16]: from sklearn.preprocessing import StandardScaler

In [17]: ss=StandardScaler()
         X_sm1=ss.fit_transform([X_sm['Time'],X_sm['Amount']])

In [18]: X_final=pd.DataFrame(X_sm1).T

In [19]: X_final.columns=['Time','Amount']

In [20]: X_sm=X_sm.drop(['Time','Amount'],axis=1)

In [21]: X1=pd.concat([X_sm,X_final],axis=1)

In [22]: X1.head()
```

Out[22]:

| | V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 | V10 | V11 | V12 | V13 | V14 | V1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | -1.359807 | -0.072781 | 2.536347 | 1.378155 | -0.338321 | 0.462388 | 0.239599 | 0.098698 | 0.363787 | 0.090794 | -0.551600 | -0.617801 | -0.991390 | -0.311169 | 1.46817 |
| 1 | 1.191857 | 0.266151 | 0.166480 | 0.448154 | 0.060018 | -0.082361 | -0.078803 | 0.085102 | -0.255425 | -0.166974 | 1.612727 | 1.065235 | 0.489095 | -0.143772 | 0.63555 |
| 2 | -1.358354 | -1.340163 | 1.773209 | 0.379780 | -0.503198 | 1.800499 | 0.791461 | 0.247676 | -1.514654 | 0.207643 | 0.624501 | 0.066084 | 0.717293 | -0.165946 | 2.34586 |
| 3 | -0.966272 | -0.185226 | 1.792993 | -0.863291 | -0.010309 | 1.247203 | 0.237609 | 0.377436 | -1.387024 | -0.054952 | -0.226487 | 0.178228 | 0.507757 | -0.287924 | -0.63141 |
| 4 | -1.158233 | 0.877737 | 1.548718 | 0.403034 | -0.407193 | 0.095921 | 0.592941 | -0.270533 | 0.817739 | 0.753074 | -0.822843 | 0.538196 | 1.345852 | -1.119670 | 0.17512 |

# MODELLING

## SIGNIFICANT VARIABLES



**Logit Regression Results**

| | coef | std err | z | P>\|z\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| const | -8.3917 | 0.249 | -33.652 | 0.000 | -8.880 | -7.903 |
| Time | -3.742e-06 | 2.26e-06 | -1.659 | 0.097 | -8.16e-06 | 6.79e-07 |
| V1 | 0.0960 | 0.042 | 2.264 | 0.024 | 0.013 | 0.179 |
| V2 | 0.0094 | 0.058 | 0.161 | 0.872 | -0.104 | 0.123 |
| V3 | -0.0079 | 0.053 | -0.149 | 0.881 | -0.112 | 0.096 |
| V4 | 0.6986 | 0.074 | 9.454 | 0.000 | 0.554 | 0.843 |
| V5 | 0.1295 | 0.067 | 1.944 | 0.052 | -0.001 | 0.260 |
| V6 | -0.1198 | 0.074 | -1.626 | 0.104 | -0.264 | 0.025 |
| V7 | -0.0969 | 0.067 | -1.453 | 0.146 | -0.228 | 0.034 |
| V8 | -0.1739 | 0.030 | -5.711 | 0.000 | -0.234 | -0.114 |

Dep. Variable: Class, No. Observations: 284807, Model: Logit, Df Residuals: 284776, Method: MLE, Df Model: 30, Date: Thu, 09 Sep 2021, Pseudo R-squ.: 0.6922, Time: 22:57:58, Log-Likelihood: -1114.8, converged: True, LL-Null: -3621.2, Covariance Type: nonrobust, LLR p-value: 0.000



| | coef | std err | z | P>\|z\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| V8 | -0.1739 | 0.030 | -5.711 | 0.000 | -0.234 | -0.114 |
| V9 | -0.2843 | 0.111 | -2.561 | 0.010 | -0.502 | -0.067 |
| V10 | -0.8176 | 0.097 | -8.432 | 0.000 | -1.008 | -0.628 |
| V11 | -0.0621 | 0.081 | -0.762 | 0.446 | -0.222 | 0.098 |
| V12 | 0.0909 | 0.087 | 1.045 | 0.296 | -0.080 | 0.261 |
| V13 | -0.3312 | 0.082 | -4.058 | 0.000 | -0.491 | -0.171 |
| V14 | -0.5571 | 0.062 | -8.949 | 0.000 | -0.679 | -0.435 |
| V15 | -0.1141 | 0.086 | -1.330 | 0.183 | -0.282 | 0.054 |
| V16 | -0.1908 | 0.125 | -1.526 | 0.127 | -0.436 | 0.054 |
| V17 | -0.0216 | 0.070 | -0.309 | 0.757 | -0.159 | 0.116 |
| V18 | -0.0131 | 0.129 | -0.102 | 0.919 | -0.266 | 0.240 |
| V19 | 0.0963 | 0.097 | 0.993 | 0.321 | -0.094 | 0.286 |
| V20 | -0.4582 | 0.082 | -5.607 | 0.000 | -0.618 | -0.298 |
| V21 | 0.3898 | 0.060 | 6.494 | 0.000 | 0.272 | 0.507 |
| V22 | 0.6297 | 0.134 | 4.707 | 0.000 | 0.367 | 0.892 |
| V23 | -0.0951 | 0.058 | -1.629 | 0.103 | -0.209 | 0.019 |
| V24 | 0.1289 | 0.147 | 0.874 | 0.382 | -0.160 | 0.418 |
| V25 | -0.0761 | 0.131 | -0.582 | 0.560 | -0.332 | 0.180 |
| V26 | 0.0195 | 0.190 | 0.103 | 0.918 | -0.352 | 0.392 |
| V27 | -0.8188 | 0.122 | -6.686 | 0.000 | -1.059 | -0.579 |
| V28 | -0.2937 | 0.088 | -3.332 | 0.001 | -0.467 | -0.121 |
| Amount | 0.0009 | 0.000 | 2.449 | 0.014 | 0.000 | 0.002 |

- For a column to be considered significant statistically, the pvalue should be less than 0.05.
- Else it is insignificant.
- So, the values considered significant are v1, v4, v5, v8, v9, v10, v20, v21, v22, v27, v28, Amount.

## BASE MODEL



The classification report shows the accuracy as 0.99 for the test data and also for the train data the accuracy is 0.99.

## KNN MODEL



For the KNN model, the accuracy for the test data is showing as 0.99 and for the train data it is showing as 0.99.

## NAÏVE BAYES MODEL



```
In [44]: ## Naïve Bayes:
         from sklearn.naive_bayes import GaussianNB

In [45]: gnb=GaussianNB()
         gnb.fit(X_train,y_train)

Out[45]: GaussianNB()

In [46]: y_pred_test_gaussian=gnb.predict(X_test)
         y_pred_train_gaussian=gnb.predict(X_train)

In [47]: print(classification_report(y_test,y_pred_test_gaussian))  ## classification report for test data

                       precision    recall  f1-score   support

                    0       0.99      0.98      0.99      7447
                    1       0.83      0.91      0.86       774

             accuracy                           0.97      8221
            macro avg       0.91      0.94      0.93      8221
         weighted avg       0.97      0.97      0.97      8221


In [48]: print(classification_report(y_train,y_pred_train_gaussian))  ## classification report for train data

                       precision    recall  f1-score   support

                    0       0.99      0.98      0.99     17465
                    1       0.83      0.94      0.88      1717

             accuracy                           0.98     19182
            macro avg       0.91      0.96      0.93     19182
         weighted avg       0.98      0.98      0.98     19182
```

For the Naïve Bayes model, the accuracy for the test data is showing as 0.97 and for the train data it is showing as 0.98.

## AdaBoost MODEL



```
In [64]: adaboost = AdaBoostClassifier()

In [65]: adaboost.fit(X_train,y_train)

Out[65]: AdaBoostClassifier()

In [66]: y_pred_train = adaboost.predict(X_train)
         print(classification_report(y_train,y_pred_train))  ## classification report for train data

                       precision    recall  f1-score   support

                    0       1.00      1.00      1.00     17465
                    1       1.00      1.00      1.00      1717

             accuracy                           1.00     19182
            macro avg       1.00      1.00      1.00     19182
         weighted avg       1.00      1.00      1.00     19182


In [67]: y_pred_test = adaboost.predict(X_test)
         print(classification_report(y_test,y_pred_test))   ## classification report for test data

                       precision    recall  f1-score   support

                    0       1.00      1.00      1.00      7447
                    1       0.99      0.99      0.99       774

             accuracy                           1.00      8221
            macro avg       1.00      1.00      1.00      8221
         weighted avg       1.00      1.00      1.00      8221
```
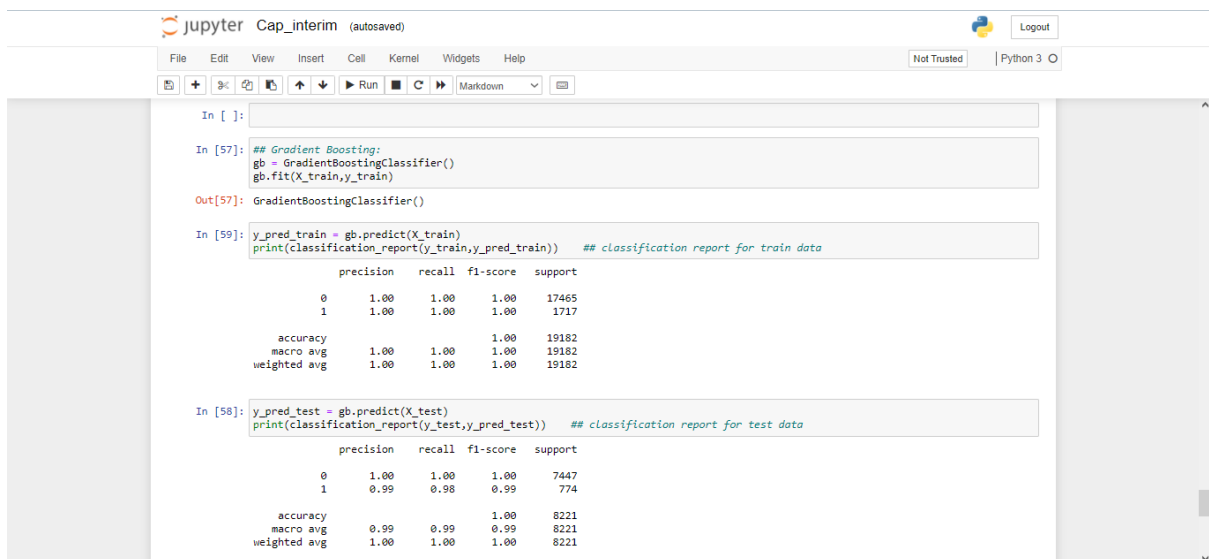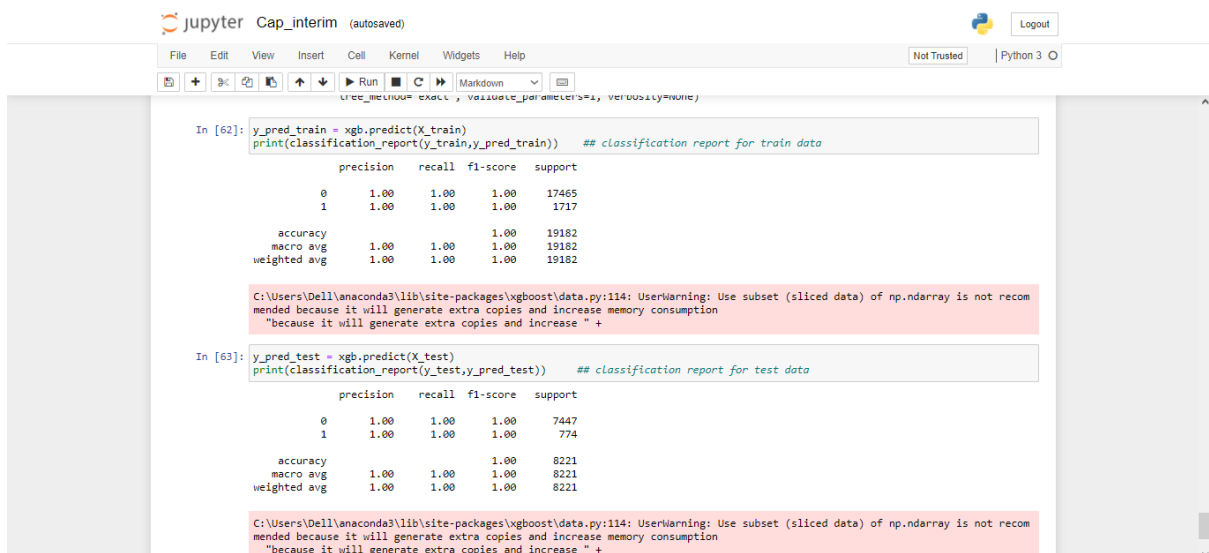
For the Ada Boost model, the accuracy for the test data is showing as 1.00 and for the train data it is showing as 1.00.

GRADIENT BOOSTING MODEL



For the Gradient Boosting model, the accuracy for the test data is showing as 1.00 and for the train data it is showing as 1.00.

XG BOOST MODEL



For the XG Boost model, the accuracy for the test data is showing as 1.00 and for the train data it is showing as 1.00.

# REFERENCE

The references used in the project are given below:

1.) Calibrating Probability with Under sampling for Unbalanced Classification by Andrea Dal Pozzolo, Olivier Caelen, Reid A. Johnson and Gianluca Bontempi.

2.) Learned lessons in credit card fraud detection from a practitioner perspective by Dal Pozzolo, Andrea; Caelen, Olivier; Le Borgne, Yann-Ael; Waterschoot, Serge; Bontempi, Gianluca.

3.) Combining Unsupervised and Supervised Learning in Credit Card Fraud Detection by Fabrizio Carcillo, Yann-Aël Le Borgne, Olivier Caelen, Frederic Oblé, Gianluca Bontempi.