

# **Project Report**

## **Proposal with Problem Statement :**

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. Identify fraudulent credit card transactions.

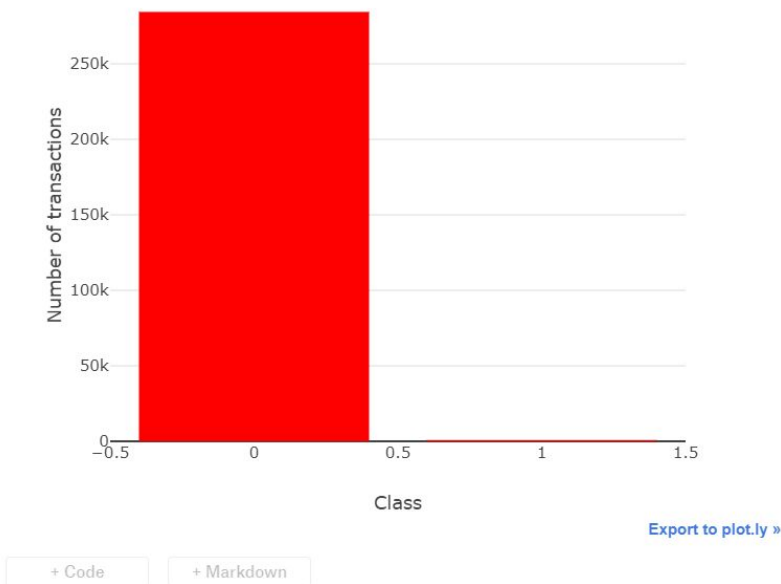
## **Data collection and Wrangling summary:**

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 numeric 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-sensitive learning. Feature 'Class' is the response variable and it takes value 1 in case of fraud and 0 otherwise.

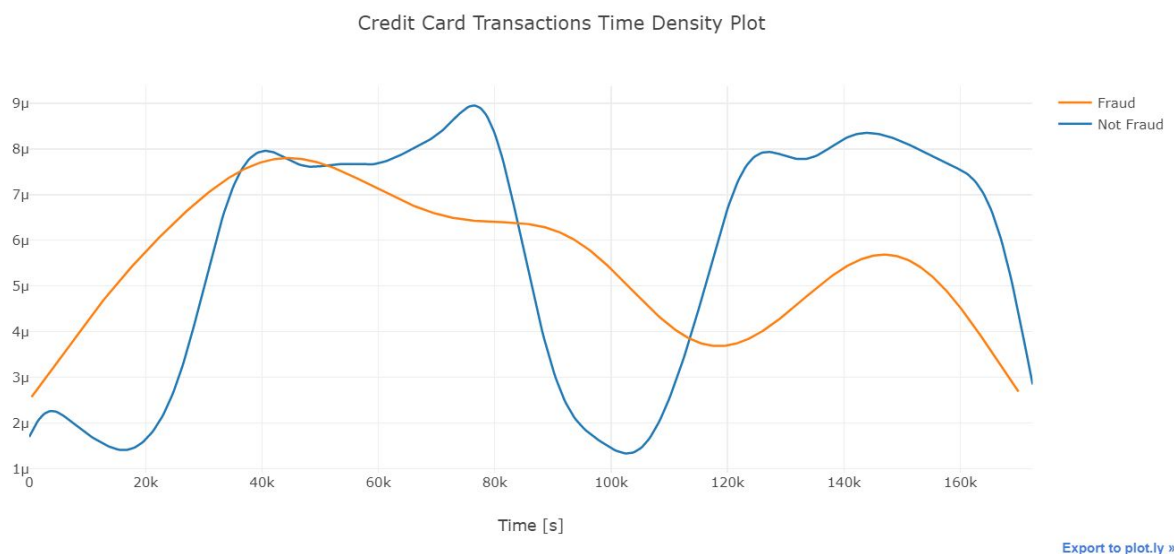
## **EDA Summary (Visualisation and Inferential stats):**

Credit Card Fraud Class - data imbalance (Not fraud = 0, Fraud = 1)

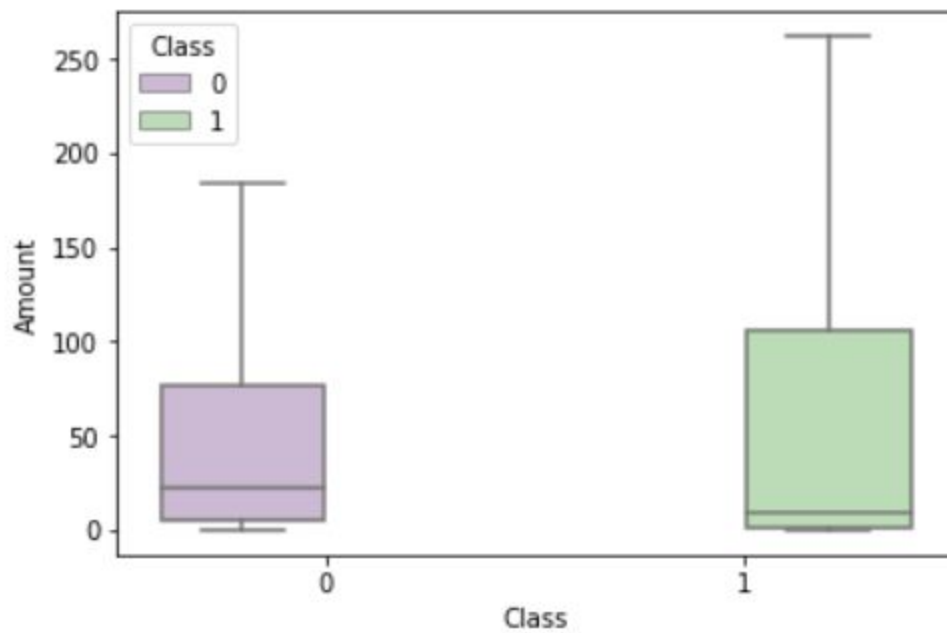


Only **492** (or **0.172%**) of transaction are fraudulent. That means the data is highly unbalanced with respect with target variable **Class**.

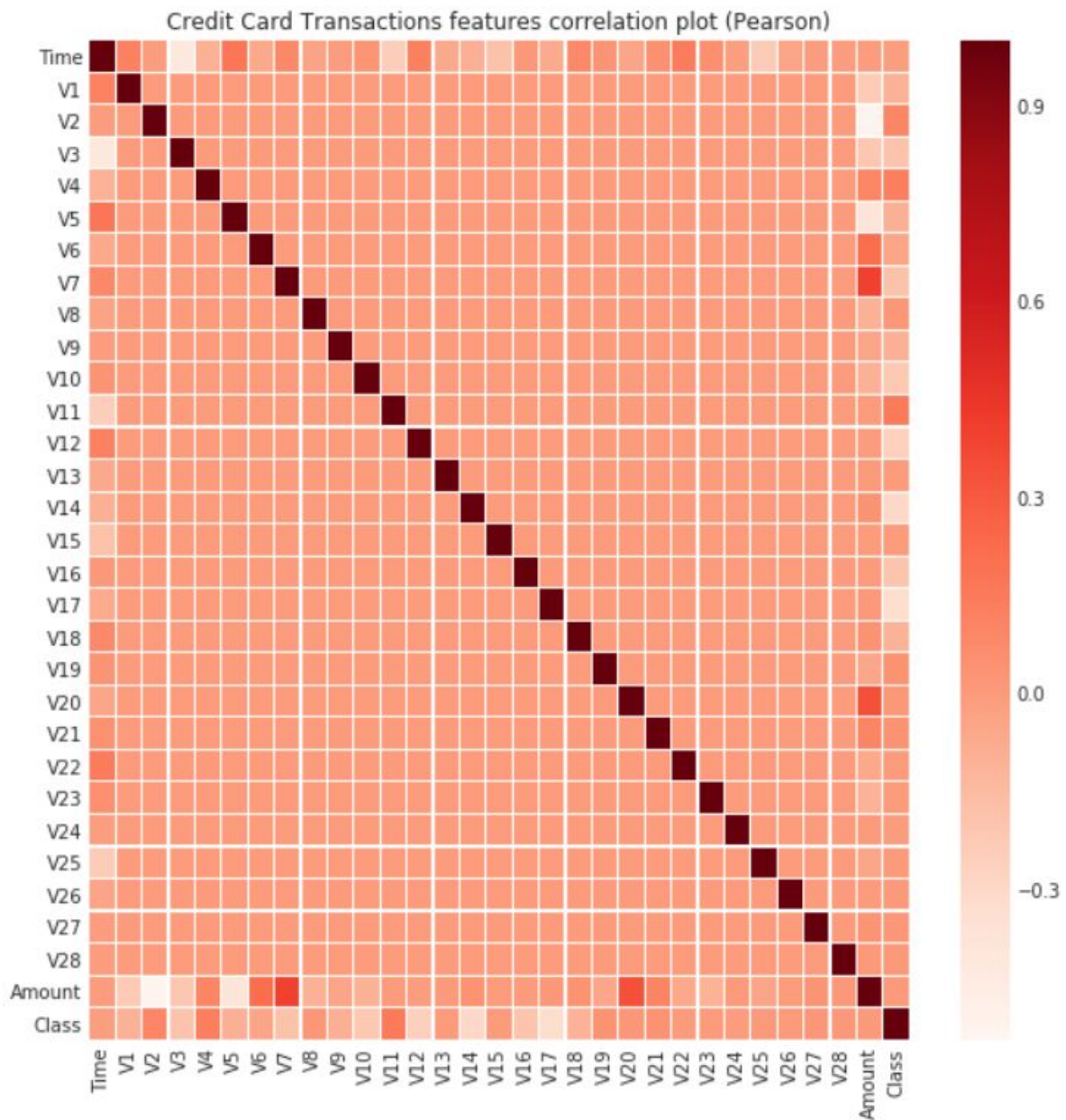
Fraudulent transactions have a distribution more even than valid transactions - are equally distributed in time, including the low real transaction times, during night in Europe timezone.



The real transaction have a larger mean value, larger Q1, smaller Q3 and Q4 and larger outliers; fraudulent transactions have a smaller Q1 and mean, larger Q4 and smaller outliers.



There is no notable correlation between features V1-V28. There are certain correlations between some of these features and Time (inverse correlation with V3) and Amount (direct correlation with V7 and V20, inverse correlation with V1 and V5).

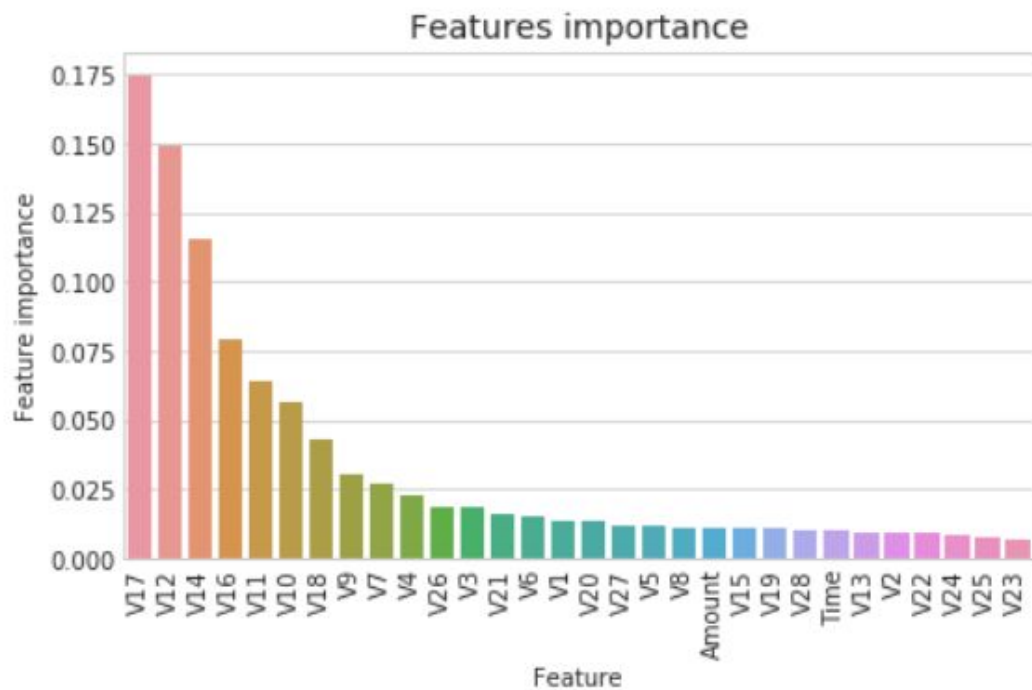


Results and in depth analysis using ML:

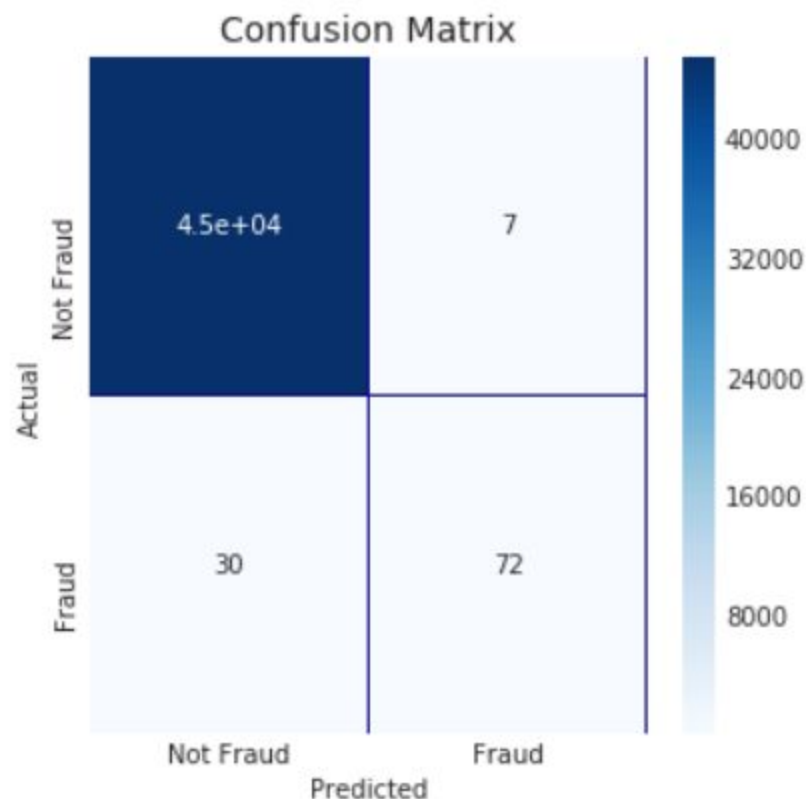
I started with Random Forest Classifier using GINI index.

```
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                        max_depth=None, max_features='auto', max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=4,
                        oob_score=False, random_state=2018, verbose=False,
                        warm_start=False)
```

Got the Feature importance as below :



Confusion Matrix :



### Type I error and Type II error

We need to clarify that confusion matrix are not a very good tool to represent the results in the case of largely unbalanced data, because we will actually need a different metrics that accounts at the same time for the **selectivity** and **specificity** of the method we are using, so that we minimize at the same time both Type I errors and Type II errors.

Null Hypothesis (H0) - The transaction is not a fraud.

Alternative Hypothesis (H1) - The transaction is a fraud.

Type I error - You reject the null hypothesis when the null hypothesis is actually true.

Type II error - You fail to reject the null hypothesis when the alternative hypothesis is true.

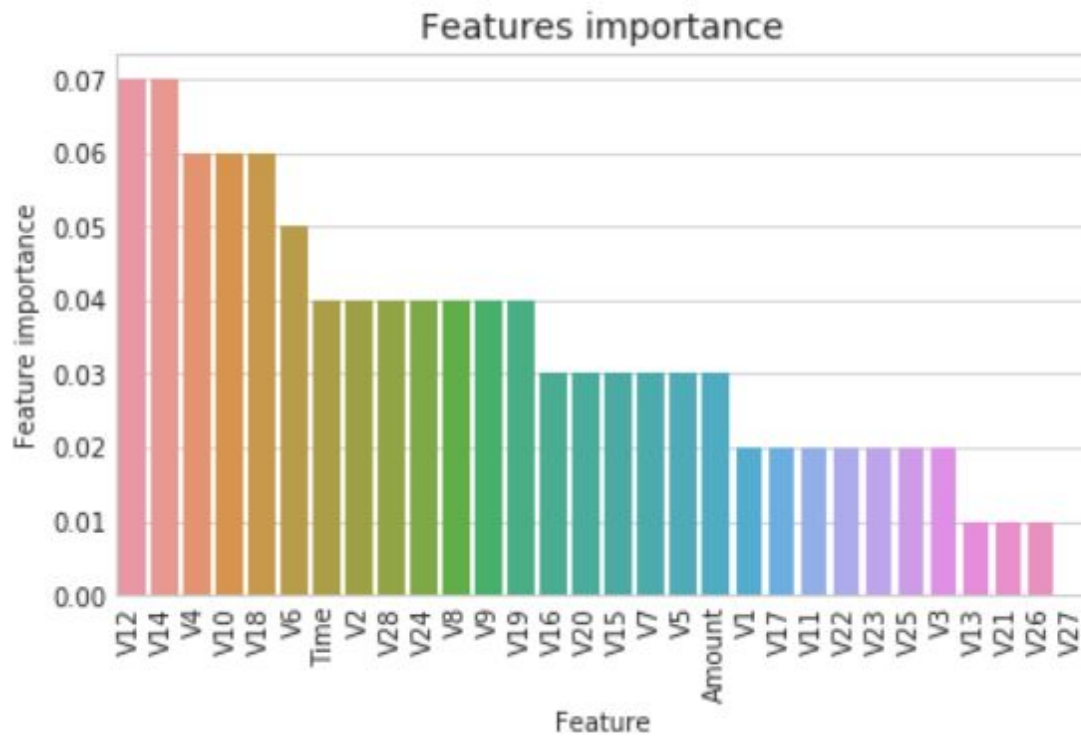
Cost of Type I error - You erroneously presume that the transaction is a fraud, and a true transaction is rejected.

Cost of Type II error - You erroneously presume that the transaction is not a fraud and a fraudulent transaction is accepted.

**The ROC-AUC score obtained with RandomForestClassifier is 0.85**

The I used ADABOOST Classifier

```
AdaBoostClassifier(algorithm='SAMME.R', base_estimator=None,  
learning_rate=0.8, n_estimators=100, random_state=2018)
```



The ROC-AUC score obtained with AdaBoostClassifier is 0.83.

Used other techniques as well . Complete summary is as below :

We investigated the data, checking for data unbalancing, visualizing the features and understanding the relationship between different features. We then investigated two predictive models. The data was split in 3 parts, a train set, a validation set and a test set. For the first three models, we only used the train and test set.

We started with **RandomForrestClassifier**, for which we obtained an AUC score of **0.85** when predicting the target for the test set.

We followed with an **AdaBoostClassifier** model, with lower AUC score (**0.83**) for prediction of the test set target values.

We then followed with an **CatBoostClassifier**, with the AUC score after training 500 iterations **0.86**.

We then experimented with a **XGBoost** model. In this case, we used the validation set for validation of the training model. The best validation score obtained was **0.984**. Then we used the model with the best training step, to predict target value from the test data; the AUC score obtained was **0.974**.

We then presented the data to a **LightGBM** model. We used both train-validation split and cross-validation to evaluate the model effectiveness to predict 'Class' value, i.e. detecting if a transaction was fraudulent. With the first method we obtained values of AUC for the validation set around **0.974**. For the test set, the score obtained was **0.946**.

With the cross-validation, we obtained an AUC score for the test prediction of **0.93**.