

**MEF UNIVERSITY**

**CREDIT CARD FRAUD DETECTION USING  
MACHINE LEARNING**

**Capstone Project**

**Tibet Erdoğan**

**İSTANBUL, 2021**



**MEF UNIVERSITY**

**CREDIT CARD FRAUD DETECTION USING  
MACHINE LEARNING**

**Capstone Project**

**Tibet Erdođan**

**Advisor: Asst. Prof. Dr. Duygu TAŞ KÜTEN**

**İSTANBUL, 2021**

## MEF UNIVERSITY

Name of the project: Credit Card Fraud Detection  
Name/Last Name of the Student: Tibet Erdoğan  
Date of Thesis Defense: 25/01/2021

I hereby state that the graduation project prepared by Tibet Erdoğan has been completed under my supervision. I accept this work as a “Graduation Project”.

25/01/2021  
Asst. Prof. Dr. Duygu TAŞ KÜTEN

I hereby state that I have examined this graduation project by Tibet Erdoğan which is accepted by his supervisor. This work is acceptable as a graduation project and the student is eligible to take the graduation project examination.

25/01/2021  
Prof. Dr. Özgür ÖZLÜK  
Director  
of  
Big Data Analytics Program

We hereby state that we have held the graduation examination of \_\_\_\_\_ and agree that the student has satisfied all requirements.

### THE EXAMINATION COMMITTEE

Committee Member

Signature

1. Asst. Prof. Dr. Duygu TAŞ KÜTEN

.....

2. Prof.Dr. Özgür ÖZLÜK

.....

## Academic Honesty Pledge

I promise not to collaborate with anyone, not to seek or accept any outside help, and not to give any help to others.

I understand that all resources in print or on the web must be explicitly cited.

In keeping with MEF University's ideals, I pledge that this work is my own and that I have neither given nor received inappropriate assistance in preparing it.

---

Name	Date	Signature
Tibet Erdoğan	25/01/2021	

# EXECUTIVE SUMMARY

## CREDIT CARD FRAUD DETECTION USING MACHINE LEARNING

Tibet Erdoğan

Advisor: Asst. Prof. Duygu TAŞ KÜTEN

JANUARY, 2021, 18 pages

This project aims to find the most efficient machine learning models to detect fraudulent transactions on credit cards.

The dataset used for this project consists of credit card transactions made by European cardholders in September 2013. This dataset presents transactions that have occurred in two days, where there are 492 frauds out of 284,807 transactions.

Machine learning methods, such as decision trees, logistic regression and random forest classifier are used to predict the fraudulent transactions. Performance of these machine learning models are compared to achieve the highest accuracy.

According to the results, it is found that the random forest classifier is the most effective model, and the SMOTE technique used to overcome the data imbalance performs better than the under-sampling technique. It is also observed that the models employed with the under-sampled data misclassify large number of non-fraud transactions as fraud. Lastly, by means of the random forest with the over-sampling technique (SMOTE), it is observed that the feature “V13” has the most important role in detecting fraud.

**Key Words:** Fraud Detection, Credit Card Fraud, Machine Learning, Logistic Regression, Decision Trees, Random Forest Method

# ÖZET

## MAKİNE ÖĞRENMESİ İLE KREDİ KARTI DOLANDIRICILIĞI TESPİTİ

Tibet ERDOĞAN

Proje Danışmanı: Dr. Öğr. Üyesi Duygu TAŞ KÜTEN

OCAK, 2021, 18 sayfa

Bu proje, uygun makine öğrenmesi modeli geliştirilerek kredi kartı dolandırıcılığını tespit etmeyi amaçlamaktadır.

Bu projede kullanılan veri kümesi 2013 yılının Eylül ayında Avrupalı müşteriler tarafından gerçekleştirilen kredi kartı işlemlerini içermektedir. İki gün içerisinde gerçekleştiren 284.807 adet işlem yer almakta olup, bunlardan 487 adeti dolandırıcılık işlemidir.

Dolandırıcılık işlemlerini tahmin etmek için lojistik regresyon, rastgele orman, karar ağaçları gibi makine öğrenmesi metotları kullanılmıştır. En iyi tahmin skoruna ulaşmak için, kullanılan makine öğrenmesi modellerinin performansı karşılaştırılmıştır.

Elde edilen sonuçlara göre rastgele orman metodunun en etkili model olduğu ve veri dengesizliğinin üstesinden gelmek için kullanılan SMOTE tekniğinin rastgele azaltma tekniğine göre daha iyi performans gösterdiği bulunmuştur. Rastgele azaltma tekniği ile oluşturulan veri kümelerine uygulanan modellerin çok sayıda dolandırıcılık dışı işlemi dolandırıcılık olarak yanlış sınıflandırdığı da görülmektedir. Son olarak, rastgele orman metodunun SMOTE tekniğiyle uygulanması sonucunda, veri kümesindeki “V13” isimli özelliğin dolandırıcılık tespitinde en önemli rolü oynadığı görülmüştür.

**Anahtar Kelimeler:** Dolandırıcılık Tespiti, Kredi Kartı Dolandırıcılığı, Makine Öğrenmesi, Lojistik Regresyon, Karar Ağacı, Rastgele Orman Metodu

## TABLE OF CONTENTS

Academic Honesty Pledge .....	v
EXECUTIVE SUMMARY .....	vi
ÖZET .....	vii
TABLE OF CONTENTS.....	viii
LIST OF FIGURES .....	ix
LIST OF TABLES.....	x
1. INTRODUCTION .....	1
1.1. Credit Card Fraud Detection Using Machine Learning: Literature Survey.....	1
1.2. Machine Learning Algorithms.....	2
2. ABOUT THE DATA.....	3
2.1. About the Data.....	3
2.2. Exploratory Data Analysis.....	3
3. PROJECT DEFINITION.....	6
4. METHODOLOGY.....	7
4.1. Feature Scaling .....	7
4.2 Under-sampling vs Over-sampling.....	7
4.3. Measures of Performance .....	8
5. RESULTS .....	9
5.1 Random Forest.....	10
5.1.1 Random Forest - Over-sampling.....	10
5.1.2 Random Forest – Under-sampling.....	11
5.2 Decision Tree.....	12
5.2.1 Decision Tree – Over-sampling.....	12
5.2.2 Decision Tree – Under-sampling.....	13
5.3 Logistic Regression.....	14
5.3.1 Logistic Regression - Over-sampling .....	14
5.3.2 Logistic Regression – Under-sampling.....	15
5.4 Feature Importance .....	16
6. CONCLUSION.....	17
REFERENCES .....	18



## LIST OF FIGURES

Figure 1 : Distribution of Amount Feature .....	4
Figure 2: Distribution of Time Feature .....	4
Figure 3 : Class Distributions according to amount feature .....	5
Figure 4: Confusion Matrix of Random Forest - Over-sampling .....	11
Figure 5: Confusion Matrix of Random Forest – Under-sampling .....	12
Figure 6: Confusion Matrix of Decision Tree - Over-sampling .....	13
Figure 7: Confusion Matrix of Decision Tree – Under-sampling .....	14
Figure 8: Confusion Matrix of Logistic Regression - Over-sampling.....	15
Figure 9: Confusion Matrix of Logistic Regression – Under-sampling .....	16
Figure 10: Feature Importance of Random Forest.....	16

## LIST OF TABLES

Table 1: Performance metrics of Models.....	10
Table 2: Performance Metrics of Random Forest - Over-sampling .....	11
Table 3: Performance Metrics of Random Forest – Under-sampling.....	12
Table 4: Performance Metrics of Decision Tree – Over-sampling.....	13
Table 5: Performance Metrics of Decision Tree – Under-sampling.....	13
Table 6: Performance Metrics of Logistic Regression – Over-sampling .....	14
Table 7: Performance Metrics of Logistic Regression – Under-sampling .....	15

# 1. INTRODUCTION

As the use of credit cards become widespread day by day, credit card fraud becomes one of the biggest problems with that banks and customers are facing. Credit card fraud occurs when fraudster accesses someone's credit or debit card. Due to its likely repetition, fraudulent transactions harm both the customer and the service provider, and thus it is important to take preventive measures.

As a result of widespread usage of credit cards, massive amount of data that contains all of the credit card transactions are stored by service providers. By means of the availability of the massive data, using machine learning models is the most effective way to detect fraudulent transactions.

This project aims to find a suitable machine learning model to detect fraudulent transactions on credit cards. Decision trees, logistic regression and random forest methods are used to predict fraudulent transactions.

The dataset used for this project consists of transactions made by credit cards in September 2013 by European cardholders. This dataset presents transactions that have occurred in two days, where there are 492 frauds out of 284,807 transactions.

## 1.1. Credit Card Fraud Detection Using Machine Learning: Literature Survey

Credit card fraud is an unauthorized use of a card, to fraudulently obtain money. Credit and debit card numbers can be stolen from unsecured websites or can be obtained in an identity theft scheme [5].

As paying with credit cards becomes a primary payment method, banks, merchants and customers have faced with many fraud cases. Some of the most common credit card fraud cases are listed below [6]:

- inception of mails of newly issued cards,
- copying card information through cloned webpages,

- phishing in which credit card number and password are hacked through emails.

According to the Interbank Card Center of Turkey (BKM), 73,856,831 credit cards and 181,116,590 debit cards are in the use as of September 2020. By means of such a massive amount of data to handle, artificial intelligence technologies have been improved to detect fraudulent transactions.

## **1.2. Machine Learning Algorithms**

Logistic regression is used to predict the target variable, and it is mostly useful for the cases predicting the presence or absence of a characteristic or outcome based on estimation variables. Logistic regression coefficients can be used to assess odds ratios for each of the independent variables in the model [7].

A decision tree is a classifier expressed as a recursive partition of the instance space and a simple representation for classifying examples [10]. This method is a supervised machine learning where the data is continuously splitted according to a certain parameter. Decision tree based models outperform the SVM models when their performances are compared over the test data sets. When the performances of algorithms are compared over the training data sets, SVM based models overfit the training data [8].

Random forest creates multitude of decision trees and integrate them together to select the most effective feature. More specifically, random forest contains different decision trees on different subsets of the dataset, and it takes the average of that different decision tree's accuracy scores to improve the performance of the model [15]. In this way, it reduces overfitting problem in decision trees. On the other hand, random forest creates a lot trees and combine all of their outputs. As a result, it needs more power and resources.

## **2. ABOUT THE DATA**

### **2.1. About the Data**

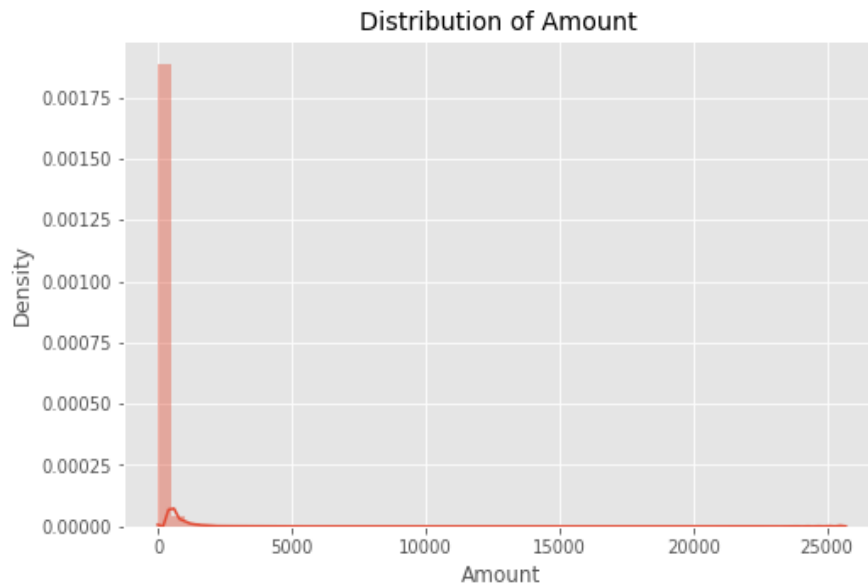
In this project, credit card transactions made by European cardholders in September 2013 are used. This dataset presents transactions that occurred in two days, where there are 492 frauds out of 284,807 transactions. Dataset has 284,807 rows and 31 columns.

The description of the data indicates that all the features has been through a Principal Component Analysis (PCA) transformation. PCA is a linear dimensionality reduction technique, and generally used on large datasets to reduce the data size and to better understand the structure of data. As a result, the dataset contains only numerical input variables [4]. Due to the protection of personal data, original features and background information cannot be provided, and thus the feature extraction and feature elimination are not implemented. All of the features are transformed with PCA except the time and the amount. To summarize, there are 28 encrypted and scaled features, and two unscaled features named as “Time” and “Amount”.

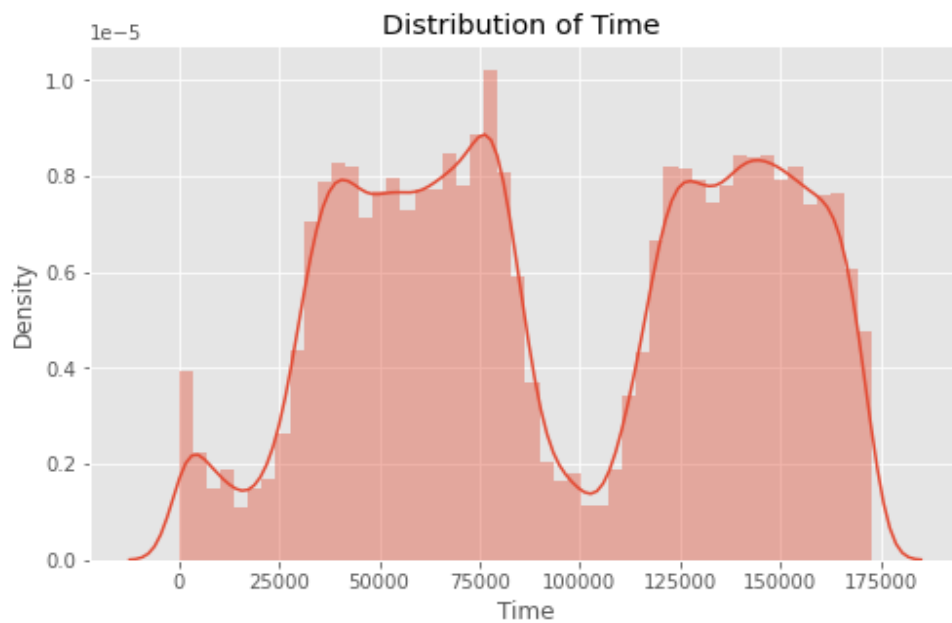
The feature entitled “Amount” is the credit card transaction amount, feature “Time” is the difference between the time of each transaction and that of the first transaction, and the feature “Class” takes the value 1 if the transaction is fraudulent and 0 otherwise.

### **2.2. Exploratory Data Analysis**

As it can be seen from Figures 1 and 2, features “Time” and “Amount” are highly skewed. These features need to be scaled according to the values in other columns. There are statistical techniques that can help us overcome skewness. In this project, the standard scaling technique is used to overcome skewness.



**Figure 1 :** Distribution of Amount Feature



**Figure 2:** Distribution of Time Feature

As it can be seen from Figure 3, this dataset contains 284,807 transactions and only 492 of that are suspicious. In other words, the number of regular transactions is larger than that of the fraudulent ones. As most of the transactions are non-fraud, the dataset is highly imbalanced. In other words, this data would lead to inaccurate overfitting and correlations. More specifically, the accuracy score may have a quite high accuracy rate since the algorithm can smoothly predict fraud ones by focusing on non-fraud transactions. There are many

techniques used to deal with such an imbalance. In this project, random under-sampling and synthetic minority over-sampling techniques (SMOTE) are used. SMOTE creates new minority class instances between existing ones. These new virtual instances created by  $k$ -nearest neighbors of examples in the minority class. Random under-sampling method balances the distribution by randomly eliminating majority class until a balanced distribution is reached [14]. Because of their success in dealing with class imbalance, SMOTE and random under-sampling feature scaling methods are used.



**Figure 3 :** Class Distributions according to amount feature

### **Under-sampling:**

A subset of 688 transactions is randomly taken from the original dataset. The sample subset contains 344 regular and 344 fraudulent transactions.

### **Over-sampling (Synthetic Minority Over-sampling Technique):**

The SMOTE creates new artificial rows in order to have an equal balance of the classes. Before we implement SMOTE, there were 284,807 transactions in the dataset and only 492 of that are fraudulent. SMOTE created 113,231 artificial rows, and thus the new dataset contains 199,019 regular and 199,019 fraudulent transactions.

### **3. PROJECT DEFINITION**

Making analyses on credit card fraud is quite important, as it affects financial institutions and their customers substantially. For that reason, several studies have been conducted to prevent credit card fraud.

Credit card frauds occur when fraudster access someone's credit or debit cards and repetitive in most of the cases. Such transactions financially damage the customers, merchants and service providers, and thus taking preventive measures are vital. Fraud methods may be changed over time as fraudsters constantly try to bypass the detection systems. For instance, establishment of a new payment system or new merchant would lead to the alteration [9].

In this capstone project, the main objective is to apply several effective machine learning algorithms for fraud detection. To achieve the objective of the project, logistic regression, random forest and decision tree algorithms are used, and prediction results are compared. A public dataset on credit card transactions has been used for this project. Due to confidentiality issues, features of dataset has been processed through Principal Component Analysis.



## 4. METHODOLOGY

### 4.1. Feature Scaling

In many machine learning algorithms, to acquire all features in a similar setting, scaling should be applied. In this way, any significant variable would not affect the performance of the model due to its large size [1]. There are two popular scaling algorithms that are the min-max normalization and the standardization.

The technique which provides linear transformation on original range of data is called the min-max normalization. The min-max normalization is a basic technique that can fit the data in a pre-defined boundary. In other words, the minimum value of the set of observed values is mapped to 0 and the maximum value of the set of observed values mapped to 1 [2].

The standardization is the method of approximating a variable to a normal variable. It is a process of shifting and rescaling the data to achieve mean of 0 and standard deviation of 1. Standardization processes make compatibility, similarity and keep measurement error at its lowest level [3].

### 4.2 Under-sampling vs Over-sampling

Data imbalance corresponds to the distributions of classes that are unequal in the dataset. In other words, data imbalance occurs when the number of instances of some specific classes significantly out numbers the instances of another classes. Under-sampling and over-sampling are the two techniques used to overcome the data imbalance. Under-sampling removes most of the samples, whereas over-sampling reproduces samples [14].

Under-sampling is a series of methods used to reconstruct dataset to improve the prediction. Since more than 99% of the transactions is labeled as normal in the considered dataset, random under-sampling would be quite effective to create a balanced training dataset. Synthetic Minority Over-sampling Technique (SMOTE) would provide a balanced dataset as well. SMOTE add new artificial rows into the dataset to have an equal balance of

the classes [11]. In this project, both the under-sampling and the over-sampling (SMOTE) techniques are used.

### 4.3. Measures of Performance

To evaluate classifier quality measures “Precision”, “Recall” and “F1 score” are used in this project. Before explaining those measures, classification of predictions is described with respect to the aim of this project as follows [13]:

- True Positive: The fraud cases that the model predicted as “fraud”.
- False Positive: The non-fraud cases that the model predicted as “fraud”.
- True Negative: The non-fraud cases that the model predicted as “non-fraud”.
- False Negative: The fraud cases that the model predicted as “non- fraud”.

Related explanations can be provided as follows [12].

**Precision:** Precision is a measure that shows the number of cases correctly predicted as positive out of the total number of cases predicted positive:

$$\text{Precision} = \text{TruePositives} / (\text{TruePositives} + \text{FalsePositives})$$

**Recall:** Recall is a measure that shows the number of cases correctly predicted as positive out of the total number of actual positive cases:

$$\text{Recall} = \text{TruePositives} / (\text{TruePositives} + \text{FalseNegatives})$$

**F1 Score:** F1 score is a performance score according to “Precision” and “Recall”. F1 score is used to analyze the balance between these two measures:

$$\text{F1 Score} = (2 * \text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

## 5. RESULTS

In this project, SMOTE and random under-sampling methods are used to overcome the issues since the dataset is imbalanced.

First, train - test split technique has been used. Train - test split technique is a process to divide the dataset into two subsets. The first subset entitled “training” is used to fit the model and the second subset entitled “test” is used to evaluate the fitted model. In this project, 70% of the dataset allocated for test, and 30% of the scaled data allocated for train of the model. The main idea of using 70% of the dataset for the test of the model is to compare different preprocessing techniques, SMOTE and under-sampling on the imbalanced data.

Many studies about credit fraud detection emphasize that using accuracy as a performance metric would lead to inaccurate results. This usually happens on unevenly distributed dataset. F1 score would give more accurate results to evaluate the model results as it takes both false positive predictions and false negative predictions. False positive and false negative predictions would be quite important to detect the credit card frauds.

Performance metric scores of all models on the considered credit card dataset with a number of feature scaling methods are shown in Table 2.

From Table 2 it can be seen that by means of the over-sampling technique, all of the algorithms used for this project perform relatively well except the logistic regression. Random forest classifier has the best performance with  $F1 = 0.88$ ,  $Recall = 0.82$  and  $Precision = 0.93$ . On the other hand, logistic regression and decision tree models, by means of the random under-sampling technique, perform quite well on detecting actual fraud cases with  $Recall = 0.92$  but produced many false positive cases which may create issues to financial institutions. Gridsearch is used to find the optimal parameters for the random forest classifier. By means of the Matthews Correlation Coefficient score, the most efficient parameters are used with the random forest classifier. It is found that the F1 score results

obtained by using the optimal parameters are not better than the results obtained by implementing the default ones.

**Table 1:** Performance metrics of Models

Fraud Detection			
Metric	Model	Score	Sampling Method
F1 Score	Random Forest	0.88	SMOTE
	Random Forest	0.14	Under-sampling
	Decision Tree	0.55	SMOTE
	Decision Tree	0.03	Under-sampling
	Logistic Regression	0.11	SMOTE
	Logistic Regression	0.09	Under-sampling
Recall Score	Random Forest	0.82	SMOTE
	Random Forest	0.84	Under-sampling
	Decision Tree	0.76	SMOTE
	Decision Tree	0.92	Under-sampling
	Logistic Regression	0.92	SMOTE
	Logistic Regression	0.88	Under-sampling
Precision Score	Random Forest	0.93	SMOTE
	Random Forest	0.08	Under-sampling
	Decision Tree	0.43	SMOTE
	Decision Tree	0.01	Under-sampling
	Logistic Regression	0.06	SMOTE
	Logistic Regression	0.05	Under-sampling

## 5.1 Random Forest

### 5.1.1 Random Forest - Over-sampling

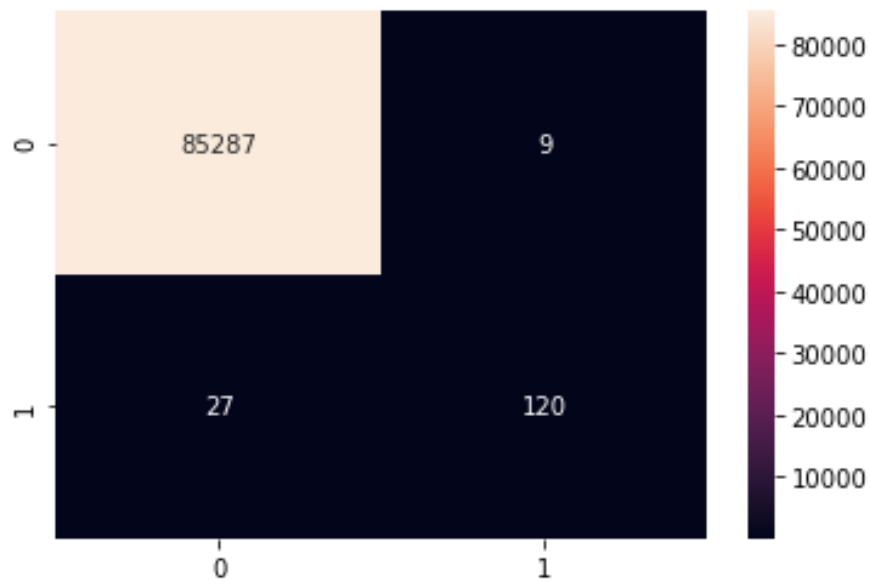
The random SMOTE method is used before classifying with the random forest. The number of trees in the model is set to 100 and all of the other model parameters are used as default. Results of the measures and the confusion matrix are presented in Table 2 and Figure 4, respectively.

The value of the precision score (0.93) indicates that the model has performed well to detect true positive cases. Recall score on the other and, indicates that some of the fraud cases could be missed when it is equal to 0.82. In case a fraud case is predicted as non-fraud, the results could harm the financial institutions. Lastly, the value of F1 score (0.88) which

shows the balance between the precision and recall scores indicates that the model performs well on detecting both the fraud and non-fraud cases.

**Table 2:** Performance Metrics of Random Forest - Over-sampling

Class	Precision	Recall	F1-Score
Non-Fraud	0.92	0.99	0.95
Fraud	0.93	0.82	0.88



**Figure 4:** Confusion Matrix of Random Forest - Over-sampling

### 5.1.2 Random Forest – Under-sampling

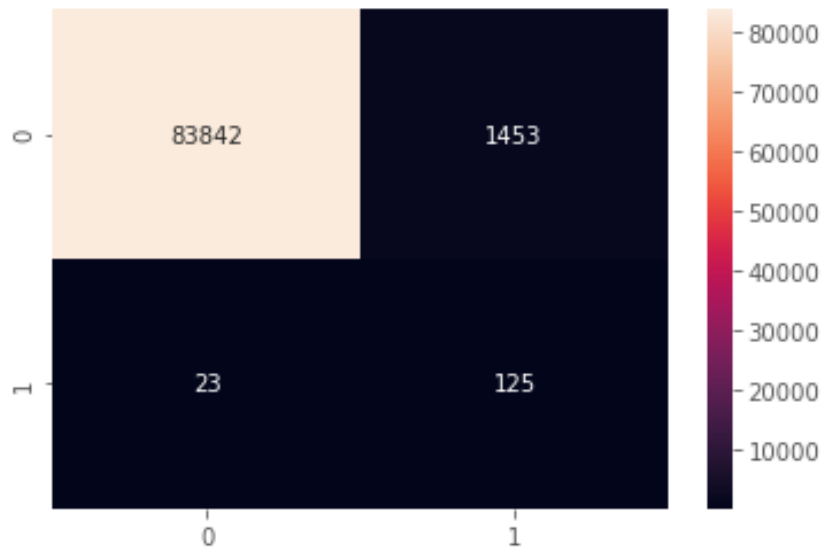
Before classifying with the random forest, the random standart under-sampling method with sampling strategy equal to 1 is implemented. The number of trees in the model is set to 100 and all of the other model parameters are used as default. Results of the measures and confusion matrix are presented in Table 3 and Figure 5, respectively.

From the results of the measures, it can be seen that, the model has predicted many regular cases as fraud where Precision is equal to 0.08. Real fraud cases are however predicted relatively well with Recall equal to 0.84. Reporting several regular cases as fraud would result as customer dissatisfaction for financial institutions. Lastly, the value of F1

score (0.14) suggests that the model cannot perform well as the model predicts many cases as false positive and false negative.

**Table 3:** Performance Metrics of Random Forest – Under-sampling

Class	Precision	Recall	F1-Score
Non-Fraud	1	0.98	0.99
Fraud	0.08	0.84	0.14



**Figure 5:** Confusion Matrix of Random Forest – Under-sampling

## 5.2 Decision Tree

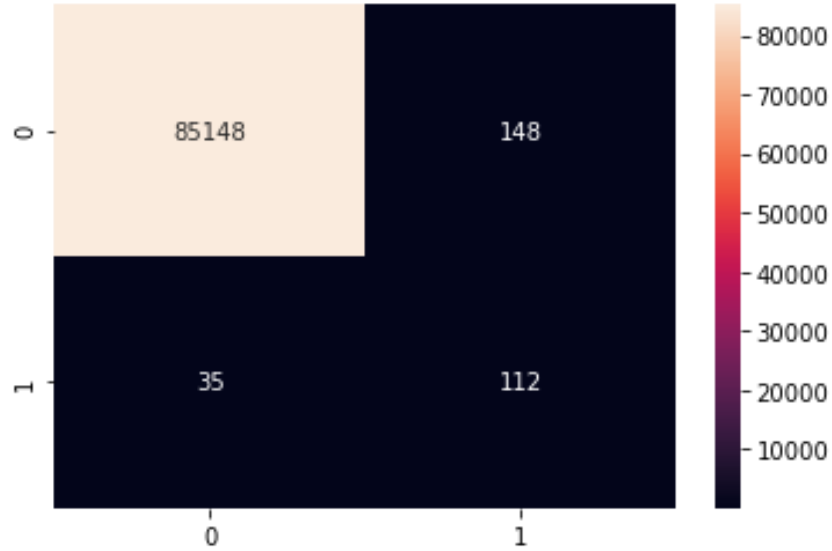
### 5.2.1 Decision Tree – Over-sampling

The random SMOTE method is used before classifying with the decision tree.. All of the model parameters are used as default. Results of measures and confusion matrix are presented in Table 4 and Figure 6, respectively.

From the results of measures, it can be seen that the model has predicted some regular cases as fraud with Precision equal to 0.43. However, real fraud cases are been predicted relatively correct with Recall equal to 0.76. The value of F1 Score (0.55) indicates that the decision tree model fitted with the SMOTE scaled dataset does not perform well as there are several cases that predicted as false positive and false negative.

**Table 4:** Performance Metrics of Decision Tree – Over-sampling

Class	Precision	Recall	F1-Score
Non-Fraud	1	1	1
Fraud	0.43	0.76	0.55



**Figure 6:** Confusion Matrix of Decision Tree - Over-sampling

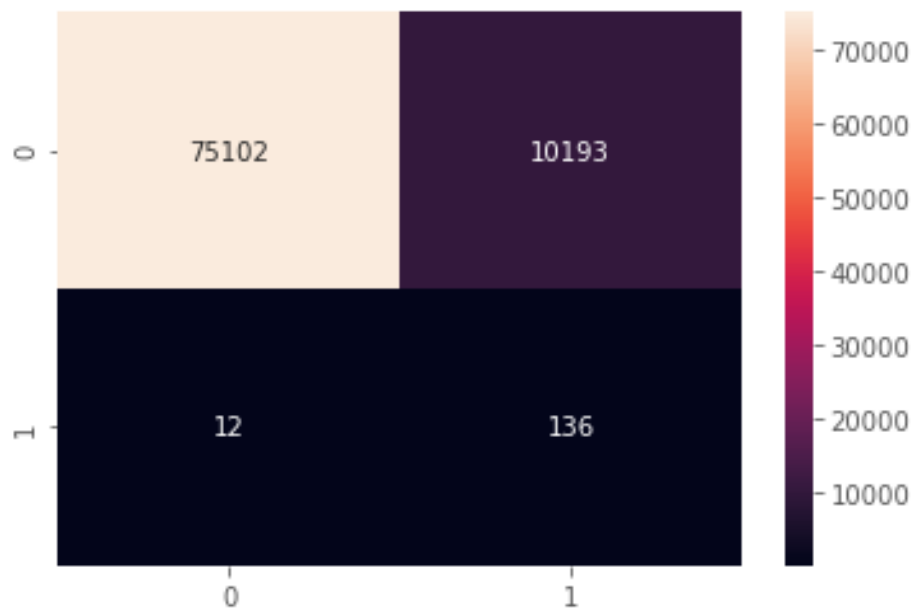
### 5.2.2 Decision Tree – Under-sampling

Before classifying with the decision tree, the random under-sampling method with sampling strategy set to 1 is used. All of the model parameters are used as default. Results of measures and confusion matrix are presented in Table 5 and Figure 7, respectively.

The model performs well by means of detecting actual fraud cases with Recall equal to 0.92. However, the model has predicted too many regular cases as fraud with Precision equal to 0.01 and F1 score equal to 0.03. The model clearly does not perform well because of its lack of precision.

**Table 5:** Performance Metrics of Decision Tree – Under-sampling

Class	Precision	Recall	F1-Score
Non-Fraud	1	0.88	0.94
Fraud	0.01	0.92	0.03



**Figure 7:** Confusion Matrix of Decision Tree – Under-sampling

### 5.3 Logistic Regression

#### 5.3.1 Logistic Regression - Over-sampling

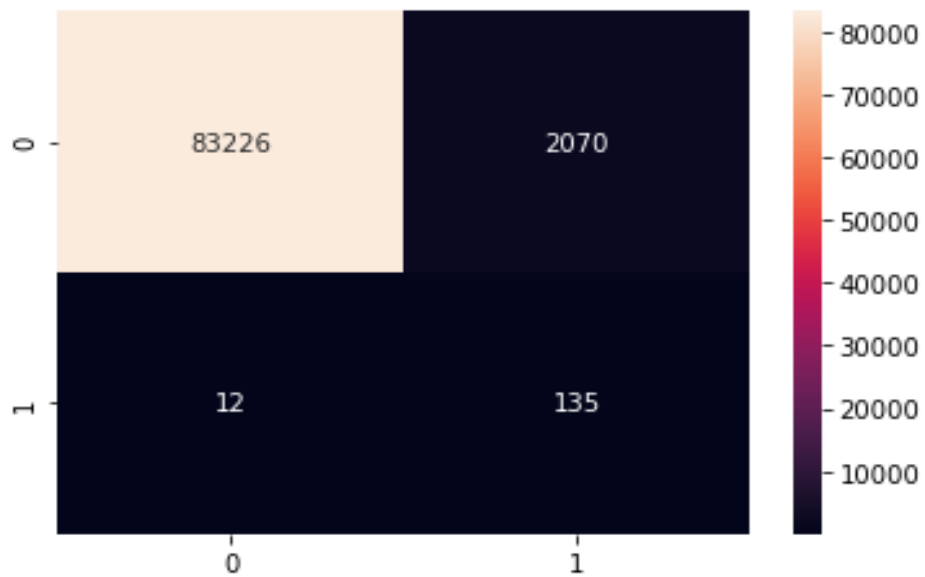
Before classifying with the logistic regression, the over-sampling method SMOTE is used. All of the model parameters are used as default. Results of measures and confusion matrix are presented in Table 6 and Figure 8, respectively.

The model performs well by means of detecting actual fraud cases with Recall equal to 0.92. However, the model has predicted too many regular cases as fraud with Precision equal to 0.06. As a result of lack of precision, F1 score is found as 0.11 and the model clearly does not perform well.

**Table 6:** Performance Metrics of Logistic Regression – Over-sampling

Class	Precision	Recall	F1-Score
Non-Fraud	1	0.98	0.99
Fraud	0.06	0.92	0.11





**Figure 8:** Confusion Matrix of Logistic Regression - Over-sampling

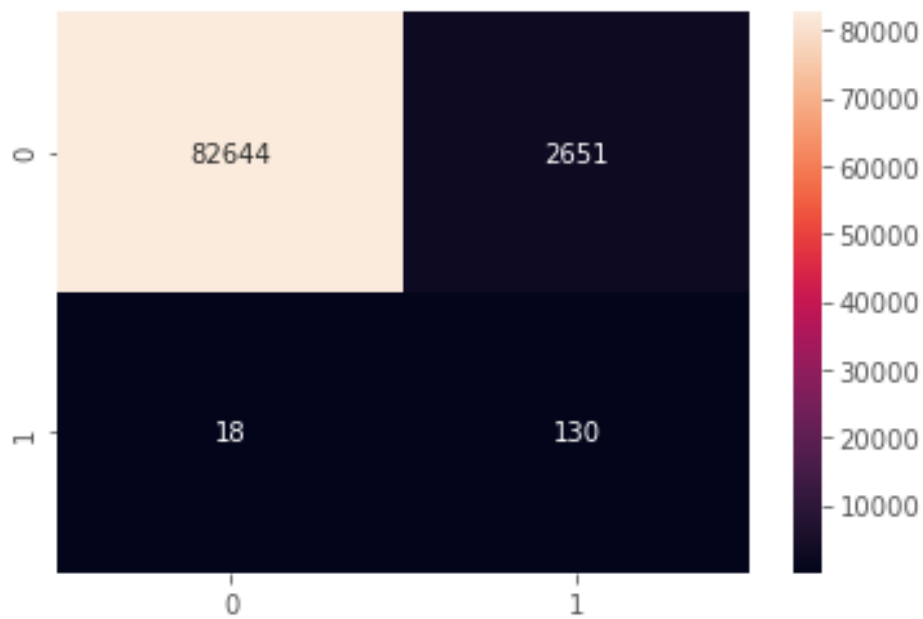
### 5.3.2 Logistic Regression – Under-sampling

Before classifying with the logistic regression, the random under-sampling method with sampling strategy equal to 1 is used. Results of measures and confusion matrix are presented in Table 7 and Figure 9, respectively.

Model perform relatively well by means of detecting actual fraud cases with Recall equal to 0.88. However, the model has predicted too many regular cases as fraud with Precision equal to 0.05. From the value of F1 score (0.09), it can be concluded that the model does not perform well because of its lack of precision.

**Table 7:** Performance Metrics of Logistic Regression – Under-sampling

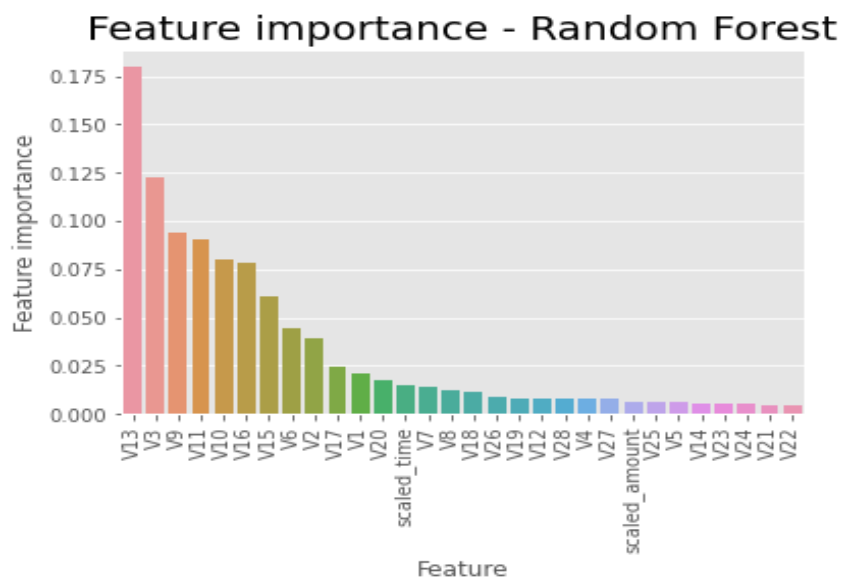
Class	Precision	Recall	F1-Score
Non-Fraud	1	0.97	0.98
Fraud	0.05	0.88	0.09



**Figure 9:** Confusion Matrix of Logistic Regression – Under-sampling

### 5.4 Feature Importance

The analysis of the feature importance conducted on Random Forest Classifier which is employed with the over-sampling technique. It can be seen from Figure 4 that by means of the random forest classifier, the feature “V13” has the most important role in detecting fraud. Also, the feature “V3” and the feature “V9” play important role in detecting fraud. It is also found that by means of decision tree and logistic regression methods, the feature “V13” play the most important role in detecting fraud too.



**Figure 10:** Feature Importance of Random Forest

## 6. CONCLUSION

The objective of this project is to find the most efficient machine learning models to detect fraudulent transactions on credit cards. In this work, different machine learning models that make classification of credit card transactions have been used on historical data. As a result, outcome of performance metrics may vary on real time data.

As the number of fraud transactions is much less than non-fraud transactions, SMOTE and random under-sampling feature scaling methods are used to deal with such an imbalance. The logistic regression, random forest and decision tree algorithms are used to find most effective performing model for the dataset subject to this project.

The random forest classifier is found as the most effective model with F1 equal to 0.88, Recall equal to 0.82 and Precision equal to 0.93. According to the results, the SMOTE technique used to overcome the data imbalance performs better than the under-sampling technique. The models employed with the under-sampled data misclassify a large number of non-fraud transactions as fraud.

Lastly, by means of the random forest, it is observed that the feature “V13” has the most important role in detecting fraud. As original features and background information are not provided further study on features cannot be implemented.

## REFERENCES

- [1] Roy, Baijayanta (2020). All about Feature Scaling Available: <https://towardsdatascience.com/all-about-feature-scaling-bcc0ad75cb35> [Accessed: 2-October-2020].
- [2] Patro, S Gopal & Sahu, Kishore Kumar. (2015). “Normalization: A Preprocessing Stage”, *International Advanced Research Journal in Science, Engineering and Technology (IARJSET)*, Vol 2, pp. 20-22.
- [3] Muralidharan, K. (2010). “A Note on Transformation, Standardization and Normalization”, *The IUP Journal of Operations Management*, Vol 9, pp. 116-122.
- [4] Brownlee, Jason. (2020) Data Preparation for Machine Learning: Data Cleaning, Feature Selection, and Data Transforms in Python.
- [5] fbi.gov, ‘Credit Card Fraud’[Online]. Available: <https://www.fbi.gov/scams-and-safety/common-scams-and-crimes/credit-card-fraud> [Accessed: 10- October- 2020].
- [6] Mitali Bansal, HCE Sonapat (2014). “Survey Paper on Credit Card Fraud Detection”, *International Advanced Research Journal in Science, Engineering and Technology (IARJSET)*, pp. 827-832.
- [7] Shen, Aihua Tong, Rencheng Deng, Yaochen. (2007). “Application of Classification Models on Credit Card Fraud Detection”, International Conference on Service Systems and Service Management, pp. 1-4.
- [8] Sahin, Yusuf and Duman, Ekrem. (2011). “Detecting credit card fraud by ANN and logistic regression”, International Symposium on Innovations in Intelligent Systems and Applications (INISTA), pp. 315-319.
- [9] Yeşilkanat, Ali and Bayram, Barış and Koroğlu, Bilge & Arslan, Seçil. (2020). “An Adaptive Approach on Credit Card Fraud Detection Using Transaction Aggregation and Word Embeddings”, *Artificial Intelligence Applications and Innovations*, pp. 3-14
- [10] Rokach, Lior and Maimon, Oded. (2005), “Decision Trees”, *The Data Mining and Knowledge Discovery Handbook*, Vol. 6, pp.165-192.
- [11] Sahayasakila.V, D. Kavya Monisha, Aishwarya, Sikhakolli VenkatavikalakshiseshsaiYasaswi (2019). “Credit Card Fraud Detection System using Smote Technique and Whale Optimization Algorithm”, *International Journal of Engineering and Advanced Technology*, Vol 8, pp. 162-164.
- [12] Ping Shung, Koo (2018). Accuracy, Precision, Recall or F1 Available: <https://towardsdatascience.com/accuracy-precision-recall-or-f1-331fb37c5cb9> [Accessed: 15-November-2020].
- [13] Suryanarayana, S and Gn, Balaji and Rao, G. (2018). “Machine Learning Approaches for Credit Card Fraud Detection”, *International Journal of Engineering and Technology (UAE)*, Vol 7.
- [14] Georgieva, Sevdalina and Markova, Maya and Pavlov, Velisar. (2019), “Using neural network for credit card fraud detection”, *AIP Conference Proceedings*, Vol 2159.
- [15] Jonnalagadda, Vaishnave and Gupta, Priya and Sen, Eesita (2019). “Credit card fraud detection using Random Forest Algorithm”, *International Journal of Advance Research, Ideas and Innovations in Technology*, Vol 5, pp 1797-1801.