# Credit Card Fraud Detection

# Using Machine Learning

Authors name: Sk Ibrahim

Department of computer science

Amity University

Lucknow 226 010, India

Email: imranshaik1902@gmail.com

**ABSTRACT:**

The physical loss or loss of sensitive credit card information raises some fraud cases known as credit card fraud detection. Many machine learning algorithms are used for detection. It is presently the most frequently occurring problem in the present world. This often occurs in both online transactions and e-commerce platforms, it generally happens when a credit card was stolen for any of the illegal purposes or even when the fraud people use the credit card information for his use. In present-day to day life, we are facing a lot of problems regarding credit card issues. The credit card fraud detection system was introduced to detect fraudulent activities. This project focuses mainly on machine learning algorithms. The algorithms used are the Random Forest algorithm (isolation forest) and the Local Outlier Factor algorithm. The results of the two algorithms are shown accurately. The isolation Forest algorithm is considered the best algorithm that is used to detect fraud.

This research shows several algorithms that can be used for classifying transactions as fraud or genuine one also a technique for Credit Card Fraud Detection is developed. As fraudsters are increasing day by day. And false transactions are done by credit card, and our goal is to detect fraud by filtering the above techniques of machine learning to get a better result.

Credit card frauds are easy and friendly targets. As we see everywhere the online payment modes have increased such as in E-commerce and many other online sites also increasing the risk for online frauds. Increase in fraud rates, different machine learning methods are being started by the researchers to detect and analyze frauds in online transactions, with an objective, in order to extract behavioral patterns from customers past transaction details. Where cardholders are gathered into different groups based on their transaction amount. Later different classifiers are trained over the groups separately. And the best methods to predict frauds can be chosen by the classifier with a better rating score.

This, followed by a feedback mechanism to solve the problem of concept drift. In this paper, we worked with the European credit card fraud dataset.

Keywords: fraud; logistic regression; Naïve Bayes; machine learning algorithms; Isolation Forest and Local Outlier Factor algorithms.

## I.    Introduction

A credit card is quite useful for almost everyone in day-to-day life. The focus of the fraud detection system is to detect fraud accurately and before the fraud is going to happen. There are some methods for detecting Credit Card Fraud using machine learning algorithms.

Credit card fraud is a growing concern in the present world with the growing fraud in government offices, corporate industries, finance industries, and many other organizations. In the present world, the high dependency on the internet is the reason for an increased rate of credit card fraud transactions but the fraud has increased not only online but also offline transactions. Though the data mining techniques are used the result is not much accurate to detect these credit card frauds. The only way to minimize these losses is the detection of the fraud using efficient algorithms which is a promising way to reduce credit card fraud. As the use of the internet is increasing, a credit card is issued by the finance company. Having a credit card means that we can borrow funds. The funds can be used for any of the purposes. When coming to the issuance of the card, the condition involved is that the cardholder will pay back the original amount they borrowed along with the additional charges they agreed to pay.

## II. Related work

With the increasing vogue of the internet, everything is available at our doorstep and convenience.[1] The increase in e-commerce has resulted in the increased usage of credit cards for online and offline payments. Though there are various benefits of using credit cards such as convenience, instant cash when it comes to security cardholders, banks, and the merchants are affected when the card is being lost and misused without the knowledge of the cardholder (Fraud activity) so this is the major loss causing by the fraudulent activities, which motivated researchers to find a solution that would detect and prevent frauds. Several methods have already been proposed and tested. Some of them are briefly reviewed below. Algorithms such as Gradient Boosting, Support Vector Machines, Decision Tree (DT), LR, and RF proved useful. In paper [5] GB, LR, RD, SVM, and a combination of certain classifiers were used, which led to the high recall of over 91% on a European dataset. High accuracy and recall were achieved only after balancing the dataset by scanning the data. In paper [6], a European dataset was also used, and a comparison was made between the models based on LR, DT, and RF. Among the three models, RF proved to be the best, with an accuracy of 95.5%, followed by DT with 94.3%, and LR with an accuracy of 90%.

## III. Machine Learning

*Machine learning and it's types:*

Machine learning is a study that a computer program can learn by itself and adapt to new data without human intervention that uses statistical model's algorithms to analyze and draw conclusions from patterns in data. The systems built on machine learning algorithms have the potential to learn from experience or historical data.

Machine learning can be divided into three categories:

a) **Supervised learning**

In supervised learning, the algorithm deals with the labeled dataset, also provides an answer key that the algorithm can use to calculate its accuracy on training data. [7] in this machine is trained on a dataset that has both the input as well as the output. So, after completing the training and the machine has acquired a certain level of a learning machine is finally deployed. Supervised learning is further classified into two classifications and regression.

Examples of Supervised learning: for example, you have a niece who is just 2 years old and is learning to speak. You want to teach her what a dog and a cat are. So, what do you do? You will show her videos of dogs and cats, or you bring a dog and a cat and show them to her in real life so that she can understand how they are different. Now there are certain things you tell her so that she can understand the differences between the 2 animals.

- Both dogs and cats have 4 legs and a tail.
- Dogs have long mouths while cats have smaller mouths.
- Dogs come in small to large sizes. But Cats, are always small.
- cat's meow while Dogs bark.

This is how u train the system first by giving some data. so that the computer will differentiate between the dogs and cats.

Applications of Supervised learning:

- Cortana, Siri, and Alexa are voice assistants that trained using our voice.

a) **Unsupervised learning**

In unsupervised learning, this is the technique in which there is no need to teach the model or supervise the model. The model itself works on its own to discover the patterns and insights that were undetected before. It deals with unlabeled data. Based on the observations the values are predicted in the future.

Examples of unsupervised learning:

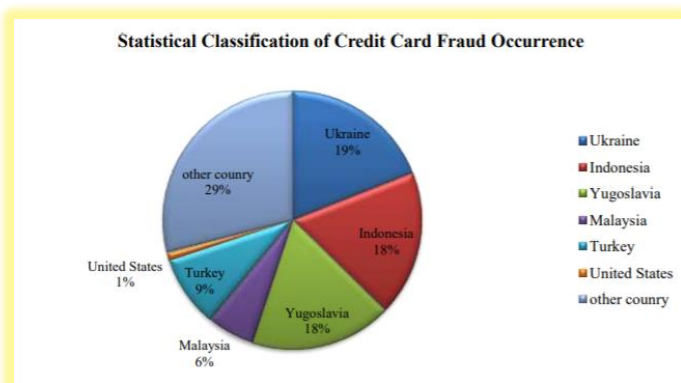- Apriori Algorithm, K-means Algorithm, Hierarchical Clustering.

Let's, take the case of a baby again and her family have a dog. She can identify this dog. A few weeks later a friend brings another dog and tries to play with the baby. Baby hasn't seen this dog earlier. But she recognizes many features (2 ears, eyes, walking on 4 legs) are like her pet dog. She identifies a new animal like a dog.

This is unsupervised learning, where you are not taught by anyone, but you can learn from the data (in this case data about a dog.) this would have been supervised learning if the friend would have told the baby that it's a dog.

## IV. Credit card fraud

Credit card fraud refers to the unauthorized use of a credit card or its information without the knowledge of the owner. Different credit card fraud tricks belong mainly to two groups of application and behavioral fraud [3]. There are two types of fraud one is duplicate fraud and another one is identity fraud. Submitting multiple applications by one user with one set of user details is called duplication fraud or different users with identical details is called identity fraud. Behavioral fraud, on the other hand, has four principal types: lost/stolen card, cardholder not present" fraud Stolen/lost card fraud occurs when fraudsters steal a credit card or get access to a lost card. Mail theft fraud occurs when the fraudster gets a credit card in the mail or personal information from the bank before reaching to actual cardholder [3]. In the former, remote transactions can be conducted using card details through the mail, phone, or the Internet. In the latter, counterfeit cards are made based on card information. Based on statistical data stated in [1] in 2012, the high risk of credit card fraud threat is faced by the countries shown in Fig.1. Ukraine has the most fraud rate with a staggering 19%, which is closely followed by Indonesia at an 18.3% fraud rate. After these two, Yugoslavia with the rate of17.8% is the riskiest country. The next highest fraud rate belongs to Turkey with 9%, Malaysia with 5.9%, and finally the United States. Other countries that are prone to credit card fraud with a rate below 1% are not demonstrated in figure 1.

Fig1. High risk countries facing credit card fraud threat



## V. Difficulties of Credit Card Fraud Detection

Fraud detection systems are snip and thin to several difficulties and challenges. An effective fraud detection technique should have abilities to address these difficulties to achieve the best performance.

- Imbalanced data. The credit card fraud detection data has imbalanced nature. It means that very small percentages of all credit card transactions are fraudulent. Therefore, this tends the detection of fraud transactions very difficult and accurate.

- Different misclassification importance: in the fraud detection task, different misclassification errors have different importance. Misclassification of a normal transaction as fraud is not as harmful as detecting a fraud transaction as normal.

- Overlapping data: many transactions may be considered fraudulent, while they are normal (false positive) and reversely, a fraudulent transaction may also seem to be false negative. Hence obtaining a low rate of a false positive and false negative is a key challenge of fraud detection systems [4, 5, and 6]

- Less Adaptability: classification algorithms are mainly faced with the problem of detecting new types of normal or fraudulent patterns. The supervised and unsupervised fraud detection systems are inefficient in detecting new patterns of normal and fraud behaviors, respectively.

- Fraud detection cost: For example, no income is obtained by stopping a fraudulent transaction of a few dollars [5, 7]. Ukraine 19% Indonesia 18% Yugoslavia 18% Malaysia 6% Turkey 9% United States 1% another country 29% Statistical Classification of Credit Card Fraud Occurrence Ukraine Indonesia Yugoslavia Malaysia Turkey United States another country 5.

- Lack of standard metrics: there is no standard evaluation criterion for assessing and comparing the results of fraud detection systems.

## VI. Materials and Methods

In this research, the credit card detection dataset was used which can be downloaded from Kaggle. [8] This dataset contains transactions that occurred in two days, made in September 2013 by European cardholders. The dataset contains 31 numerical features. Since some of the input variables contain financial information, the PCA transformation of these input variables was performed to keep these data anonymous. Three of the given features weren't transformed. Feature "Time and Amount " shows the time between the first transaction and every other transaction in the dataset and the amount of the transactions made by credit card. Feature "Class" represents the label and takes only 2 values: value 0 in case of normal transaction and 1 in case of fraud.

The dataset contains a total of 284,807 transactions of which 492 transactions were frauds and the rest were normal. Considering the numbers, we can see that this dataset is highly imbalanced, where only 0.173% of transactions are labeled as frauds.

In this research project, we have used two algorithms to detect the outliers there are isolation forest algorithm and the local outlier factor. These are some of the newest techniques to detect anomalies.

## VII. Algorithms used

(a) Isolation forest

It is the most common technique and works on the principle of the decision tree algorithm. [9] Anomalies are found as those instances of data that do not conform to the defined normal profile. However, the isolation forest does not work on the above process. Isolation forest works as an unsupervised machine learning algorithm. It identifies anomalies by isolating outliers in the data It detects the anomalies fast and it requires less memory compared to other anomaly detection algorithms. Also, this is one of the best advantages of using the isolation forest.

At first, it isolates the outliers found in it by randomly selecting a feature from the given data set of features and then randomly selecting a split value between the maximum and minimum ranges of the selected feature. Isolation forest works on the principle of recursion. This algorithm also works on the recursive method as it recursively generates partitions on the datasets by

randomly selecting a feature and then randomly selecting a split value for the feature.

(b) Local Outlier Factor

The local outlier factor (LOF) [10] is an algorithm that helps to identify the outliers present in the dataset. When a point is considered as an outlier based on its local neighborhood, it is a local outlier. LOF identifies an outlier considering the density of the neighborhood. [11] The Direct Outlier (LOF) operator is available in Data Transformation > Data Cleansing > Outlier Detection. The output of the Local Outlier Factor operator contains the example set along with a numeric outlier score. The LOF algorithm does not explicitly label a data point as an outlier; instead, the score is exposed to the user.[12] The number of neighbors considered, (parameter neighbors) is typically chosen 1) greater than the minimum number of objects 2) smaller than the maximum number of close-by objects that can potentially be local outliers.

## VIII. Results and discussion

Some of the results that we get after using these two algorithms: basically, at first we can see how many classes are there with respective the frequency
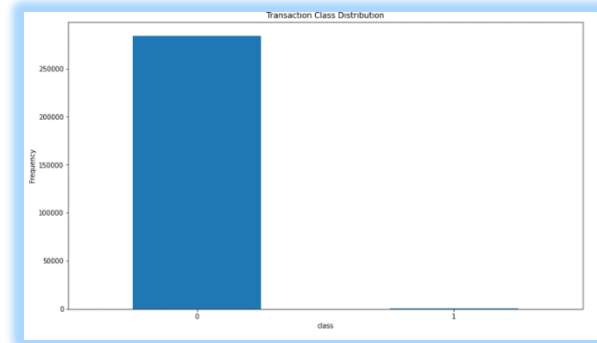


Fig 2: transaction class distribution

Here, we come to know that normal transaction is more than 250000, but the fraudulent transactions are quite less we can see in the figure.

Also, we can get to know the total cases of fraud and normal, so there are 284315 normal cases and 492 fraud cases, and also we can describe them both by Fraud.Amount.describe(), Normal.Amount.describe()

```
In [11]: print('Fraud cases:', format(len(Fraud)))
         print('Normal cases:', format(len(Normal)))

         Fraud cases: 492
         Normal cases: 284315

In [12]: 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
```

Fig 3. *Total fraud and normal cases*

As we can see, this is an imbalance dataset, so we are using the algorithms like Isolation Forest and Local Outlier Factor. Also, we have done some visual representations to get more understand



```
In [15]: 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();
```
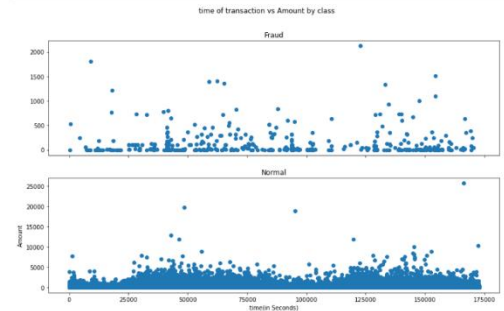
Fig 4. *time of transaction vs amount by class*

So, we can see there are lot of transaction with respective time but fraudulent are less.

Apart of that we can also do correlation and try to find out how all the features are with respective the class variable. So, for this we can create the heat map to get more points visually
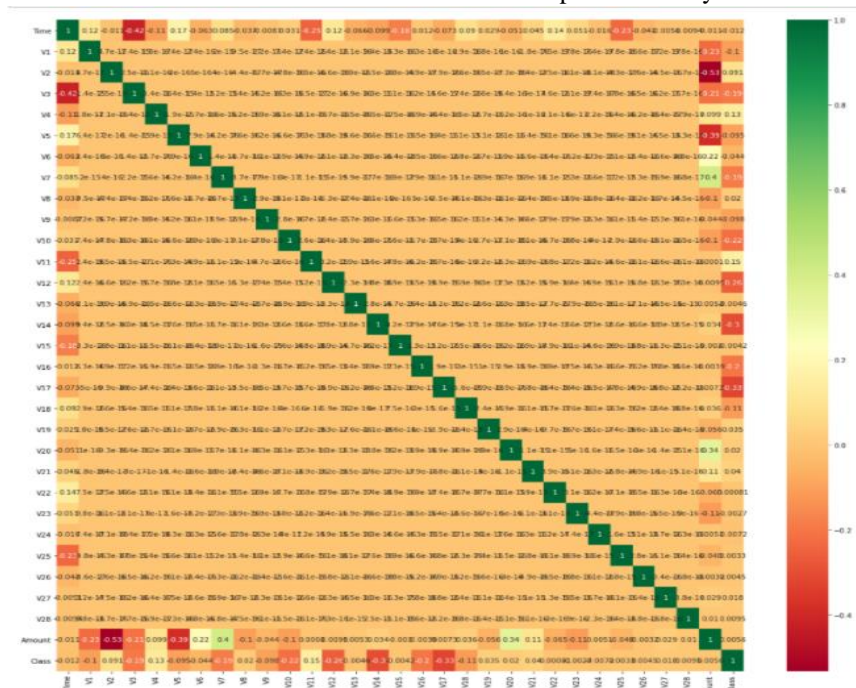


Fig 5. *Heat map of the given data*

After this we have created the dependent and independent features to apply the models because this is an imbalanced data set so it can be taken care by algorithms. On the other hand, normal observations require more conditions to isolate. Therefore, we can calculate the anomaly score:
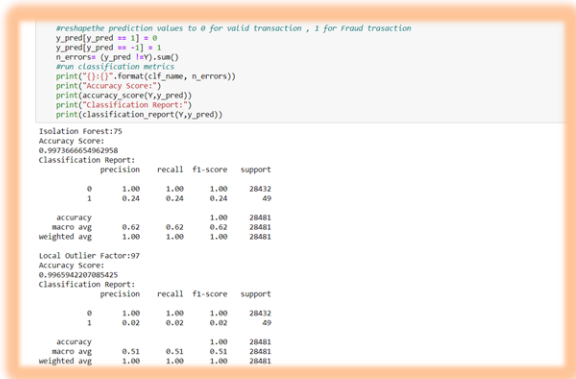
Fig 6: *accuracy and score*

*Code of the project:* Github (2019) Credit-Card-Fraudlent [online] available at:
https://github.com/krishnaik06/Credit-Card-Fraudlent/blob/master/Anamoly%20Detection.ipynb

**Observations:**

IX. Isolation Forest detected 73 errors versus Local Outlier Factor detecting 97 errors vs. SVM detecting 8516 errors

X. Isolation Forest has a 99.74% more accurate than LOF of 99.65% and SVM of 70.09

XI. 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 the LOF detection rate of just 2 % and SVM of 0%.

XII. So overall Isolation Forest Method performed much better in determining the fraud cases which is around 30%.

XIII. We can also improve on this accuracy by increasing the sample size or using 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.

**References**

[1] International Journal of Engineering Research & Technology (IJERT) ISSN: 2278-0181 Volume 8, Issue 08 Special Issue - 2020 Published by, www.ijert.org NCICCT - 2020 C [accessed 11 mar 2021]

[2]KhyatiChaudhary,JyotiYadav, BhawnaMallick, " A review of Fraud Detection Techniques: Credit Card", International Journal of Computer Applications Volume 45– No.1 2012.

[3] Michael Edward Edge, Pedro R, Falcone Sampaio, journal of computers and security, Vol. 28,

[4] Linda Delamaire, Hussein Abdou, John Pointon, "Credit card fraud and detection techniques: a review", Banks and Bank Systems, Volume 4, Issue 2, 2009.

[5] Salvatore J. Stolfo, David W. Fan, Wenke Lee and Andreas L. Prodromidis; "Credit Card Fraud Detection Using Meta-Learning: Issues and Initial Results"; Department of Computer ScienceColumbia University; 1997.

[6] Maes S. Tuyls K. Vanschoenwinkel B. and Manderick B.; "Credit Card Fraud Detection Using Bayesian and Neural Networks"; Vrije University Brussel – Belgium; 2002.

[7] Athey, S. (2018). The impact of machine learning on economics: An agenda University of Chicago Press.

[8] Andreas L. Prodromidis and Salvatore J. Stolfo; "Agent-Based Distributed Learning Applied to Fraud Detection"; Department of Computer Science-Columbia University; 2000.

[9] Salvatore J. Stolfo, Wei Fan, Wenke Lee and Andreas L. Prodromidis; "Cost-based Modeling for

Fraud and Intrusion Detection: Results from the JAM Project"; 0-7695-0490-6/99, 1999 IEEE

[10] Kaggle.com. (2019). Credit Card Fraud Detection. [online] Available at: https://www.kaggle.com/mlg-ulb/creditcardfraud [Accessed 10 Jan. 2019].

[11] Prakash verma (2020) Detecting anomalies using tree-based algorithms [online] available at: https://heartbeat.fritz.ai/isolation-forest-algorithm-for-anomaly-detection-2a4abd347a5# [accessed 10 march 2020]

[12] vaibhav jayaswal (2020) Algorithm for outlier identification [online] available at : https://towardsdatascience.com/local-outlier-factor-lof-algorithm-for-outlier-identification-8efb887d9843 [accessed 10 march 2020 ]

[13]Vijay Kotu, Bala Deshpande PhD,
in Predictive Analytics and Data Mining, 2015

[14] Github (2019) Credit-Card-Fraudlent [online] available at: https://github.com/krishnaik06/Credit-Card-Fraudlent/blob/master/Anamoly%20Detection.ipynb [accessed 10 mar 2021]