

Design of a Model in Machine Learning For Credit Card Fraud Detection

Abderrahim Bouchouq, Wang Peiliang
Engineering College of Huzhou Normal University

E-mail: abdobour97@gmail.com

Abstract

In e-commerce, fraud has become a major problem, and much effort is being invested into identifying and preventing it. Currently, fraud detection and prevention systems are only able to detect and prevent a small percentage of fraudulent transactions, resulting in billions of dollars in losses. Because online transactions are likely to grow dramatically in the coming year, better fraud detection and prevention is critical. We provide a data-driven strategy for estimating the likelihood of a fraudulent or legal transaction based on machine learning techniques applied to large data sets. To predict the likelihood of a customer's next transaction being fraudulent, the algorithm was trained using past e-commerce credit card transaction data. Random Forest, Support Vector Machine, Gradient Boost, and combinations of these are used to compare the performance of supervised machine learning approaches. The problem of class imbalance is considered before the model is trained on a classifier, and methods such as oversampling and data pre-processing are employed.

Keywords: machine learning, model, fraud, credit card

DOI: 10.7176/CEIS/13-2-06

Publication date: April 30th 2022

1.0 Introduction

Credit card fraud detection is a prominent problem in finance research, with considerable economic implications. While traditional data analysis methods have been employed in the past, the similarities between this and other difficulties, such as the development of recommendation systems and diagnostic/prognostic medical tools, suggest that a complex network approach could give considerable benefits. It is feasible to find unauthorized instances in a real card transaction data set. It's based on a recently developed network reconstruction method that enables the creation of representations of a single instance's divergence from a reference group. As a result, credit card companies must be able to detect fraudulent credit card transactions so that customers are not charged for items they did not purchase. This is a list of fraudulent transaction statistics in China. When compared to 2017, the number of prosecutions for funding fraud increased by about 8.17 percent in 2019. The number of financial fraud prosecutions accepted by Chinese courts decreased between 2016 and 2020. The term "financial fraud" refers to crimes committed with the use of credit cards, insurance, securities, or other financial instruments. As part of the synthesis project in Big Data and Statistical Learning, we did the following study on one of today's most pressing banking issues. As a result, every bank and financial services provider in China and abroad is working to protect credit card operations to assure a more stable and trusted transaction system. Because fraud can be exceedingly damaging to both clients and service providers, fraud detection is a critical topic to work on. Engineers interested in finance may find that constructing fraud detection models is one of our daily project duties, resulting in the subject's striking aforementioned desire. Assume you've been engaged to assist a credit card company in detecting possible fraud cases so that customers aren't charged for items they didn't purchase. You're given a dataset containing transactions between people, as well as information on whether or not they're fake, and you have to tell them apart. This is the case we're going to deal with. Our ultimate goal is to solve this problem by creating classification models that can classify and differentiate fraud transactions.

1.2 Motivation for the Paper

From bridges to satellites, all modern engineering products require some form of monitoring to ensure that they function as intended. This monitoring could be required for safety reasons, or it could be used to improve operating efficiency by reducing energy consumption or assisting with maintenance planning to optimize system usage. Traditionally, someone with particular training has performed this monitoring. This permits them to examine the system and make health assessments. This could be for a variety of reasons in a variety of situations. Even a seasoned professional may not be able to fully comprehend the system. The average time between failures may be so long that having a specialist on hand to keep track of it isn't cost-effective. There's a chance that no specialists

are available. The system may be too new for the expert to have gained sufficient experience with it. The risk of a diagnosing error cannot be accepted since the implications of failure are so severe. Complex monitoring and diagnostic systems are required because the business environment may be too hostile for human examination

2.0 Review of Literatures

2.1 Review Other Models

Detecting fraudulent transactions can be done in a variety of ways. It's exceedingly tough to detect fraud, and it's only possible after it's happened. Because fraudulent transactions are minor in comparison to total transactions, this is the case. In their paper [8], the authors compared seven methods for detecting such transactions. ANN had the best results across the board, with an accuracy of 99.71 percent, a detection rate of 99.68 percent, and a false alarm rate of 0.12%. ANN is the most efficient, even though it takes the most time and processing power to train. SVM has a maximum false alarm rate of 5.2% and a detection rate of 85.45%, which is not comparable to other superior techniques. Fuzzy logic has the lowest detection rate, at 77.8 percent. With an accuracy of 97.93 percent, a detection rate of 98.52 percent, and a false alarm rate of 2.19 percent, decision trees are weighted in favor of complexity during training. Random forest is a categorical and numerical data decision tree regression and classification tool [6]. The authors employed a random forest and an SVM classifier to detect fraudulent transactions in the dataset. Pre-processing was done to avoid missing data and scale feature values. The authors discovered that SVM did not perform well with imbalanced data when compared to random forest classifiers. Another benefit of using the random forest technique was that, because it used a subset of data with different decision trees, adding more data points did not affect the model. Because each tree has a very small chance of influencing others, bias and overfitting are reduced. Random forest classification was used by Mohankumar and Karuppasamy [15] to detect fraudulent credit card transactions. The PCA approach was used to mask the values in the dataset. Feature values were adjusted to reduce volatility across characteristics. The SMOTE algorithm was used to balance the data. In the balanced data, there are 175000 classes. A random forest classifier is used to categorize data points into binary categories. According to the paper's findings, the precision-recall curve has an analogous value of around 0.85. The random forest classifier has been one of the most widely used techniques in e-commerce to detect credit card fraud due to its flexibility and scalability for large datasets. The computational resources required to train the random forest model are modest as compared to superior state-of-the-art techniques like ANN. ANN is not commonly employed in real-time e-commerce applications due to computational and time constraints. There is a significant problem with a class imbalance in the existing datasets that mislead research [17]. In their research, the authors looked at data balance for efficient analysis, regression, and classification challenges. The primary methodologies they looked into were random oversampling and undersampling, statistical oversampling and undersampling, SMOTE, Feature Selection, Hybrid Sampling, Cost-effective Learning, and Ensemble Learning. The SMOTE approach, as well as feature selection, were found to be commonly used in research papers [9, 14, 16]. These two tactics produce the best results when it comes to balancing obstacles in data analysis.

Razooqi et al. [19] advised using fuzzy logic to change weights to use a genetic algorithm with ANN, which led to even better results and a lower FN rate. Although the amount of time spent training increased drastically, the end result for ANN was far better. Maes et al. [12] created a Bayesian network for anticipating financial data labels. Even with small datasets, the model provided good results and utilizing ANN to alter the network's parameters cut training time in half. Shirgave et al. [20] introduced a supervised learning random forest to classify alerts as bogus or authentic, paving the way for the implementation of a semi-supervised machine learning method to classify alerts. According to Lakshmi and Kavila [10], a random forest classifier surpasses choice trees and logistic regression in terms of accuracy, with 90.0, 94.3, and 95.5 for logistic regression, decision tree, and random forest classifier, respectively. The random forest classifier outperforms logistic regression and decision trees in a comparison of the three approaches. Furthermore, as Havarapu Bhanusri et al. [2] point out, we can't discern the difference between fraudulent and lawful transactions using machine learning approaches for the current dataset, necessitating additional research. The choice of algorithms for a credit card fraud detection system, according to Sorournejad et al. [21], should minimize False Positive and False Negative rates while maximizing True Positive and True Negative rates and assuring a respectable detection rate. Integrating Genetic Algorithms (GAs) with ANNs for credit card fraud detection can improve engine performance, according to Carsten [3].

2.3 Payments Fraud

In the United States and around the world, payment fraud is a critical and growing problem. In the United States, fraud on non-cash payment systems totaled more than \$8 billion in 2015, up 37% from 2012 [18]. Financial institutions and payment operators are increasingly relying on machine learning algorithms to build efficient and effective fraud detection systems [24]. I implement and analyze the performance of various machine learning

models, such as logistic regression, random forests, and neural networks, using a big dataset from Kaggle. Around 300,000 credit card transactions were recorded across Europe over two days.

99.8% of transactions in the Kaggle dataset are regarded as genuine, whereas 0.2 percent are declared fraudulent. As a result of class imbalance, standard models may struggle to distinguish between the dominant and minority classes [3]. As part of this work, I analyze and assess various approaches to dealing with this problem, including sampling techniques such as under-sampling the majority class and over-sampling the minority class. Due to privacy concerns, it is impossible to construct meaningful correlations between most of the variables in the dataset. As a result, rather than inference, my research focuses on performance prediction.

Machine learning has a long history of being used to detect fraud in the payments industry. While banks and payment companies regularly use fraud algorithms, Bhattacharyya et al. point out that there are few research on the use of machine learning techniques for payment fraud detection [4], partially due to the sensitive nature of the data. According to their findings, random forests, despite being underutilized, may outperform more known approaches. The most critical aspect in a neural network's fraud classification ability, according to Roy et al. [5], is network size. According to Chaudhary et al., no single method is optimum across all performance parameters; each has its own set of strengths and limitations [6]. There is a lot of machine learning research on class imbalance. Some of the tactics used to address this problem include oversampling the minority class, undersampling the majority class, creating synthetic minority samples, and changing the relative cost of misclassifying the minority and majority classes. According to Japkowicz and Stephen, the effects of class imbalance vary on the degree of imbalance, sample size, and classifier used, with sampling procedures hurting the performance of particular classifiers [7]. Oversampling can lead to overfitting, and synthetic data creation adds noise that can reduce prediction performance, according to Cui et al. [8].

For a variety of reasons, detecting credit card theft is difficult. One of these is the likelihood of changing buying patterns over time. Because the approaches to effectively learn on the training sample may not correlate to the distribution of the testing set, the monitoring system may need adaptation as a result of this dataset shift or concept drift. Concept drift and dataset shift have previously been discussed in the literature. Others attempted to adapt to concept drift, while others attempted to define it. Abdallah & al. [9] define idea drift as a phenomenon in which the underlying model (or concept) evolves: buying behavior may change over time, and fraudsters' techniques may alter. Given the features of the testing set transactions, the decision function learned on the training set may become out-of-date and no longer represent the conditional distribution of the target variable (y) (X). They propose two approaches to dealing with this issue: the developing approach involves learners who keep up with the data stream, and the regulated approach involves detecting concept drift and intervening to change the hybrid technique when it occurs (usually retraining the learner).

2.3 Neural networks

Neural networks are based on the human brain, and their ability to learn has aided scientists and engineers in solving a range of problems. For identifying credit card fraud, C. Wang et al. [8] proposed a whale method based on neural networks. [10] Presented a cost-sensitive neural network-based approach for detecting credit card fraud. [16] utilized the Support Vector Machine to detect credit card fraud. To divide data into many categories or clusters, classification and clustering procedures are used. These techniques can be used to a wide range of problems. [5] identified credit card fraud using a meta classifier on a big dataset. To detect credit card fraud, [13] employed a system based on partitioning and clustering methods. [14] employed skewed data and a variety of classification algorithms to detect credit card fraud. [17] proposed a credit card fraud detection algorithm based on a classification algorithm.

Machine learning algorithms are artificial intelligence (AI) technologies that are used to solve problems involving massive amounts of data in a range of industries. Several research employed machine learning [4][6][7][22][15] and deep learning [23][15] techniques to detect credit card fraud. However, further research is needed, as is the application of machine learning algorithms to detect fraud in credit card transactions. Here are a few examples of how machine learning algorithms can be used.

3.0 Model Development

3.1 Approaches: *Explanatory Data Analysis*

Logistic Regression is a supervised classification approach that is used to model the likelihood that Y belongs to a specific category. The method of logistic regression is used to estimate discrete values from a set of independent variables. It assists you in predicting the likelihood of an event occurring by fitting data to a valid function. Its output value is between 0 and 1 because it forecasts probability.

Linear Discriminant Analysis (LDA) is a dimensional reduction technique for supervised classification problems. It's used to represent group differences, such as separating two or more classes.

K Nearest Neighbors (KNN) is a simple classification technique that identifies the query's closest neighbors and uses those neighbors to determine the query's class.

Support A discriminate classifier explicitly described by a separating hyperactive plane is known as a vector classifier. In other words, the algorithm produces an ideal hyperplane that categorizes fresh samples given labeled training data (supervised learning).

The Random Forest Classifier is a classification approach that employs the following five steps: (1) Choose "k" features at random from a total of "m" features (where $k < m$). (2) Using the optimal split point, calculate the node "d" among the "k" characteristics. (3) Using the best split, split the node into daughter nodes. (4) Repeat steps 1–3 until the "l" number of nodes is attained. (5) Build a forest by repeating steps 1 through 4 "n" times to create "n" trees.

3.2 Fraud Detection Framework

The architecture of the proposed methodology is depicted in Figure 1. As the first step in the Normalize Inputs block, the training dataset is normalized using the min-max scaling strategy in Equation (4) [31]. The scaling ensures that all of the input values fall inside a particular range. In the GA Feature Selection block, the normalized data from the Normalize Inputs block is used to implement the GA method. For each iteration of the GA Feature Selection block, the GA offers a candidate attribute vector v_n , which is utilized to train the models in the Training block, which is represented by the Training data and Train the model's blocks.

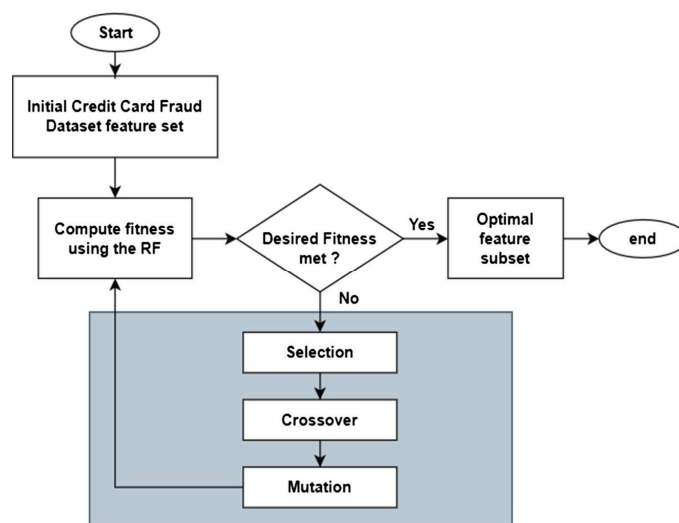


Figure 1 Flow chat of the detections System

3.3 Dataset

The credit card fraud transaction data set we're using comes from Online, and it contains 28315 transaction details, 0.5 fractions of them were found to be fake. The goal is to develop a classifier model that can be used to detect anomalies. The data collection contains only numerical input after PCA transformation. The major parts are features V1, V2, and V28, except for 'Time' and 'Amount,' which have not been adjusted using PCA.' 'Class,' the response variable, has a value of 1 when there is a scam and 0 if there's not.

Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14	V15	V16	V17	V18	V19	V20	V21	V22	V23	V24	V25	V26	V27	V28	Amount	Class
0	-1.359807	-0.072781	0.536346	0.738155	-0.338321	0.462387	0.239586	0.098699	0.363787	0.090794	-0.5516	-0.617801	-0.99139	-0.311169	1.68177	-0.47404	0.207912	0.025796	0.403993	0.251421	-0.018307	0.277837	-0.110474	0.066928	1.285394	-0.189115	0.133584	-0.021053	149.62	0
0	1.181875	0.266150	0.166480	1.448154	0.060117	-0.082361	-0.078803	0.085107	-0.235425	-0.166974	0.127267	0.065235	0.488095	-0.143772	0.635581	-0.463917	-0.114805	-0.183361	-0.145783	-0.069083	-0.225775	-0.638672	0.101288	-0.339846	0.167104	0.125845	-0.008983	0.014724	2.69	0
1	-1.358354	-1.340163	0.732083	0.179796	-0.503188	0.800498	0.781461	0.247678	-1.514654	0.207649	0.624015	0.066087	0.217297	-0.165946	0.345864	-2.890083	-1.099084	-0.121359	-2.261857	0.524979	0.247982	0.771679	0.909412	-0.689281	-0.327642	-0.139097	-0.055253	-0.059752	378.66	0
1	-0.966272	-1.185226	0.793993	-0.862391	-0.010309	0.247032	0.237608	0.377439	-1.387024	-0.054852	-0.226487	0.178282	0.507759	-0.287924	-0.631418	-1.059647	-0.684083	1.965775	-1.232622	-0.208038	-0.10830	0.055736	-0.190321	-1.175575	0.647376	-0.221929	0.062728	0.014576	123.5	0
2	-1.158233	0.877368	1.548718	0.403039	-0.407193	0.095925	0.592940	-0.270533	0.817793	0.753074	-0.822843	0.581956	1.345816	-1.119670	0.175121	-0.451449	-0.237033	-0.038195	0.803489	0.408542	-0.009431	0.798278	-0.137458	0.141267	-0.200110	0.502292	0.219422	0.215151	69.99	0
2	-0.425966	0.905253	1.141109	-1.168252	0.209894	-0.029728	0.476209	0.260314	-0.568671	-0.371407	1.341262	0.359898	-0.358091	-0.137134	0.517618	0.401725	-0.058133	0.068053	-0.033194	0.084677	-0.208254	-0.559825	-0.026398	-0.371427	-0.332794	0.105914	0.253844	0.081003	3.67	0
4	1.229657	0.140035	0.043708	1.202812	-0.191880	0.277801	-0.005159	0.081219	0.46496	-0.099254	-1.410907	-0.153826	-0.751063	-0.167372	0.050146	-0.436387	0.002820	-0.611987	-0.045575	-0.219631	-0.167716	-0.27071	-0.154104	-0.780055	0.750169	-0.257237	0.034507	0.005167	4.89	0
7	-0.644269	1.179635	0.074380	-0.492199	0.948934	0.428185	1.206314	-3.807864	0.615374	1.249376	-0.619480	0.291474	1.757964	-1.322865	0.686125	-0.076127	-1.222127	-0.358222	0.324507	-0.156742	0.943463	-1.015455	0.057503	-0.649709	-0.415267	-0.051634	-1.206921	-1.085339	40.8	0
7	-0.894286	0.286157	-1.131892	-0.271526	0.609597	0.721818	0.370145	0.851084	-0.392048	-0.41043	-0.705117	-1.104452	-0.286254	0.074354	-0.328783	-0.210077	-0.499780	1.187649	0.570328	0.052737	-0.073425	-0.268092	-0.204233	0.015918	0.373204	-0.344157	0.011747	0.142404	93.2	0
9	-0.338021	1.119594	0.043666	-0.222187	0.493608	-0.246761	0.651832	0.069538	-0.736727	-0.368846	0.017614	0.836386	0.006843	-0.443231	0.502191	0.739452	-0.540980	0.476673	0.451773	0.203715	-0.246914	-0.633753	-0.120794	-0.38505	-0.069733	0.094188	0.246193	0.083076	3.68	0
10	-0.449838	-1.176339	0.913858	-1.375661	-1.971381	-0.629152	-1.423136	0.048459	-1.720408	0.626659	1.199649	-0.67144	-0.513947	-0.095045	0.203904	0.019675	0.253147	0.854348	-0.221365	-0.387226	-0.009302	0.313894	0.027740	0.500513	0.251367	-0.129478	0.042849	0.016233	7.8	0
10	0.384978	0.1616095	-0.8743	-0.094019	0.924584	0.317027	0.470454	0.538247	-0.558895	0.309754	-0.259116	-0.326143	-0.090047	0.362834	0.928903	-0.129487	-0.809979	0.359985	0.707668	1.259916	0.049923	0.238425	0.099129	0.996710	-0.767315	-0.45499	0.026455	0.042421	9.99	0
10	1.249597	-1.221637	0.383902	-1.234899	-1.485459	-0.75323	-0.689405	-0.227487	-1.094611	1.333793	0.227662	-0.242682	0.205416	-0.317631	0.725675	-0.815612	0.879364	-0.847789	-0.683193	-0.102756	-0.231809	-0.483265	0.084667	0.392809	0.161134	-0.35499	0.026455	0.042421	121.5	0
11	0.693760	0.287721	0.826127	0.225204	-0.178388	0.375437	-0.096717	0.115981	-0.221083	0.602304	-0.773657	0.323872	-0.011076	-0.178485	-0.635564	-0.198925	1.240054	-0.880486	-0.882916	-0.153197	-0.036876	0.074124	-0.071407	0.047438	0.548247	0.104942	0.021491	0.021293	27.5	0
12	-2.791855	-0.327771	1.641792	1.074727	-0.136388	0.807595	-0.422911	-1.807107	0.755719	1.151087	0.844555	0.792944	0.370481	-0.317485	-0.635564	-0.198925	1.240054	-0.880486	-0.882916	-0.153197	-0.036876	0.074124	-0.071407	0.047438	0.548247	0.104942	0.021491	0.021293	58.8	0
12	-0.752417	0.345485	0.293229	-1.468643	-1.158394	-0.07785	-0.008581	0.003003	-0.436167	0.747738	-0.793981	-0.770407	1.047027	-1.066004	1.009335	0.601136	-0.279265	-0.419984	0.325233	0.263458	0.499625	0.336305	-0.250573	-0.003084	-0.039124	-0.087086	-0.180580	0.129341	15.99	0

Figure 2 Data set Extraction from sources

Machine learning and data mining techniques have been widely used to detect credit card fraud. On the other hand, purchase behavior and fraudulent tactics might change over time. This phenomenon is known as dataset shift [11] or concept drift [12] in the realm of fraud detection. In this study, we present a method for measuring the dataset shift at our face-to-face credit card transactions dataset (cardholder in the shop) day by day. In practice, we compare and contrast the days and assess the classification's usefulness. The more precise the classification, the more varied the buying behavior across time, and vice versa. As a consequence, a distance matrix describing the dataset shift is obtained. After agglomerative clustering of the distance matrix, we notice that the dataset shift pattern matches the calendar events for this time period (holidays, weekends, etc). The credit card fraud detection challenge then incorporates this dataset shift information as a new feature. As a consequence, detection has improved slightly.

Detecting fraudulent transactions in bank accounts is the main issue that our fraud detection model may help us with. The banking sector is being beset by a significant fraud problem. This is a prevalent problem. Internet fraudsters have mastered the skill of hijacking online sessions by obtaining client credentials and employing malware to steal cash from unwary account holders, and no country is fully secure from them. To train the system, you require well-labeled data in supervised learning. It means that some data has already been labeled with the proper response. It's similar to learning that takes place in the presence of a teacher or supervisor. Unexpected data outputs can be predicted using a supervised learning system, which learns from labeled training data. To be more specific, classification algorithms. However, I will conduct data analysis by: processing, solving classification issues using XGBoost, Random forest, KNN, Logistic regression, SVM, and Decision tree, and selecting the most efficient model for the job at hand. Six distinct classification models may be built: Decision Tree, K-Nearest Neighbors (KNN), Logistic Regression, Support Vector Machine (SVM), Random Forest, and XGBoost. Although there are many more models that may be employed, these are the most widely used models for classification problems. All of these models may be built using the approaches provided by the packages. For the XGBoost model, we'll merely utilize the xgboost package.

Credit card fraud detection is the most prevalent, and therefore the most expensive, problem in the modern world. The rise in both the number of online transactions and the number of e-commerce platforms is to blame. Credit card fraud happens when a credit card is stolen and used for an illegal purpose, or when a fraudster uses credit card information for his gain. Credit card problems are becoming increasingly prevalent in our culture. In order to detect fraudulent activities, a credit card fraud detection system has to be implemented. The focus of this study will be on machine learning techniques and their applications. The algorithms used were the random forest algorithm and the Adaboost algorithm. The performance of the two algorithms is assessed using their accuracy, precision, recall, and F1-score. As a starting point, the confusion matrix is used to create the ROC curve. Several algorithms are examined, including the Random Forest and the Adaboost, and the algorithm with the best accuracy, precision, recall, and F1-score is chosen as the best approach for identifying fraud.

3.4 Data Processing Stages

Data Cleaning - One of the most significant tasks in the data cleaning process is to fill in the missing information. There are various techniques to address this problem, such as disregarding the complete tuple, but most of them are likely to skew the data. Filling them was no longer a problem since the source file, which included legitimate transactions, did not contain any entries with missing data. Tuples with no meaning were eliminated from the files since they do not contribute to the production of useful data and do not skew the data. In addition, adjustments such as deleting redundant columns and splitting the date-time column into two were made.

Data Integration - Because the fraudulent and authentic record files were in two different files, the two data sources were combined before they were exposed to further alteration.

Data Transformation - Here, all of the category data was condensed into a numerical representation that could be understood. The transactional dataset includes a variety of data kinds and ranges. As a result, data transformation includes data normalization. Data normalization reduces the numeric range of attribute data to a manageable size.

Data Reduction - Dimension reduction is the approach utilized for this. We must avoid the possibility of learning incorrect data patterns, and the chosen characteristics must remove the fraud domain's irrelevant elements and attributes [10]. PCA stands for principal component analysis, and it is a well-known transform technique. From the standpoint of numerical analysis, this strategy addresses the feature selection problem. By determining the appropriate number of principal components, PCA was able to effectively execute feature selection.

4.0 Model Evaluation

In terms of effectiveness, some algorithms exceed others. However, to create a more appropriate and acceptable model, it is deemed vital, to begin with, a data pre-analysis method. To note, if we had the option to work with the anonymized data, our model would surely have more full knowledge of the credit card fraud issue. Financial parties can't do their analyses using their unique private information for the customers' privacy, so better knowledge of the data factors would be much appreciated and will assist in the process of classifying the various input variables. Exploratory Data Analysis (EDA) is a vital element in the analytics process as it helps us to comprehend our data. It is vital to study data collection before doing formal analysis. Without a competent EDA, no initial models should be produced. This will help us to better evaluate the patterns in the data, reveal outliers or odd events, and discover intriguing links between variables, among other things. In our instance, the EDA will aid us in better defining our data and preparing it for statistical learning approaches such as those outlined in the preview poster.

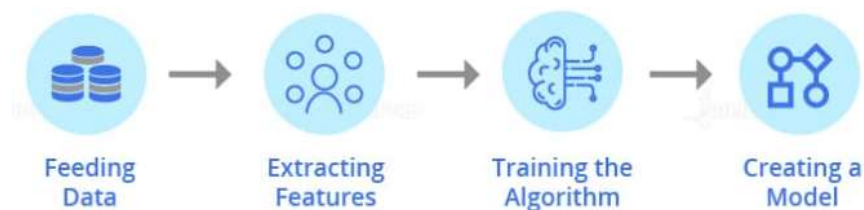


Figure 3 Component of the System and machine Learning Training Model

This paper presents an approach for identifying fraud at the interface. For fraud detection, the recommended solution utilizes unequal, substantially skewed transactions and a modeling methodology. The dataset for credit card fraud detection utilized here is the machine learning Kaggle dataset, which features highly skewed data. The qualities that are assessed are 1 for fraud and 0 for non-fraud. In the banking business, fraud detection analysis was a vital device. Artificial neural networks are presently the least effective technology for identifying credit card fraud. Problems in a way and excessive false positives plague the present technique for identifying fraud. In such circumstances, this research study leverages the collaboration of fully convolutional units to develop a model for spotting credit card fraud that is exceptionally accurate.

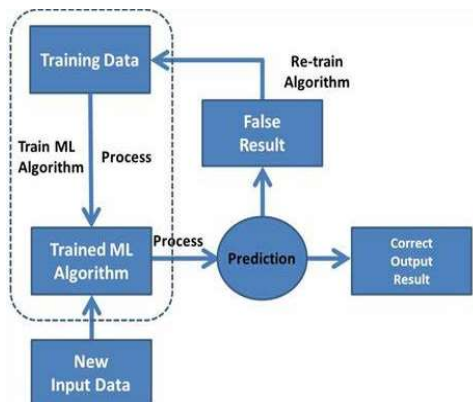


Figure 4: Flow chart of Training Model

Furthermore, for conclusions to be generalizable to the target population, the evaluation prior fraud probability must mirror the population's prior fraud probability (the naturally occurring prior fraud probability).

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

Likelihood
Class Prior Probability
Posterior Probability
Predictor Prior Probability

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Figure 5 Predictions equations

Credit card fraud is common and causes substantial financial losses [1]. The number of online transactions has exploded, with online credit card transactions accounting for a significant share of them. As a result, credit card fraud detection software is highly valued by banks and financial organizations. Fraudulent transactions can take many different shapes and fall into several different categories. This research looks at four forms of fraud that occur in real-world transactions. Each fraud is treated using a series of machine learning models, with the best answer picked after a thorough examination. This evaluation provides a step-by-step guide to selecting a successful algorithm depending on the kind of fraud, as well as a relevant performance measure. Another crucial feature of our project is real-time credit card fraud detection. To do so, we use machine learning models and an API module to do predictive analytics to assess if a transaction is legitimate or not. We also look at a novel method for coping with skewed data distribution. The data we used in our study came from a financial institution, according to a confidential disclosure agreement.

Credit cards are frequently used for online banking. In recent years, there have been several reports of credit card fraud. Fraud done using a credit card is notoriously difficult to spot and prosecute. Machine Learning (ML) is a type of Artificial Intelligence (AI) that is used in research and engineering to solve several problems. In this research, machine learning algorithms are applied to a data set of credit card frauds, and the performance of three machine learning approaches for detecting credit card thefts is compared. The Random Forest machine learning method has the best accuracy when compared to the Decision Tree and XGBOOST approaches.

5.0 Conclusion

To summarize, the work on this project has been a challenging step forward in my grasp of large data and statistical learning. I got the excellent chance to learn about a variety of data approaches, including Exploratory Data Analysis and Principal Components Analysis, as well as create seven classification algorithms from which we identified which one was the most reliable in terms of classification forecast. Furthermore, tree-based algorithms appear to be better suited to our Credit Card Fraud detection problem. Furthermore, more insightful and enhanced datasets will help our research in a deeper investigation of the problem and more thorough characterization of the target demographic. However, due to the wide range of data that may be employed, resolving the class disparity is a difficult problem. One such example is fraud. There is no one-size-fits-all solution to class disparity; instead,

several alternatives should be tested to find the optimal one. Furthermore, because certain algorithms take a long time to execute, it may be essential to make concessions and give up some data, just as people do. Finally, while we concentrated on Supervised Learning, other areas such as Unsupervised Learning or even Deep Learning may be effective in tackling these challenges and identifying fraud.

References

- [1] Artificial Neural Networks- Encyclopedia of Physical Science and Technology https://www.academia.edu/15726358/Artificial_Neural_Networks, Last Visited, 2021.
- [2] Bhanusri A., Valli K., Jyothi P., Sai G., Rohith R., and Subash S., "Credit Card Fraud Detection Using Machine Learning Algorithms," *Journal of Research in Humanities and Social Science*, vol. 8, no. 2, pp. 04-11, 2020.
- [3] Classification Accuracy is Not Enough: More Performance Measures You Can Use, Machine Learning Mastery. <https://machinelearningmastery.com/classification-accuracy-is-not-enough-more-performance-measures-you-can-use>, Last Visited, 2021.
- [4] Credit Card Fraud Detection: Anonymized credit card transactions labeled as fraudulent or genuine, <https://www.kaggle.com/mlgulb/creditcardfraud>, Last Visited, 2021.
- [5] Devi M., Janani B., Gayathri S., and Indira N., "Credit Card Fraud Detection using Random Forest Technique," *International Research Journal of Engineering and Technology*, vol. 06, no. 3, pp. 6662-6666, 2019.
- [6] Entropy: How Decision Trees Make Decisions, <https://towardsdatascience.com/entropy-howdecision-trees-make-decisions-2946b9c18c8>, Last Visited, 2021.
- [7] Jain N., Tiwari N., Dubey S., and Jain S., "A Comparative Analysis of Various Credit Card Fraud Detection Techniques," *International Journal of Recent Technology and Engineering*, vol. 7, no. 5S2, pp. 402-407, 2019.
- [8] Kalra M. and Patni J., "Playing Doom with Deep Reinforcement Learning," *International Journal of Computer Applications*, vol. 1, pp.14-20, 2019.
- [9] Lakshmi S. and Kavila S., "Machine Learning for Credit Card Fraud Detection System," *International Journal of Applied Engineering Research*, vol. 13, no. 24, pp. 16819-16824, 2018.
- [10] Logistic Regression, https://en.wikipedia.org/wiki/Logistic_regression, Last Visited, 2021.
- [11] Mishra P., Patel V., Mittal P., and Patni J., "Algorithm Analysis Tool Based on Execution Time Input Instance-based Runtime Performance Benchmarking," *International Journal of Computer Applications*, pp. 27-30, 2018.
- [12] Mohankumar B. and Karuppasamy K., "Credit Card Fraud Detection Using Random Forest Technique," *International Journal of Innovative Research in Science, Engineering and Technology*, vol. 8, no. 4, pp. 4128-4135, 2019.
- [13] Patni J., Billus S., Billus S., and Singh R., "Feature-Based Opinion Mining and Managed Machine Learning with Sentimental Classification Models," *International Journal of Engineering and Advanced Technology*, vol. 9, no. 2, pp. 3992-3998, 2020.
- [14] Pavithra P. and Babu S., "Data Mining Techniques for Handling Imbalanced Datasets: A Review," *International Journal of Scientific Research and Engineering Development*, vol. 2, no. 3, 2018.
- [15] Racz A., Bajusz D., and Heberger K., "Effect of Dataset Size and Train/Test Split Ratios in QSAR/QSPR Multiclass Classification," *Molecules*, vol. 26, no. 4, 2021.
- [16] Razoogi T., Khurana P., Raahemifar K., and Abhari A., "Credit Card Fraud Detection Using Fuzzy Logic and Neural Networks," in *Proceedings of the 19th Communications and Networking Symposium*, San Diego, pp. 1-5, 2016.
- [17] Shirgave S., Awati C., More R., and Patil S., "A Review on Credit Card Fraud Detection Using Machine Learning," *International Journal of Scientific and Technology Research*, vol. 8, no. 10, pp. 1217-1220, 2019.
- [18] Sorounejad S., Atani Z., and Monadjemi A., "A Survey of Credit Card Fraud Detection Techniques: Data and Technique Oriented Perspective," <https://arxiv.org/abs/1611.06439>, Cornell University, Last Visited, 2021.
- [19] Support-Vector Machine, https://en.wikipedia.org/wiki/Supportvector_machine, Last Visited, 2021.
- [20] Uqaili I. and Ahsan S., "Machine Learning-Based Prediction of Complex Bugs in Source Code," *The International Arab Journal of Information Technology*, vol. 17, no. 1, pp. 26-37, 2020.
- [21] Board of Governors of the Federal Reserve System. *Federal Reserve Payments Study finds U.S. payments fraud a small but growing fraction of overall payments*. 2018. Web. 11 June 2019.

- [22] Roy, A., Sun, J., Mahoney, R., Alonzi, L., Adams, S., and Beling, P. Deep learning detecting fraud in credit card transactions. Systems and Information Engineering Design Symposium (SIEDS) (2018), pp. 129–134.
- [23] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie. Class-balanced loss is based on the effective number of samples. CoRR, abs/1901.05555, 2019.
- [24] Andrea Dal Pozzolo, Olivier Caelen, Reid A. Johnson and Gianluca Bontempi. Calibrating Probability with Undersampling for Unbalanced Classification. In Symposium on Computational Intelligence and Data Mining (CIDM), IEEE, 2015.