

Plastic Money Deceit Detection Using Machine Learning

S¹ Kalaiyarasi, I² Jayalani Shirin S S³ Koushika, R⁴ Abinaya, M⁵ Maheswari

Associate Professor⁵, UG Students¹²³⁴

¹²³⁴⁵Department of Computer Science and Engineering, ¹²³⁴⁵Panimalar Engineering College
¹²³⁴⁵Chennai

Abstract—Fraudulent credit card transactions must be detected by credit card agencies in order to prevent consumers from being invoiced for items which were not purchased. Such troubles may be tackled with facts of technology and its significance, alongside with gadget mastering, cannot be overstated. The purpose of our system is to exemplify the model of its records set using device studying and detecting fraud transaction. The plastic money deceit detection problem consists of designing over credit card transactions with the records which we come across that might result in false transaction. This version is used to perceive whether an basic monetary exchange is deceit or no longer. Our aim is to properly detect illegal transaction by reducing the wrong deceit classifications. In this method, we've got focused on reading and pre-processing records units as properly as the distribution of inconsistency detection methods in addition to a few device gaining knowledge of algorithms additionally, and finding out which set of rules gives the exceptional accuracy.

Keywords—credit card fraud, deceit activities, Random Forest, Adaboost, k-nearest neighbour, logistic regression, xgboost.

I. INTRODUCTION

Credit card fraud [5] means unauthorized operation of an account this is used to make transactions without the actual proprietor of the account or the bank authority's knowledge. We want to take important precautions whilst doing those transactions to avoid these frauds. Additionally, the financial institution government want to apply the modern technologies to are expecting those frauds on the way to alert their clients in advance. credit card count[11] was one of the majority used products which are modelled to make some transactions such as gasoline, groceries, TVs, touring, buying payments and so on because of non-availability of funds at that example. The transaction[9]s made by use of credit card had become popular in both developed as well as undeveloped or uncivilized countries also in order to provide a transaction move at ease of use. Fraud detection approach[4] (for our dataset) is to expect the transactions which can be made through the account holders which are simply executed by way of other people who has access to the account. That is a completely complicated trouble that needs the interest of the account holder in addition to the financial institution authorities so that their other customers need now not suffer from the identical problem. But this hassle has a trouble of class imbalance. The quantity of proper transactions finished with the aid of a purchaser might be a long way higher than the fraud transactions occurred or maybe be 0. Also, the purchaser can do a transaction that deviates from his previous transactions that can be misinterpreted as a fraud transaction. Additionally, [9]the price requests dispatched are checked via automated tools that confirms which request want to be showed. those algorithms test these requests and file suspicious requests to experts who function behind and they in flip look into them through contacting the owners of the accounts whether the transactions are proper or no longer. The fine way[10] to locate incase a transaction is illegal or not we want to locate the purchasing means of the consumer by means of the usage of current facts and use system gaining knowledge to discover whether fit or not.

sorts of deceit:

- on-line and Offline
- Card crimes,
- facts phishing
- utility deceit
- Telecommunication Fraud

II. LITERATURE SURVEY

Fraud is illegal or criminal deception[15] aimed at obtaining financial or non-public gain. it's miles a planned act this is towards the regulation, rule or coverage with an intention to reap unofficial monetary aids. Several literary works touching on detecting fraud transaction on this area have been introduced before and are delivered for consumer usage. A whole analysis done with the beneficial aid of Clifton Phua and his pals have discovered out that strategies hired on this vicinity encompass data mining packages, automatic deceit detection, hostile detection. In every brilliant research, Suman, studies philomath, GJUS&T at Hisar HCE provided various methods that can be used for detecting fraud transaction. A commensable [12]studies area became provided with the avail of many authors, they used some of the outlier detection principles to accurately detect false monetary exchange in the already existing statistical credit card dataset of one sure financial Institutions. Many several supervised and unsupervised machine learning algorithms[6] were used in order to detect the fraud transactions by performing it with a several dataset. The first and foremost aim is to overcome the three level of challenges that is to work with class inequality issue, and to deal with the tagged and untagged data set and to increase the efficiency of our proposed system to be a well developed system which has the capability to handle more number of transaction which leads the system to work proper. The goal of data analytics is to identify hidden patterns and apply them to make informed decisions under various circumstances. [2] With the development of contemporary technology, credit card fraud has increased significantly and has become an easy target for fraud. Credit card fraud costs the economy billions of dollars every year and is a major problem. Economic fraud[1] is a risk that is always increasing and has a variety of effects on the economic sector. The identification of credit card fraud in online transactions was made possible in large part because to data mining. Due to two main factors, detecting credit card fraud, which is a data mining challenge, will get more difficult: first, the profiles of legitimate

and fraudulent behavioural change. In order to build several types of algorithm designs and discover the typical behaviours of fraudulent transactions, this work employed device mastering strategies[8]. Records cleansing, factor construction, extraction of features, and version education are some of the main components. This effort corrects and summaries the material in the next part so that a few typos and unintentional errors won't affect the results. Exceptional case mining is a type of data mining that is mostly used in the internet and financial industries. It helps [7]with identifying items that could be indifferent to the main scheme, i.e., transactions that aren't real. They have estimated the difference between the attribute's determined value and the behaviour of the customer using the cost of all those qualities. An authentic card transaction metrics set can be understood using eccentric techniques like hybrid facts mining/complex society class set of rules, which are primarily based on community search algorithm that consequences developing portrayal of the distraction of one instance from a source institution and have generally proven effective on medium-sized online transactions. The necessary steps have been taken to move forward from a fundamentally nascent element. There have been[13] efforts made to improve the interaction between alerts and comments in the event of a fraudulent transaction. In the event of a fraudulent transaction, the authorized device would be informed and a disclaimer for the ongoing transaction would be provided. One of the strategies that began to throw light on this area was the Artificial Genetic Algorithm, which prevented fraud by taking a novel path. It showed.

III. PROPOSED METHODOLOGY

In this system we're going to design a model which detects all the 5 specified fraud discovery algorithms for a single given data set, and find delicacy rate for all of those to specify the better algorithm which gives further delicacy rate. On addition to that we're also using some of other machine learning algorithms to design a better model with advanced accuracy rate. We were importing several modules similar as numpy, pandas, seaborn, matplotlib for our criteria and numerous further algorithm related modules also. The innovation we bring out on to the already existing systems, is we were using both fraud detection algorithm as well as machine learning algorithm to detect which gives best accuracy rate.

The reason behind choosing the best three algorithm from all fraud detection algorithm is specified as below,

- Use of Logistic Regression here is ,to specify the binary values which determine whether it is true (fraud)or false, by implementing correct insertion of data into its existing logistic function.
- Random forest is deals with construction of more number of decision trees which will give the accurate execution ,moreover it is quite less explainable than decision tree.
- KNN is used ,since it is a simple algorithm which specifies by make use of the majority vote if its nearest neighbours.

The reason behind choosing of boosting algorithm rather than choosing other machine learning algorithms is specified as below,

- Adaboost is used since here all weak classifiers were treated in a well manner ,in our system we were introduced trace of reassigning value in the algorithm and the classification report for both original and modified were displayed as the results obtained.
- XGboost is used because it works well with imbalanced data, here in our system we were using some set of weaker models such as logistic regression ,k nearest neighbour and random forest algorithm by combining these algorithms and treating up well this weaker models will helpful in obtaining our aimed results.

1 SYSTEM ARCHITECTURE

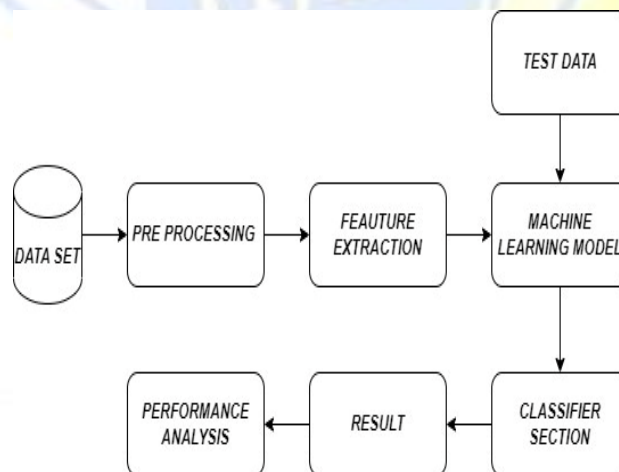


FIGURE -1- System Architecture or workflow

2 ALGORITHM DESCRIPTION

2.1 RANDOM FOREST

RF is an aggregation style and is regarded as association learning for identifying foundational concepts and regression. To understand uneven patterns, deep woods are utilized. Nevertheless, RF uses a mean of the variance of its figure, which may be decreased using this technique. if equal portions of the training sample are examined by deep trees. A random pattern is chosen and replaced with the training set that fits for those samples by the training statistics($p = p1, pn$) with responses($Q = q1,, qn$) and bagging(X times). The following list of guidelines is provided as an algorithm ,Random Forest.

The Random Forest Algorithm's operation is described in the phases that follow:

Step1: Choose Random selection from a specified data collection or training set in step 1.

Step 2: For each training set of data, this algorithm will build a decision tree.

Step 3: Voting will be conducted using an average of the decision tree.

Step 4: Lastly, choose the prediction result that received the most votes as the predicted result.

There are more specified features for using this random forest algorithm, the most specified one is as it deals with working with more number of decision trees which leads in predicting the result accurately.

2.2 LOGISTIC REGRESSION

The straightforward procedure of logistic regression calculates the frequency of an occurrence by estimating the interaction between a single established double variable and unprejudiced variables. The regulating parameter C regulates the switch-off between maintaining the interpretation's simplicity and increasing complexity (overfitting) (underfitting). The interpretation becomes more difficult and the energy of the rule is lost for high values of C , overfitting the data. The Randomised seek $CV()$ procedure is tailored for the most crucial roles in the datasets that are unique, standard, and have the highest degree of uniqueness using the parameter " C ". This model is started and then tailored to the trainings, as mentioned in the technique, as soon as the factor " C " is determined for each dataset. We will employ the same strategy we used for earlier Regression motifs when applying Logistic Regression using Python. The methods for doing so are listed below,

- Data Pre-processing step
- fitting Logistic Regression to the Training set
- predicting the impact of the test
- Test accuracy of the result(Creation of Confusion matrix)
- visualising the test set result.

By using the above specified steps we can completely achieve its efficiency in correctly predicting the accuracy rate for the given specified dataset.

2.3 K NEAREST NEIGHBOUR

The quantity in the information that we require to analyze is referred to as the item or the recognized Variable in the context of supervised learning, which is the literacy that takes into account the volume or result that we desire or expect within the training data (labelled data). Next, we build the interpretation of how the "k- Neighbors Classifier" model operates for the k- Closest Neighbors (KNN) and choose the number "five" as the NN's representative. The number of the "n-neighbors" is called randomly, but it may be determined appreciatively by repeating a number of values, assessed using appropriate methods, and putting the predicted values into the "knn- yhat" variable. The pseudocode for administering the KNN with set of rules from scrape

- cargo the training information.
- put together statistics by means of scaling, missing value remedy, and dimensionality reducing.
- discover the top- quality figure for okay
- anticipate a order value for brand new information
 - o Calculate distance(X, X_i) from $i = 1, 2, \dots, n$.
 - o where $X =$ new data factor, $X_i =$ education records, distance as per your favoured distance metric.
 - o kind those distances in growing order with corresponding train records.
 - o From this taken care of table, pick out the zenith ' k ' rows.
 - o discover the maximum frequent class from those named ' k ' rows

2.4 XGBOOST

The gradient boosting framework is used by the selection-tree-based completely group ML set of rules known as XG boost. Artificial neural networks therefore perform better all the other techniques or frameworks when used with unstructured data that has prediction problems (text, among other things). The XGB Classifier is the XGB boom version for classes. That could fit into our dataset on education. Models are appropriate for use with the academic API and the fit () functionality. The constructor's version can surpass the parameters for instructing the model. We now employ suitable defaults.

Three clear phases make up the boosting ensemble approach:

- The aim of variable y is described by a preliminary model F_0 . A fresh edition h_1 is beneficial to the regression from the preceding phase; This version is likely related to a residual $(y - F_0)$.
- In order to create F_1 , the enhanced model of F_0 , F_0 and h_1 are now combined
- To increase F_1 's overall performance, we should update after F_1 's residuals and build a new version F_2 ; this will be done form iterations, till the residuals were decreased as much as feasible.

2.5 ADABOOST

AdaBoost, sometimes called Adaptive Boosting, is a method for device learning that is applied like an Ensemble learning. Because of this want wood with best 1 split, goal trees with just tier are the most popular estimator used with AdaBoost. These are also known as "choice Stumps" trees. This approach creates a model while assigning equal weights to each piece of input. Then, it gives points that would have been classified incorrectly a higher weight. The new model now puts more importance on all the elements with improved weights. It will continue using these educational models until and until fewer errors are made.

The procedure is as follows,

Step 1: On the basis of the education data alone and just using the weighted samples, a weak classifier (such as a decision stump) is created. The weights of each pattern indicate how important accurate labelling is in this situation. The samples are first delivered with identical weights for the first stump.

Step 2: to develop a selection strategy for each variable and observe how well each one categorizes samples according to its intended use. For illustration, we examine Age, Junk Food Consumption, and Exercise in the diagram below. For each individual stump, we would examine how many samples were rightly or wrongly identified as meritorious or unsuitable.

Step 3: In order to label the erroneously categorized samples effectively inside the next choice stump, greater weight is given to them. Every classifier is also given a weight based on how accurate it is; hence, a classifier with a high accuracy is given a high weight.

Step 4: Continue as in Step 2, repeating the process until all data points have been effectively classified or the most production stage has been achieved.

3 GRAPHICAL REPRESENTATION OF OUR SYSTEM

3.1 Class

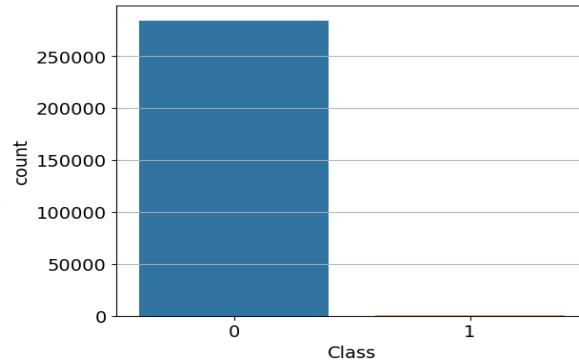


FIGURE 2-Representation of class

3.2 Time

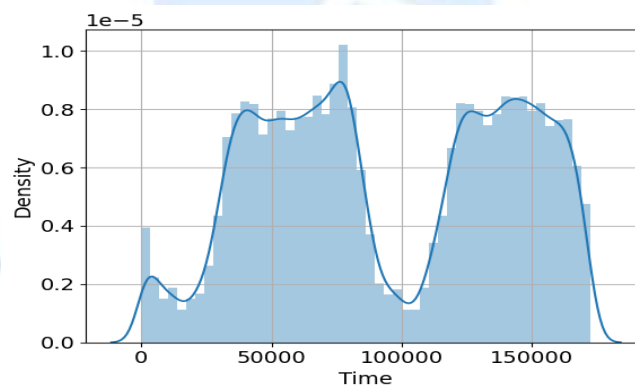


FIGURE 3-Representaion of time

3.3 Amount

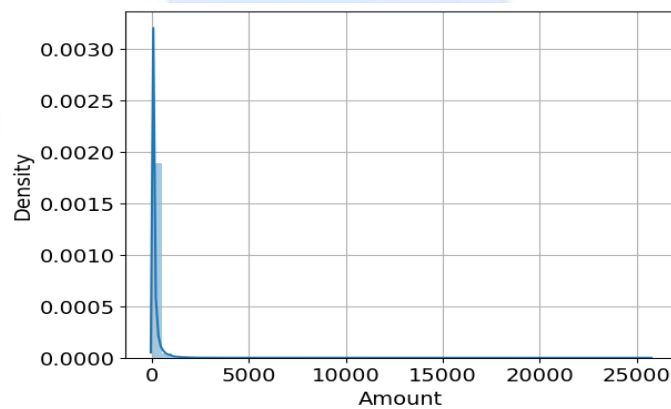


FIGURE 4-Representation of Amount

3.4 Comparison between the amount ,time and class

As of now we have separately seen about the different classification graphs based upon class, time and amount for our specified system which has been treated with some dataset. The above classification graphs have been derived based upon our used monetary exchange dataset. The following classification of graph is the comparison between the class , time and amount one at a time, which results in showing the performance of our proposed system from the already existing system. The graph is specified as below,

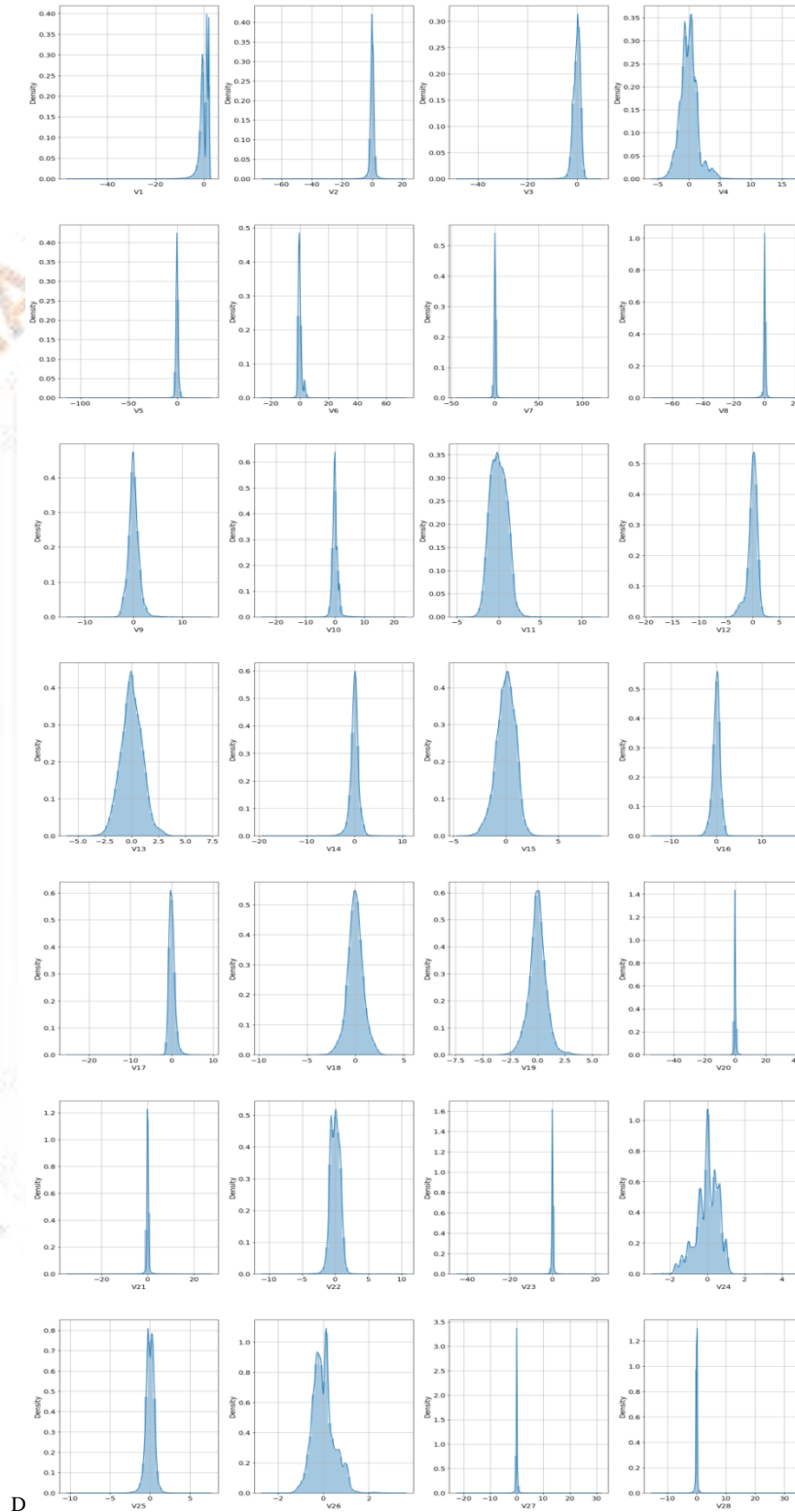


FIGURE 5- Comparison between classes,time,amount

IV. DATA SET

This dataset consists of the transactions that are made using credit cards by EU cardholders in September 2013. This dataset also, contains transactions that come off within a few days, with 492 of the 284,807 transactions being fraudulent. The record is particularly unbalanced, with zero awesomeness (cheat), Out of 172% of entire transactions. It consists of the most useful numeric input variables that can be the final result of the PCA transformation. Annoyingly, due to confidentiality issues, authentic traits and additional genetic information cannot be provided in the record. Features V1, V2, ... V28 are the key components obtained with PCA and are currently being transformed with PCA. Not the most effective skills are "time" and "quantity".

A time characteristic consists of the number of seconds elapsed between each transaction in the data set and the primary transaction. The function "quantity" is the transaction amount. This option is available as an instance-dependent fee-based survey. The awesomeness trait is a response variable that takes a value of 1 if cheating and 0 otherwise. Given the size imbalance, we recommend using the area under the accuracy forgetting curve (AUPRC) to measure accuracy.

V. RESEARCH DETAILS AND DISCUSSION

In the evaluation part of the system, we evaluated the model using the following specific metrics:

1)Accuracy

$$\text{Accuracy} = \frac{(TN+TP)}{(TP+FP+TN+FN)}$$

2)Precision

$$\text{Precision} = \frac{TP}{(TP+FP)}$$

3)Recall

$$\text{Recall} = \frac{TP}{(TP+FN)}$$

4)F1-Score

$$\text{F1 Score} = \frac{2}{(1/\text{precision}) + (1/\text{recall})}$$

Terms involved,

- TN=True Negative
- TP=True Positive
- FN=False Negative
- FP=False Positive

Classification report for random forest algorithm,

	precision	recall	f1-score	support
0	1.00	1.00	1.00	71079
1	0.94	0.78	0.85	123
accuracy			1.00	71202
macro avg	0.97	0.89	0.93	71202
weighted avg	1.00	1.00	1.00	71202

Classification report for Logistic Regression algorithm,

	Precision	recall	f1-score	support
0	1.00	1.00	1.00	71079
1	0.85	0.63	0.72	123
accuracy			1.00	71202
macro avg	0.92	0.81	0.86	71202
weighted avg	1.00	1.00	1.00	71202

Classification report for K-Nearest -Neighbour algorithm,

	Precision	recall	f1-score	support
0	1.00	1.00	1.00	71079
1	0.94	0.78	0.85	123
accuracy			1.00	71202
macro avg	0.97	0.89	0.93	71202
weighted avg	1.00	1.00	1.00	71202

Classification report for XGBoost algorithm,

	Precision	recall	f1-score	support
0	1.00	1.00	1.00	71079
1	0.94	0.79	0.86	123
accuracy			1.00	71202
macro avg	0.97	0.89	0.93	71202
weighted avg	1.00	1.00	1.00	71202

Classification report for ADABOOST algorithm,

Original Classification Report:

	precision	recall	f1-score	support
0	1.00	0.50	0.67	71079
1	0.00	0.45	0.00	123
accuracy			0.50	71202
macro avg	0.50	0.47	0.33	71202
weighted avg	1.00	0.50	0.67	71202

Modified Classification Report:

	precision	recall	f1-score	support
0	1.00	0.50	0.67	71079
1	0.00	0.45	0.00	123
accuracy			0.50	71202
macro avg	0.50	0.47	0.33	71202
weighted avg	1.00	0.50	0.67	71202

THE RESULT ANALYSIS OF OUR SYSTEM IS TABULATED AS BELOW,

TABLE-I

Algorithms Specified in a	F1 Score Detected in Already Existing Systems	F1 Score Which We Detected in Our Proposed System
Random Forest	0.80	0.83
Logistic Regression	0.67	0.72
K Nearest Neighbour	0.75	0.82
ADABOOST	-	0.0030
XGBoost	-	0.86

TABLE-II

ALGORIHMS USED	F1 SCORE
Random Forest	0.828333
Logistic Regression	0.719626
K-Nearest Neighbour	0.817777
ADABOOST	0.003079
XGBoost	0.858407

VI CONCLUSION

Credit card score fraud is the largest scam currently prevalent in the world. This document describes how credit card fraud occurred and investigated these frauds using a dataset consisting of real-world transactions. Using proprietary fraud detection and machine learning algorithms, we found a way to expect fraudulent transactions in our dataset to end up using the XGBoost algorithm, which received excellent F1 score ratings.

References

- [1] Awoyemi, John O., et al. "Credit Card Fraud Detection Using Machine Learning Techniques: A Comparative Analysis." 2017 International Conference on Computing Networking and Informatics (ICCNI), 2017, doi:10.1109/iccni.2017.8123782
- [2] Mohammed, Emad, and Behrouz Far. "Supervised Machine Learning Algorithms for Credit Card Fraudulent Transaction Detection: A Comparative Study." IEEE Annals of the History of Computing, IEEE, 1 July 2018, doi.ieee-computersociety.org/10.1109/IRI.2018.00025.
- [3] F. H. Glancy and S. B. Yadav, "A computational model for financial reporting fraud detection", *Decision Support Syst.*, vol. 50, no. 3, pp. 595-601, 2018.
- [4] Lakshmi S V S S1 ,Selvani, "Machine Learning For Credit Card Fraud Detection System" Deepthi Kavila2 International Journal of Applied Engineering Research ISSN 0973-4562 Volume 13, Number 24 (2018) pp. 16819-16824R.
- [5] Anuruddha Thennakoon, Chee Bhagyani, Sasitha Premadasa, Shalitha Mihiranga, Nuwan Kuruwitaarachchi, "Real-time Credit Card Fraud Detection Using Machine Learning" Publisher: IEEE, 2019 9th International Conference on Cloud Computing, Data Science & Engineering (Confluence), **Date of Conference:** 10-11 January 2019, **Date Added to IEEE Xplore:** 29 July 2019, **ISBN Information:**, **INSPEC**, **Accession Number:** 18868933, **DOI:** 10.1109/CONFLUENCE.2019.8776942.
- [6] Vaishnavi Nath Dornadulaa* , Geetha Sa, Credit Card Fraud Detection using Machine Learning Algorithms," International Conference on recent Trends on Advanced Computing", 2019, ICRTAC 2019 , Vellore Institute of Technology, Chennai-600127, India.
- [7] S P Maniraj , Aditya Saini , Shadab Ahmed , Swarna Deep Sarkar, "Credit Card Fraud Detection using Machine Learning and Data Science", International Journal of Engineering and Technical Research 08(09) September 2019 , DOI:10.17577/IJERTV8IS090031_.
- [8] Yiheng Wei, Yu Qi, Qianyu Ma, Chengyang Shen, Chen Fang," Fraud Detection by Machine Learning", 2020 2nd International Conference on Machine Learning, Big Data and Business Intelligence (MLBDBI) | 978-1-7281-9638-1/20/\$31.00 ©2020 IEEE | DOI: 10.1109/MLBDBI51377.2020.00025.
- [9] IEEE Xplore Part Number: CFP20K74-ART; ISBN: 978-1-7281-4876-2 | DEC 2020 Department Of Information Technology Velagapudi Ramakrishna Siddhartha Engineering College Vijayawada, India.
- [10] D. Tanouz, R Raja Subramanian, D. Eswar, G V Parameswara Reddy, A. Ranjith kumar, CH V N M praneeth," Credit Card Fraud Detection Using Machine Learning" , Proceedings of the Fifth International Conference on Intelligent Computing and Control Systems (ICICCS 2021) IEEE Xplore Part Number: CFP21K74-ART; ISBN: 978-0-7381-1327-2, 2021 5th International Conference on Intelligent Computing and Control Systems (ICICCS) | 978-1-6654-1272-8/21/\$31.00 ©2021 IEEE | DOI: 10.1109/ICICCS51141.2021.9432308.
- [11] Pooja Tiwari, Simran Mehta, Nishtha Sakhuja, Jitendra Kumar ,," Credit Card Fraud Detection using Machine Learning: A Study Technical Report", arXiv:2108.10005v1 [cs.AI] 23 Aug 2021.
- [12] 2020 4th International Conference on Intelligent Computing and Control Systems (ICICCS) Date of Conference: 13-15 May 2020 INSPEC, Accession Number: 19711073, DOI: 10.1109/ICICCS48265.2020.9121114 Publisher: IEEE.
- [13] L. Mukhanov, B. Shchukin" Credit Card Fraud Detection System ", V. Filippov Institute of Control Sciences Moscow, MAR 2021 L. Mukhanov Insitute of Electronic Controlling Machines Moscow, Russia B. Shchukin Institute of Control Sciences Moscow, Russia Email: tsh@cyber.mephi.ru.
- [14] Ruttala Sailusha, V. Gnaneswar, R. Ramesh, G. Ramakoteswara Rao, "Credit Card Fraud Detection Using Machine Learning", 2020 4th International Conference on Intelligent Computing and Control Systems (ICICCS). doi:10.1109/iciccs48265.2020.9121114 Credit Card Fraud Detection using Deep and Machine Learning", IEEE Publisher , 2022 International Conference on Applied Artificial Intelligence, and Computing (ICAAIC) DOI: 10.1109/ICAAIC53929.2022.
- [15] Bora mehar, Sri satya teja , Boomireddy munendra, Mr. S. Gunasekaran, "Research Paper on Credit Card Fraud Detection", International Research Journal of Engineering and Technology (IRJET) e-ISSN: 2395-0056 Volume: 09 Issue: 03 | Mar 2022 www.irjet.net p-ISSN: 2395-0072.