

Indonesian Journal of Computer Science

ISSN 2549-7286 (*online*) Jln. Khatib Sulaiman Dalam No. 1, Padang, Indonesia Website: ijcs.stmikindonesia.ac.id | E-mail: ijcs@stmikindonesia.ac.id

Credit Card Fraud Detection using KNN, Random Forest and Logistic Regression Algorithms : A Comparative Analysis

Vaman A. Saeed¹, Adnan M. Abdulazeez²

vaman.saeed@dpu.edu.krd

¹Technical College of Duhok, Information Technology Department, Duhok Polytechnic University, Duhok 42001, Kurdistan Region of Iraq ²Technical College of Engineering-Duhok, Duhok Polytechnic University, Duhok 42001, Kurdistan Region, Iraq

Article Information	Abstract
Submitted : 22 Jan 2024 Reviewed: 24 Jan 2024 Accepted : 10 Feb 2024	Because credit cards are utilized so frequently, fraud appears to be a significant concern in the credit card industry. It is challenging to quantify the effects of misrepresentation. Globally, credit card fraud has cost institutions and consumers billions of dollars. Despite the existence of numerous anti-fraud mechanisms, fraudsters continue to seek out novel methods and strategies to commit fraud. An additional challenge in the estimation of credit card fraud loss is that the magnitude of unreported or undetected forgeries cannot be determined, only losses associated with those frauds that have been detected can be measured. Implementing effective fraud detection algorithms through the utilization of machine-learning techniques is crucial in order to mitigate these losses and provide support to fraud investigators. This paper presents a machine learning-based method for the detection of credit card fraud. Three methodologies are implemented on the raw and pre-processed data. Python is used to implement the work. By comparing the accuracy-based performance evaluations of k-nearest neighbor and logistic regression with Random Forest, it is determined that the former exhibits superior performance in the basis of accuracy which quantifies the proportion of accurate predictions made by the model out of all predictions.
Keywords	
Credit Card, Fraud detection, Random Forest, K-Nearest Neighbor, Logistic regression	

A. Introduction

Since the inception of online trade, fraud has persisted. Online shopping saw a dramatic uptick due to the COVID-19 epidemic, which presented con artists with a fresh chance. A whopping 38% of all reported scams in 2020 were related to internet shopping, up from 24% before to the epidemic. While that number has gone down after the crisis passed, the sector is still reeling from the effects of security breaches; in 2022, online payment fraud cost businesses over \$37 billion. Because of this, experts predict that the market for tools to identify and prevent fraud in online transactions will double from 2021 to 2026, reaching a value of over \$45 billion(Zhang et al., 2023). As can be clearly observed from graph in Figure (1).



Figure 1. Global losses from credit card fraud

The use of credit cards has significantly expanded in recent years all over the world. People increasingly believe in the concept of being cashless and are wholly reliant on online purchases. With the rise of credit card, online shopping has become much more convenient and accessible. Still, fraudulent credit card transactions rake in millions of dollars annually[2]. One simple target is credit card theft. You may remove a substantial sum quickly and safely without the owner knowing. Credit card fraud has become more widespread as cashless transactions have grown in popularity [3]. Because fraudsters are persistent in their efforts to pass off fraudulent transactions as real, detecting fraud is a daunting and tough undertaking [4]. In terms of ease and convenience, payment cards are hard to beat. Credit card fraud is on the rise as a result of more people using their cards for online purchases in particular. Similarly in the business world, the increase brings up financial risk and unpredictability. Problems with consumer data confidentiality make it hard to get actual transaction records, which are necessary for building good prediction models for fraud detection [5].

Modern technology has made great strides in machine learning, which not only replaces human experts but also works on massive datasets that would be inaccessible to them under normal circumstances. It is possible to accomplish fraud tracking in any way; the only limiting factor is the datasets available. Anomalies should always be recognized in supervised training. Over the previous few decades, several supervised approaches have been employed to detect instances of credit card fraud. It appears that very unbalanced databases are the main problem in using ML to detect fraud. Creating a fraud prevention approach that accurately identifies fraudulent behavior while minimizing false positives is a major problem for investigators [6]. In this research, we present a machine learning-based credit card fraud detection method that successfully analyses the results of various ML models, such as Logistic Regression, Random Forest, and K-Nearest Neighbor, on credit card fraud datasets that are extremely skewed. This study applies a suite of machine learning methods to credit card fraud datasets to glean predictive information. Since the class label is readily available and machine learning classification is often thought to be the optimum answer, it is categorized as supervised learning. A few examples of important classification algorithms include K-nearest neighbor (KNN), Random Forest, and Logistic Regression [7].

This study aims to provide light on the merits, benefits, and comparative analysis of three models—Linear Regression, K-Nearest Neighbor, and Random Forest—in detecting credit fraud. An established assessment criterion assessed by accuracy rates is used to verify their prediction robustness and accuracy.

Payment cards provide a straightforward and practical approach to conducting transactions. As the use of payment cards increases, particularly for online purchases, the incidence of fraudulent activities also rises. Financial risk and uncertainty are introduced into the commercial sector due to the increase. However, authentic transaction records that could aid in the construction of efficient predictive models for fraud detection are challenging to acquire, primarily due to concerns regarding the privacy of consumer data [8]. An approach is proposed by [9] to assess the performance of logistic regression, k-nearest neighbour, and naïve Bayes on credit card fraud data that is highly biased. An evaluation is conducted to determine which machine learning model is most effective in addressing each instance of fraud. Three methodologies are implemented on the unprocessed and pre-processed data. Python is used to implement the work. In assessing the effectiveness of the methods, precision, time consumption, and balanced classification rate are considered. Based on the comparative outcomes, logistic regression exhibits superior performance in comparison to naïve Bayes and knearest neighbour methods [10]. [11] applies an assortment of machine learning algorithms, including logistic regression, naïve Bayes, random forest, and ensemble classifiers employing the boosting technique, to an unbalanced dataset. A comprehensive analysis is conducted on the current and proposed models pertaining to the detection of credit card fraud. A comparative study is also undertaken on these methodologies. So, the data are subjected to various classification models, and the performance of each model is assessed using quantitative metrics including accuracy, precision, recall, f1 score, support, and confusion matrix. Our study concludes with an explanation of the optimal classifier, which is one that is trained and evaluated using supervised techniques and yields superior results. The primary objective of this approach [12] was to examine and ascertain the most effective classification algorithm for credit card fraud detection using benchmark datasets. Random Forest has been determined to possess the highest degree of accuracy in comparison to alternative classifiers. Since the two datasets utilized in this study are unbalanced, a balanced set is also employed to facilitate a more accurate comparison of the algorithms. The dataset is balanced by employing the Synthetic Minority Oversampling Technique (SMOT). By comparing the outcomes using the programming languages Weka and Python. The results of the investigation demonstrate that the methodology is, in fact, extremely useful for any practical application. This study [13] presents a feedback system-based credit card fraud detection mechanism that is effective and is based on machine learning methodology. The implementation of a feedback approach in the classifier improves

both the detection rate and cost-effectiveness. The efficacy of various classifier strategies—including random forest, tree classifiers, artificial neural networks, support vector machine, Nave Bayes, logistic regression, and gradient boosting classifiers—was subsequently evaluated on a credit card fraud data set that was slightly skewed. Historically, the efficacy of methods has been assessed solely based on the performance evaluation metrics for various classifiers: precision, recall, F1-score, accuracy, and FPR percentage. A machine learning model is trained using a variety of algorithms, including logistic regression, random forest, support vector machine (SVM), and neural networks, on the basis of the provided dataset [14]. By employing a comparative analysis of the F_1 score, they successfully forecasted the algorithm that would most effectively fulfil their objective for the identical. With a F_1 score of 0.91, the study concluded that Artificial Neural Network (ANN) performed the best. The evaluation of the effectiveness of the employed techniques (KNN, Naîve Bayes, Logistic Regression, Chebyshev Functional Link Artificial Neural Network (CFLANN), Multi-Layer Perceptron, and Decision Trees) is conducted using a range of accuracy metrics [15].

The papers mentioned earlier offer a comprehensive compilation of a wide range of machine-learning models and methodologies utilized in the domain of credit card fraud detection. They demonstrate the progression of techniques and the continuous endeavours to enhance the accuracy of predictions in the face of fraudulent transactions.

B. Research Methodology

As a result of the rapid evolution of fraud patterns, a proactive approach to fraud detection must be evaluated. In this regard, machine learning has experienced significant growth in recent years. Consequently, the integration of machine-learning techniques into fraud detection algorithms is crucial for minimizing such losses and aiding fraud investigators. Pre-processing techniques and a range of models, such as Logistic Regression, K-Nearest Neighbors, and Random Forest, have been utilized to predict instances of fraud [16]. The data flow of our proposed methodology is detailed in Figure 1.

In order to retrieve, manipulate, model, and evaluate data, the study utilized an extensive collection of tools and libraries from the Python ecosystem. The data collection portion in our proposed method "Online payments big dataset for fraud detection" takes advantage of the extensive collection of datasets available on 'Kaggle' that users can examine and analyze [17]. The datasets are



Figure 2. Data flow of Proposed Model

utilized for the purposes of modeling, testing, and debugging [18]. The credit card dataset is divided into two parts: a training set and a testing set, with a total of 1048,567 sets. We selected the proportions of 80% and 20%, respectively, for this study. In order to facilitate data pre-processing, manipulation, and model development, essential libraries such as 'pandas' played a crucial role in organizing, cleansing, and converting the dataset into a format that was appropriate for modeling. The scikit-learn library offered a wide range of machine learning algorithms, including Random Forest, Linear Regression, and KNN, which could be utilized to implement and assess models. A cohesive framework was established by combining these tools to facilitate data management, model construction, and performance assessment.

C. Pre-Processing

Preprocessing categorical data is accomplished with the OneHotEncoder in machine learning. Categorical data comprises designations or categories for which a numeric representation is not inherent. A significant number of machine learning algorithms, particularly those that perform numerical computations, demand that input features be represented numerically. One-hot encoding is a method for representing categorical variables in a manner that facilitates the effective discovery of patterns and relationships in the data by machine learning models. A multitude of machine learning

algorithms and libraries, including those found in scikit-learn, are specifically engineered to operate without difficulty when confronted with numerical data. Encoding categorical data one-hot ensures that the algorithm correctly interprets the data and promotes compatibility.

1. Random Forest

Random Forest is a classification algorithm comprised of a collection of tree-structured classifiers that each apply a unit vote at input x for the most popular class. The independent random vectors in each tree are identically distributed. The training test is used to generate a random vector that is distinct from the preceding random vectors of the same distribution. An upper bound is extracted for Random Forests in order to calculate the generalization error in terms of the interdependence of individual classifiers and the precision of the random vectors [19]. The flowchart of Random Forest is depicted in Figure 2. The fundamental procedures of the Random Forest consist of selecting a random sample from the dataset and constructing a prediction tree from each tree. Vote on the final prediction of the predication tree and select the one that receives the most votes [8].

prediction Y^{for} a given input X in Random Forest is obtained by averaging predictions of individual trees T:

$$\widehat{Y} = \frac{1}{N} \sum_{i=1}^{N} T_i(X) \tag{1}$$

where N is the number of trees in the forest. The hyperparameters of the Random Forest model, such as the number of trees, maximum depth, and minimum samples per leaf, are optimized to minimize the mean squared error (MSE) or maximize the coefficient of determination (accuracy).



Figure 3. Data Random Forest general example [19]

2. The K-Nearest Neighbour (KNN)

For discrete and continuous label data problems, the KNN algorithm is a type of supervised machine learning algorithm utilized to designate a class to a new data point. In order to predict the label by calculating the similarities between the input data and the training instance, KNN retains the training data [20]. KNN consists of two primary steps: determining the adjacent neighbor and computing the gaps. The object that receives the greatest number of ballots from any entity in regard to their class is referred to as a prediction [19].

3. Logistic Regression:

is a statistical method employed to define the correlation between one or more independent variables and a binary dependent variable using the formula: dependent variable using the formula:

$$l = \log p \left(\frac{p_y}{1 - p_y}\right) \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$
(2)

where

l: log-odds, p: is the base of the logarithm

 β n : are parameters of the model

py : is the probability of the event.

The term "target variable" is used to refer to a dependent variable in machine learning. Predictory variables or features are another name for independent variables.

D. Results and Discussion

The models' performance was assessed utilizing the Kaggle Fraud Labels dataset, which comprised the time period from November 7th, 2021 to November 31st, 2022. Three distinct models—Logistic regression, K-Nearest Neighbor, and Random Forest—were implemented. Every model underwent extensive training and testing, with accuracy being the primary focus of evaluation metrics.

When assessing the performance of machine learning models, accuracy is a frequently employed metric in the testing process. The precision formula is as follows:

$$Acc = \frac{Nomber \ of \ correct \ predictions}{Total \ nomber \ of \ predictions} * 100$$
(3)

[19]. Let's break down the components of this formula:

Number of Correct Predictions: This represents the number of occurrences in which the outcome was accurately predicted by the machine learning model. This would indicate that the model's prediction for a classification assignment corresponds to the true class or label of the input data.

Total Number of Predictions: This is the sum of every prediction that the model generates, irrespective of their accuracy. The accuracy metric offers a direct and uncomplicated means of comprehending the comprehensive correctness of a model with respect to all classes or outcomes. The value is denoted as a percentage, which spans from 0% (indicating no accurate predictions) to 100% (indicating all predictions are accurate).

Nevertheless, it is critical to acknowledge that accuracy might not consistently be the most appropriate metric, particularly in circumstances involving an imbalance of classes. An instance where one class is significantly more prevalent than others could result in a model with a high accuracy that predicts the majority class for every occurrence, despite failing to generate meaningful predictions for the minority classes. Other evaluation metrics, such as precision, recall, F1 score, or area under the ROC curve, may be more informative in such situations.

In essence, accuracy signifies the comprehensive truthfulness of a machine learning model and is computed through the division of the count of accurate predictions by the overall count of predictions.

Models	Accuracy
Random Forest	0.9997
K-Nearest Neighbor	0.9993
Logistic regression	0.9991

Table 1. Accuracy Compression for Proposed Models

The Random Forest algorithm is a robust ensemble learning technique that aggregates the forecasts generated by numerous decision trees. The Random Forest model demonstrates a notable capacity to differentiate between authentic and fraudulent credit card transactions in the provided dataset, as evidenced by its exceptionally high accuracy rate of 99.97%.

The K-Nearest Neighbor algorithm is a straightforward yet efficient method for classifying instances according to the majority class of their k-nearest neighbors. Strong performance is exhibited by the K-Nearest Neighbor model, which accurately classifies transactions and recognizes patterns with a 99.93% success rate. Logistic Regression is a frequently employed linear model in the domain of binary classification. Based on its accuracy of 99.91%, it can be concluded that the Logistic Regression model demonstrates remarkable proficiency in differentiating authentic credit card transactions from fraudulent ones.

The accuracy of all three models—Random Forest, K-Nearest Neighbor, and Logistic Regression—in the proposed method for detecting credit card fraud is exceptionally high. The preponderance of instances for which the models generate accurate predictions are predicted by these encouraging results. On the contrary, it is imperative to take into account additional metrics including precision, recall, and F1 score, and conceivably delve deeper into the dataset, in order to ascertain that the models can effectively detect and manage fraudulent transactions, the minority class, in addition to relying solely on the prevalence of legitimate transactions (the majority class). Further evaluation should be given to the practical implementation of these models, taking into account variables such as interpretability and computational efficiency.

E. Conclusion

Credit card fraud undoubtedly represents an instance of fraudulent deceit. The task of fraud identification initially appears to be a complex and skill-intensive challenge, but that perception changes when machine learning algorithms are introduced. One significant limitation of each technique is that their efficacy in different environments cannot be guaranteed. They yield improved outcomes exclusively with a specific category of datasets while producing subpar or unsatisfactory results with all others. Although some methods, such as random forest, have high detection rates and provide accurate results, they are prohibitively expensive to train. While some, such as K-Nearest Neighbor, perform exceptionally well with tiny data sets, they are not scalable with regard to large datasets. However, with unsampled raw data, certain methods such as logistic regression produce more precise results.

In the proposed methodology for detecting credit card fraud, the Random Forest, K-Nearest Neighbor, and Logistic Regression models have exhibited remarkable accuracy.

Random Forest Model, with its utmost accuracy, was determined to be superior to K-Nearest Neighbor and logistic regression methods through a comparison of all three approaches. Potential future enhancements to performance may involve the construction of a fraud detection model that integrates various deep learning techniques with machine learning algorithms.

F. References

- [1] D. Zhang *et al.*, "Understanding fraudulent returns and mitigation strategies in multichannel retailing," *Journal of retailing and consumer services*, vol. 70, p. 103145, 2023.
- [2] J. O. Awoyemi, A. O. Adetunmbi, and S. A. Oluwadare, "Credit card fraud detection using machine learning techniques: A comparative analysis," in *2017 international conference on computing networking and informatics (ICCNI)*, IEEE, 2017, pp. 1–9.
- [3] R. P. S. Rani, A. Durgabhavani, R. R. I. Malar, and K. Singh, "SMOTE: Credit Card Fraud Detection Using Supervised Machine Learning Methods".
- [4] V. N. Dornadula and S. Geetha, "Credit card fraud detection using machine learning algorithms," *Procedia Comput Sci*, vol. 165, pp. 631–641, 2019.
- [5] X. Niu, L. Wang, and X. Yang, "A comparison study of credit card fraud detection: Supervised versus unsupervised," *arXiv preprint arXiv:1904.10604*, 2019.
- [6] N. K. Trivedi, S. Simaiya, U. K. Lilhore, and S. K. Sharma, "An efficient credit card fraud detection model based on machine learning methods," *International Journal of Advanced Science and Technology*, vol. 29, no. 5, pp. 3414–3424, 2020.
- [7] M. H. Aung, P. T. Seluka, J. T. R. Fuata, M. J. Tikoisuva, M. S. Cabealawa, and R. Nand, "Random Forest Classifier for Detecting Credit Card Fraud based on Performance Metrics," in 2020 IEEE Asia-Pacific Conference on Computer Science and Data Engineering (CSDE), IEEE, 2020, pp. 1–6.
- [8] Y. Sahin and E. Duman, "Detecting credit card fraud by ANN and logistic regression," in *2011 international symposium on innovations in intelligent systems and applications*, IEEE, 2011, pp. 315–319.
- [9] M. U. Safa and R. M. Ganga, "Credit Card Fraud Detection Using Machine Learning," *International Journal of Research in Engineering, Science and Management*, vol. 2, no. 11, pp. 372–374, 2019.
- [10] V. N. Dornadula and S. Geetha, "Credit card fraud detection using machine learning algorithms," *Procedia Comput Sci*, vol. 165, pp. 631–641, 2019.
- [11] A. Bhanusri, K. R. S. Valli, P. Jyothi, G. V. Sai, and R. Rohith, "Credit card fraud detection using Machine learning algorithms," *Journal of Research in Humanities and Social Science*, vol. 8, no. 2, pp. 4–11, 2020.
- [12] M. H. Aung, P. T. Seluka, J. T. R. Fuata, M. J. Tikoisuva, M. S. Cabealawa, and R. Nand, "Random Forest Classifier for Detecting Credit Card Fraud based on Performance Metrics," in 2020 IEEE Asia-Pacific Conference on Computer Science and Data Engineering (CSDE), IEEE, 2020, pp. 1–6.
- [13] N. K. Trivedi, S. Simaiya, U. K. Lilhore, and S. K. Sharma, "An efficient credit card fraud detection model based on machine learning methods,"

International Journal of Advanced Science and Technology, vol. 29, no. 5, pp. 3414–3424, 2020.

- [14] P. Sharma, S. Banerjee, D. Tiwari, and J. C. Patni, "Machine learning model for credit card fraud detection-a comparative analysis.," *Int. Arab J. Inf. Technol.*, vol. 18, no. 6, pp. 789–796, 2021.
- [15] D. Dighe, S. Patil, and S. Kokate, "Detection of credit card fraud transactions using machine learning algorithms and neural networks: A comparative study," in 2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBEA), IEEE, 2018, pp. 1–6.
- [16] S. V. Suryanarayana, G. N. Balaji, and G. V. Rao, "Machine learning approaches for credit card fraud detection," *Int. J. Eng. Technol*, vol. 7, no. 2, pp. 917–920, 2018.
- [17] "fraud dataset," https://www.kaggle.com/code/tarunjonwal/onlinepayment-fraud-detection.
- [18] Vaman A. Saeed and Adnan M. Abdulazeez2 and Adnan M. Abdulazeez2, "online-payments-fraud-detection-dataset."
- [19] N. M. Abdulkareem and A. M. Abdulazeez, "Science and Business," *International Journal*, vol. 5, no. 2, pp. 128–142, 2021.
- [20] N. M. Abdulkareem, A. M. Abdulazeez, D. Q. Zeebaree, and D. A. Hasan, "COVID-19 world vaccination progress using machine learning classification algorithms," *Qubahan Academic Journal*, vol. 1, no. 2, pp. 100–105, 2021.