

<https://doi.org/10.48047/AFJBS.6.7.2024.2226-2232>



African Journal of Biological Sciences

Journal homepage: <http://www.afjbs.com>



Research Paper

Open Access

## Privacy and Security in Smart Card using AI and ML Techniques

M. Supriya<sup>1</sup> Dr M.Suresh Babu<sup>2</sup> Dr. A. Pranayanath Reddy<sup>3</sup>

<sup>1,2,3</sup> Department of CSE, Teegala Krishna Reddy Engineering College, Hyderabad, Telangana, India.  
supriya.mougiligidda@gmail.com, sureshcse@tkrec.ac.in, a.pranayanath@tkrec.ac.in

### Article History

Volume 6, Issue 7, 2024

Received: 29 Apr 2024

Accepted : 10 JUN 2024

doi:10.48047/AFJBS.6.7.2024.2226-2232

### Abstract

Credit card fraud detection using AI and machine learning techniques is a vital application in the financial industry, aimed at safeguarding both customers and financial institutions from fraudulent activities. Enhancing model performance relies heavily on creating relevant features from raw data. For example, features such as transaction frequency, average transaction amount, and time since the last transaction can be calculated. Anomaly detection methods like Isolation Forests, One-Class SVM, and auto encoders help identify transactions that deviate from a cardholder's normal behavior. This detection task is typically framed as a binary classification problem, classifying transactions as either "fraudulent" or "legitimate." Common algorithms employed include Logistic Regression, Decision Trees, Random Forests, Support Vector Machines, and Neural Networks. While effective fraud detection mechanisms and methods like poly encryption and randomization are commonly used to secure transactions, they are not infallible.

**Index Terms:** Artificial Intelligence, SVM, Random Forest, Logistic Regression, Anomaly detection.

### 1. INTRODUCTION

This paper showcase a novel approach to credit card fraud detection by integrating feature selection with various ML classifiers, including Random Forests, Decision Trees, Logistic Regression, Naive Bayes and Artificial Neural Networks. A genetic algorithm is employed for feature selection, optimizing the model's feature set to enhance accuracy and efficiency. The primary goal is to improve fraud detection in e-commerce and e-payment systems, which are increasingly susceptible to fraudulent activities due to the exponential growth of intelligent devices and seamless connectivity.

Effective fraud detection mechanisms, such as poly encryption and randomization, are commonly used to secure transactions but are not foolproof. This

research leverages machine learning, a subfield of artificial intelligence, to address these challenges.

A critical aspect of this study is the importance of reproducibility. Many published works on fraud detection are challenging to reproduce due to the confidential nature of credit card transaction information. Reproducibility is essential for scientific research, ensuring that findings are reliable and verifiable.

The research explores the application of supervised machine learning algorithms, including Random Forests, Decision Trees, Artificial Neural Networks, Naive Bayes, and Logistic Regression, for credit card fraud detection. These algorithms are trained and tested using a labeled dataset, where each transaction is marked as "fraudulent" or "legitimate."

Techniques such as oversampling, under sampling, or using synthetic data (SMOTE) can help address this imbalance. It's important to understand why the model classifies a transaction as fraudulent. This can help in investigations and decision-making. Techniques like SHAP (SHapley Additive exPlanations) can be used for model interpretability.

Feature Selection with Genetic Algorithm (GA): The paper emphasizes the importance of selecting the right features for credit card fraud detection. Feature selection can lead to more efficient and accurate models. The GA algorithms are a type of maximization technique which are used to predict the most equivalent features for fraud detection.

## 2. ML CLASSIFIERS

Decision trees are a machine learning algorithm used for classification, structuring the decision-making process into a tree-like format based on feature values. Random forests enhance prediction accuracy and mitigate overfitting by combining multiple decision trees through an ensemble learning method. Logistic regression, a straightforward classification algorithm, is widely employed for binary classification problems such as fraud detection. Neural networks, particularly deep learning models, are highly versatile and effective for recognizing complex patterns. Neural networks, especially deep learning models, are highly versatile for complex pattern recognition. They can learn intricate patterns in the data. Naive Bayes is a probabilistic classifier that is particularly efficient and effective when dealing with text classification and similar problems. It can also be useful for fraud detection.

Feature Selection with Genetic Algorithm: To address the challenge of a high-dimensional feature space in credit card fraud datasets, the research employs a feature selection algorithm based on Genetic Algorithm (GA). The GA uses Random Forest as its fitness function. Different parameters and input variables, automatically handle missing values, and resist noise in the data.

Overall, the research aims to develop a solution for identifying credit card fraud with high accuracy. The use of feature selection through Genetic Algorithm and the selection of various machine learning algorithms show a comprehensive approach to tackling the problem of credit card fraud detection, which is crucial in the context of the increasing risk of fraud in online transactions.

Logistic Regression (LR) is a commonly used ML method, primarily employed for binary classification tasks. It is particularly effective when you need to predict one of two possible outcomes (e.g., yes/no, true/false, spam/ham).

Binary Classification: Logistic Regression is ideal for binary classification problems, modeling the probability that an input belongs to a specific class (typically the positive class) within a range of [0, 1].

Logit Transformation and Linear Function: A type of generalized linear model called logistic regression uses a logistic function to transform the input features after applying a linear function to them. The linear combination of features is represented by a probability value between 0 and 1 by this logistic function.

Sigmoid (Logistic) Function: The logistic function is an S-shaped curve, often referred to as the sigmoid function. It has the mathematical form:

$$F = \frac{1}{(1 + e^{-z})}$$

where "z" is the linear combination of input features and coefficients. Logistic Regression produces a probabilistic output. To make binary decisions, a threshold (usually 0.5) is applied to the predicted probabilities.

Coefficient Estimation: In training a logistic regression model, the algorithm estimates coefficients (weights) for each input feature. These coefficients determine the impact of each feature on the predicted probability.

Maximum Likelihood Estimation: Given the model, logistic regression maximizes the test data to determine the best coefficients. Regularization procedures, like L1 (Rope) and L2 (Edge), can be applied to forestall overfitting and further develop the model's speculation execution.

Interpretability: Logistic Regression models are often favored for their interpretability. You can assess the importance of individual features by examining their coefficients. Logistic Regression is widely used in various domains, including healthcare (e.g., disease prediction), marketing (e.g., customer churn prediction), and credit scoring (e.g., assessing credit risk).

While Logistic Regression is ideal for binary classification, we apply to multiclass classification problems using techniques like one-vs-all (OvA). This method is simple, interpretable, and can provide a good baseline model for many classification tasks. Decision Trees (DT) and Random Forest (RF) are both machine learning methods used for regression and classification tasks.

### 2.1 Decision Trees (DT)

A Decision Tree consists of various types of nodes, including:

Root Node: Primitive node where the decision-making process begins.

Decision Node: These nodes represent points in the tree where a choice or decision is made based on a specific feature's value.

Leaf Node: Leaf nodes are the terminal points in a decision tree, representing the final decision or outcome. The goal is to create a tree that can accurately predict or classify new, unseen data. Decision Trees are renowned for their simplicity and interpretability.

Random Forest (RF): Random Forest is an ensemble learning method that enhances prediction accuracy and reduces overfitting by leveraging multiple Decision Trees. Decision Trees are trained on various subsets of the data (bootstrapped samples) and a random subset of features in a Random Forest, which introduces diversity into the trees. For predictions, each ensemble tree makes a prediction, and the majority vote (for classification) or averaging (for regression) is used to make the final decision.

The Random Forest approach aims to reduce variance and increase the robustness of predictions compared to a single Decision Tree. It is particularly effective with high-dimensional data and complex classification or regression problems. Random Forests are known for their excellent generalization performance and resistance to overfitting. While the text mentions a "mathematical definition of the RF," it does not provide the actual mathematical equation.

**First Study:**

Algorithms: (LR), (DT), (SVM), (RF).  
 Dataset: Highly imbalanced dataset of European cardholders from 2013.  
 Evaluation Metric: Classification Accuracy.  
 Results: RF achieved the highest accuracy at 98.60%

Table 1. Results of First Study

Classification	Accuracy %
RF	98.60
LR	97.70
SVM	97.50
DT	95.50

Suggested that advanced pre-processing techniques could further improve classifier performance.

**Second Study (Varmedja et al.):**

Algorithms: RF, NB and MLP.  
 Dataset: Dataset dealing with class imbalance.  
 Pre-processing: Used Synthetic Minority Oversampling Technique (SMOTE).  
 Results: RF mainly with deception detection accuracy of 99.96%.

Table 2. Results of Second Study

Classification	Accuracy %
RF	99.96
NB	99.23
MLP	99.93

This study suggests further investigation into feature selection methods to improve the performance of various ML algorithms.

**Third Study**

Feature selection is essential for enhancing the accuracy, efficiency, and interpretability of predictive models. The algorithms under consideration are Decision Tree (DT), k-Nearest Neighbor (KNN), Logistic Regression (LR), Random Forest (RF), and Naive Bayes (NB). It is proposed

that more research should be conducted to understand how different feature selection techniques affect these algorithms and improve their performance.

Dataset: Highly imbalanced dataset of European cardholders.

Evaluation Metric: Precision.

Results: KNN achieved the highest precision at 91.11%,

Table 3. Results of Third Study

Classification	Accuracy %
KNN	91.11
RF	89.77
LR	87.5
DT	85.11
NB	69.52

**Fourth Study (Awoyemi et al.):**

Algorithms: NB, KNN, LR.  
 Dataset: European cardholders' credit card fraud dataset with an imbalanced nature.  
 Pre-processing: Used a hybrid sampling technique.  
 Evaluation Metric: Accuracy.  
 Results: NB achieved the highest accuracy at 97.92%,

Table 4. Results of Fourth Study

Classification	Accuracy %
NB	97.92
KNN	97.69
LR	54.86

Noted that feature selection was not explored in this study.

**Fifth Study:**

Algorithms: Decision Tree (DT), Logistic Regression (LR), Isolation Forest (IF).  
 Dataset: European credit cardholder fraud dataset, dealing with imbalance.  
 Pre-processing: Utilized SMOTE.  
 Evaluation Metric: Accuracy.  
 Results: LR achieved the highest accuracy at 97.18%,

Table 5. Results of Fifth Study

Classification	Accuracy %
LR	97.18
DT	97.08
IF	58.83

**Sixth Study (Manjeevan et al.):**

Algorithms: Genetic Algorithm (GA) combined with Random Forest, GA-ANN, and GA-DT.  
 Evaluation Metric: Accuracy.  
 Results: GA-DT achieved the highest accuracy at 81.97%,

Table 6. Results of Sixth Study

Classification	Accuracy %
GA-DT	81.97
GA- ANN	81.82
GA RF	77.95

These studies demonstrate the importance of choosing appropriate machine learning algorithms and handling class imbalance when developing credit card fraud detection systems. The results vary among the studies, and different algorithms perform better in different contexts.

Additionally, feature selection and advanced pre-processing techniques are suggested as ways to potentially improve model performance in these fraud detection tasks.

**Dataset Content:**

European cardholders over a two-day period in October 2014 collected the dataset comprises transactions done by credit card. It includes a total of 284,827 transactions.

**Class Imbalance:**

A significant challenge in this dataset is class imbalance. Specifically, only 0.172% of the transactions are labeled as fraudulent (class 1), while the vast majority are non-fraudulent (class 0).

**Features:**

The dataset comprises 30 features, denoted as V1, V2, ..., V28, Time, and Amount. Last column in dataset represents the class or type of transaction. Transactions labeled as "1" are considered fraudulent, while those labeled as "0" are non-fraudulent.

**Feature Names Not Provided:**

Features V1 to V28 are intentionally unnamed for privacy, safeguarding techniques and integrity reasons, which is a common practice in credit card fraud datasets to protect sensitive information.

**Handling Class Imbalance:**

SMOTE is a straightforward that creates engineered tests for the minority class by interjecting between existing data of interest, hence adjusting the class dispersion. Several studies with references have used this dataset for fraud detection. Class imbalance is a significant concern in such fraud detection tasks because the number of non-fraudulent transactions far exceeds the number of fraudulent ones. The application of SMOTE is an approach to address this imbalance issue, which aims to improve the in identifying fraudulent transactions.

**Synthetic Minority Oversampling Technique**



Fig 1: SMOTE process [9]

Kasongo[7] - Intrusion Detection Systems:

Implemented a Genetic Algorithm (GA)-based FS to enhance the machine learning models in intrusion detection systems. Results showed that applying Random Forest Under the Curve (AUC) of 0.98, indicating effective intrusion detection.

### 3. FRAUD DETECTION IN E-BANKING

Fraud detection in e-banking, or electronic banking, is crucial for maintaining the security and trust of online financial transactions. With the increasing popularity of online banking, there has also been a rise in various types of fraud, such as phishing, identity theft, account takeover, and credit card fraud. To combat these threats, e-banking institutions employ various fraud detection methods and technologies. Real-time transaction monitoring is essential to identify and flag suspicious activities as they occur. Algorithms and ML models can be used to analyze transaction patterns, detect anomalies, such as unusually large or frequent transactions, transactions from unusual locations, or multiple failed login attempts. Behavioral analytics involves building profiles of account holders based on their typical usage patterns. Any deviations from these patterns can trigger alerts. For eg, if a customer suddenly begins making global transactions when they've never done so before, it could be a sign of fraud. They can identify unusual attitude and flag potentially fraudulent transactions. Tracking the geographic location of a user's device can help identify suspicious transactions originating from unfamiliar locations, especially for international transactions. Effective fraud detection in e-banking requires a combination of technological solutions, user education, and proactive measures. It's an ongoing process that evolves with the changing landscape of cyber threats.

#### 3.1 Multi-Factor Authentication (MFA)

Implementing MFA for login and sensitive transactions adds an additional security. This typically involves something the user knows, something the user has (a token or an intelligent device app), and something the user is (biometric verification). Device fingerprinting assigns a unique identifier to a user's device. If a request or transaction originates from an unfamiliar or suspicious device, it can trigger a review or require additional authentication.



Fig 2. Multi factor authentication Source: www.nist.gov

These examples illustrate the importance of feature selection in various domains and its ability to enhance the performance of ML models. By considering the parameters and maximizing the feature space, models can achieve higher accuracy and effectiveness in solving specific problems such as intrusion detection, disease prediction, and fraud detection.

Genetic Algorithm (GA) feature selection is implemented in the context of the Random Forest

(RF) classifier for solving optimization tasks. It highlights the following key points:

Evolutionary Inspired Algorithms :

EAs are a class of maximization algorithms inspired by the objectives of evolution and natural selection. They maintain a population of potential solutions, evolve these solutions through variations, and use fitness measures to evaluate the quality of each solution.

Population:

In EAs, a population consists of a group of potential solutions, each referred to as an individual.

Fitness:

An individual within the population. The fitness measure assesses how well an individual performs in solving the optimization problem at hand.

Variation:

Individuals in the population evolve through variations, which are inspired by biological gene evolution. Using Random Forest (RF) as the Fitness Method. Random Forest (RF) is selected because it effectively addresses the overfitting issue often associated with regular Decision Trees (DTs). RF performs well even on datasets with class imbalance and eliminates the need for normalization.

Fitness Measure for Feature Selection:

The fitness function evaluates a candidate solution (a feature vector) by assessing its fitness. The fitness is measured based on the accuracy achieved by the specific attribute vector during the testing process of the RF classifier within the Genetic Algorithm (GA).

#### **Algorithm for Implementing RF in the GA:**

Algorithm 1 in the text provides a pseudocode implementation of the fitness function used in the GA.

- Data is divided into training and testing subsets.
- An RF classifier instance is created.
- RF instance is trained.
- Predictions are stored.
- Evaluation is conducted, with accuracy as the main performance metric.

The overall objective of this approach is to use a GA to select the most relevant features for a machine learning model, particularly the RF classifier. This can help improve model performance and optimize feature selection for specific tasks.

Algorithm 2 outlines the pseudo code for the context of feature selection.

Initialization Phase:

In this step, we reduce errors and clean the data with Extract transform load.

For the various components that are used in the computation, variables are defined:

List A: consists of the names of all dataset features.

y variable: Addresses the objective variable.

B Array: A vacant cluster assigned to store the ideal element names.

Variable k: indicates how many times a candidate feature vector must be computed in total. Age of

Introductory Populace (Stage 1): The primitive data of feature names is generated and stored in list A.

Computation of Candidate Feature Vectors (Steps 2 and 3):

Algorithm 2 is computed iteratively to generate candidate feature vectors.

Fitness Calculation (Step 4):

The fitness value, q, is computed. This value feature vector is optimal or not.

Iteration and Evolution (Step 5):

If it is not optimal then (i.e., q is not satisfactory), the algorithm proceeds to the next steps for evolution.

The crossover (k-point crossover, where k=1) is performed.

Mutation is applied.

Fitness is computed again.

Convergence and Termination:

The entire process (Steps 2 to 5) is conducted until the algorithm reaches out.

Algorithm 2 represents a process for selecting an optimal set of features by iteratively generating candidate feature vectors, evaluating their fitness, and evolving the feature set to improve performance. The convergence point is reached when the algorithm no longer improves its accuracy after a set number of iterations (k). This process aims to identify the most relevant features for the specific task and maximize the performance of the machine learning model. The text describes the architecture of a proposed fraud detection methodology, as depicted in Figure 3. This methodology includes several key components and steps. Here's a breakdown of the components and their functions:

Normalize Inputs Block:

The "Normalize Inputs" block is the methodology's first step. The training dataset is normalized using the min-max scaling technique in this block.

Min-Max Scaling:

Min-max scaling is a strategy used to change mathematical highlights to a predefined range (typically somewhere in the range of 0 and 1) in view of the base and greatest upsides of the element. A feature called "f" is normalized using the min-max scaling formula within the range [0, 1], where "min(f)" denotes the feature's lowest value and "max(f)" its highest value.

GA Feature Selection Block:

The "GA Feature Selection" block implements the Genetic Algorithm (GA), using the normalized data obtained from the "Normalize Inputs" block.

Training Block:

The task of training machine learning models falls under the purview of the "Training Block." GA generates a candidate attribute vector known as "vn" at each iteration of the "GA Feature Selection" block. The machine learning models are then trained using this vector.

**Testing Block:**

The "Trained Model" block is replaced with the "Test Data" block. Each machine learning model undergoes examination for each "vn" generated by the GA until a feasible output is obtained. This fraud detection framework encompasses preprocessing (normalization), feature selection using a Genetic Algorithm, training machine learning models, and testing the models. The objective of the framework is to optimize feature selection for improved performance in credit card fraud detection. The normalization step ensures uniformity in input values, while the GA aids in selecting the most pertinent features for the task. The models are trained and tested using the selected feature vectors to achieve desired outcomes. The Results and Discussions section presents the outcomes of experiments conducted in two steps. In the first step, a classification process was executed using different feature vectors denoted as  $F=\{v1, v2, v3, v4, v5\}$ . For each feature vector in F, various machine learning techniques (RF, DT, ANN, NB, and LR) were trained and tested.

**Performance of Different Feature Vectors:**

Using different feature vectors, the performance of the models varied.

For v1, RF demonstrated superior precision, while ANN achieved the highest test accuracy (TAC) of 99.94%. With a precision of 99.93 percent, RF had the best results with v2. Additionally, for v3, RF likewise got the most elevated misrepresentation recognition precision of 99.94%. For v4, DT accomplished an exactness of 99.1% with an accuracy of 81.17%. With RF achieving a precision of 95.34 percent and an accuracy of 99.98%, v5 produced the best results.

**Comparative Analysis:**

In terms of precision and accuracy, v5 performed the best of the various feature vectors. Prominently, the Naive Bayes (NB) technique reliably failed to meet expectations with regards to review, accuracy, and F1-Score across these examinations.

Using a synthetic dataset that was made public, additional experiments were carried out to confirm the effectiveness of the proposed method. User information, transaction details, and the target variable "Is Fraud," which indicates whether a transaction is fraudulent, are among the features in this information set. There were a significant number of legitimate and fewer fraudulent credit card transactions in the dataset. Dataset Features:

The dataset included various features such as user information, transaction details, and other attributes. The features included user information (User, Card), transaction details, merchant information like name, type, pincode, errors, and the target variable "Is Fraud."

**Experimented Methods:**

Experiments considered the following machine learning methods: Decision Tree (DT), Random

Forest (RF), Naive Bayes (NB), Artificial Neural Network (ANN), and Logistic Regression (LR).

These results highlight the effectiveness of the proposed method in credit card fraud detection, even on a synthetic dataset with a substantial class imbalance. The GAFS process, combined with different machine learning models, demonstrated high accuracy and efficiency in identifying fraudulent transactions. The GA-RF model using feature vector v5 achieved an exceptional overall accuracy of 99.98% on the fraud dataset. The GA-DT model using feature vector v1 achieved an impressive accuracy of 99.92%. These results surpassed the performance of existing methods.

**Validation on Synthetic Dataset:**

The proposed framework underwent validation using a synthetic credit card fraud dataset. In these experiments, the GA-DT attained a flawless AUC of 1 and 100% accuracy. Additionally, the GA-ANN exhibited exceptional performance, achieving an AUC of 0.94 and also 100% accuracy.

**4. CONCLUSION**

Inform customers about safe internet based rehearses and draw in with them on the off chance that a thought deceitful exchange happens. This area of monetary security is critical and consistently advancing because of the rising events of monetary misrepresentation in online business and e-installment frameworks. The proposed framework's adaptability and effectiveness in a variety of scenarios could be evaluated by examining it on various datasets in subsequent endeavors. When combined with the Genetic Algorithm for feature selection, the results demonstrate that the proposed method is effective at detecting credit card fraud. The foundation for future research and application in the field of fraud detection is laid by this framework, which demonstrates the capacity to surpass the methods that are currently in use.

**REFERENCES**

- [1] Iwasokun GB, Omomule TG, Akinyede RO. Encryption and tokenization-based system for credit card information security. *Int J Cyber Sec Digital Forensics*. 2018;7(3):283–93.
- [2] Burkov A. *The hundred-page machine learning book*. 2019;1:3–5.
- [3] Maniraj SP, Saini A, Ahmed S, Sarkar D. Credit card fraud detection using machine learning and data science. *Int J Eng Res* 2019; 8(09).
- [4] Dornadula VN, Geetha S. Credit card fraud detection using machine learning algorithms. *Proc Comput Sci*. 2019;165:631–41.
- [5] Thennakoon, Anuruddha, et al. Real-time credit card fraud detection using machine learning. In: 2019 9th international conference on cloud computing, data science & engineering (Confluence). IEEE; 2019.

- [6] Privacy-preserving tax fraud detection in the cloud with realistic data volumes – UGC Sponsored Seminar on MoFS – Department of Commerce, S.V.University – March 2017.
- [7] Sydney Mambwe Kasongo, Yanxia Sun, “A deep learning method with wrapper based feature extraction for wireless intrusion detection system”, *Computers & Security*, Volume 92, May 2020.
- [8] Data Protection and Privacy Preservation using Anonymisation and Pseudonamisation" – IIMT – New Delhi.
- [9] A Comparative study on Privacy Preservation techniques using Fisher Yates, Mondrian and Datafly algorithms - *International Journal of Scientific & Engineering Research* Volume 11, Issue 2, February-2020.
- [10] [online]  
<https://medium.com/@parthdholakiya180/smote-synthetic-minority-over-sampling-technique-4d5a5d69d720>
- [11] Security and Challenges in Privacy Preservation of unstructured data using Pseudonymization and Data masking techniques – *IJSER* - ISSN 2229-5518 – April 2019.
- [12] Privacy preservation and Privacy by Design techniques in Big Data.- *International Journal of Computer Sciences and Engineering* (ISSN: 2347-2693), Vol.7, Issue.4, April 2019.