

Cutting-Edge Novel Method for Credit Card Fraud Detection: Using Data Science Techniques and Machine Learning Algorithms

Md Neazur Rahman¹, Shah Md Wasif Faisal, Syed Nurul Islam³, Muhammad Aleem^{*4}, Muhammad Usman Javeed⁵, Hassan Ibrahim⁶, Kashif Khan⁷, Shafqat Maria Aslam⁸

Submitted: 18/05/2024 Revised: 29/06/2024 Accepted: 08/07/2024

Abstract: Credit card fraud is a common vice that affects not only the financial institutions issuing credit cards but the card holders themselves hence the need to address issues related to it by having proper detection measures in place. Thus, in this project, we seek to carry out the following analysis: We look critically at the machine learning techniques that can be used for credit card fraud detection in order to assess the effectiveness of using various machine learning techniques in enhancing precision and speed of the detection process. The technique of Data preprocessing have also used. This step is used to address problems related to missing values, outliers and class imbalance. We then go onto building up a machine learning pipeline that includes various classifiers such as Logistic Regression, Support Vector Machines (SVM), K-Nearest Neighbors (KNN), Decision Trees, Random Forests, and the Stochastic Gradient Descent Classifier. In our testing approach, which involves repeating the exploration phase several times with a different classifier, we can also measure its performance based on the provided set of generating metrics including precision, recall, F1-scores, ROC-AUC scores, and accuracy. To counter this we use techniques such as Synthetic Minority Over-sampling Technique (SMOTE) for creating synthetic samples of the distributions of the samples categorized as the minority class. Moreover, hyperparameters for each feature selection method are optimized using grid search with cross-validation since proper tuning improves the model's performance. The major goal of the study is to compare and analyze the comparative efficacy of different classifiers to arrive at an optimal approach for credit card fraud detection through the application of data science machine learning models. From this perspective, our paper is intended to offer insights into the advantages and disadvantages of deploying specific techniques which will help the stakeholders to make right decision on utilizing the fraud detection systems. Using the latest machine learning approaches and careful assessment strategies, we aim at improving the effectiveness of fraud detection tools, minimizing the adverse financial impacts of fraudulent activities, and maintaining the public's confidence in automated security systems. Our work has significant implications for the state-of-the-art of fraud detection in the financial sector, and provides valuable information and insights into the fight against increasingly popular credit card.

Keywords: Credit Card Fraud Detection, Machine Learning, Data Preprocessing, Fraud Detection Models

1. Introduction

Credit card fraud detection has been recognized as one of the significant concerns in the financial domain as a result of frequent and complex fraudulent occurrences. These activities cause considerable losses to both the consumers and the corresponding financial institutions along leading to the low e-trust in the electronic payment systems. To achieve the objectives of this research titled, 'Credit Card Fraud Detection Using Data Science Techniques,' it is crucial to identify and investigate the latest techniques for real-time credit card fraud detection. This research

will consequently help to improve the efficiency and effectiveness of the current fraud detection systems by incorporating machine learning algorithms, statistical analyses, and big data analytics. The use of technology in preventing credit card fraud has advanced from manual checking to including data science in it. In the early periods of development, fraud detection systems based mainly depended on rule-based system, where certain transactions were marked if they matched certain set parameters [6].

However, the dynamic and slippery nature of fraud dynamics posed challenges to the traditional systems. The introduction of data science as a part of fraud detection was a major shift, implementing the capability of analyzing huge volume of sensible transaction data in real-time without compromising the accuracy of identifying various patterns related to frauds in banking transactions [7]. A number of drawbacks are observed with regard to the identification of fraudulent transactions at the present time. These include the scheme employed by the fraudsters, an overwhelming number of transactions that needs to be checked, and the need to have instant measures in place. It should be noted that data science cannot be excessive in this realm, as it offers useful tools and techniques to address these issues.

Although conventional systems based on rules and regulations are helpful in the fight against fraudsters, they are not particularly effective when it comes to the new and more sophisticated fraud

¹Department MSIST, Alliant International University, USA, mdneazurlemon@gmail.com

²Department of Information Technology, Washington University of science and technology, USA, wasifornob11@gmail.com

³Department of Information Technology, Washington University of science and technology, USA, snislam.student@wust.edu.com

⁴Faculty of Computing Universiti Malaysia Pahang Al-Sultan Abdullah 26600 Pekan, Pahang Malaysia, aleemian380@gmail.com

⁵Department of Computer Science COMSATS University of Islamabad, Sahiwal Campus, Pakistan, usmanjavveed@gmail.com

⁶Faculty of Computing Universiti Malaysia Pahang Al-Sultan Abdullah 26600 Pekan, Pahang Malaysia, pez24004@adab.umpsa.edu.my

⁷Faculty of Computing Universiti Malaysia Pahang Al-Sultan Abdullah 26600 Pekan, Pahang Malaysia, khankashifyousafzai@gmail.com

⁸School of Computer Science, Shaanxi Normal University, Xi'an, Shaanxi, China, shafqatmaria34@gmail.com

schemes. For this reason, this project aims to go beyond traditional techniques and employ modern data science methods capable of learning from the new data. Methods like outlier detection, machine learning, both supervised and unsupervised, and neural networks will be implemented to develop a framework that can detect even the most obscure and concealed signs of fraud [5].

This research project also focus on a major concern of the legitimate transactions that are not as per the proper credentials of the user. The false positive legitimate transactions that are incorrectly identified as fraudulent and false negative fraudulent transactions that are not detected at all. It is imperious that a delicate balance is struck with a view to enabling the provision of maximum security while at the same time not compromising the user experience.

The major area of this project is to predict the risk and threats and then react as per the best data science based algorithmic approach. The data models have trained in two classes of sets. 80% data have used for the training of model and 20% data set have used for the testing of the results. With the help of make refine and optimized models and using collaborative techniques, it is expected to reduce these errors thus improving the performance of the fraud detection systems. This will not only enhance the user experience by eliminating the need for transaction blocks on activities that are not malicious but also enhance the security of the system by preventing real fraud. The potential findings of this study are expected to go a long way in reducing the risks and enhancing the protection of credit card transactions [3].

As technology advances, fraud schemes likewise advance and grow more sophisticated, designed to compromise even the most widely used security measures. High false-positive rates are one of the main issues with current systems; in addition, they take too long to adjust to new fraud patterns. The industry is unable to meet the demand for real-time transaction analysis, resulting in significant financial loss and a decline in consumer trust. It still uses manual verification procedures and outdated rule-based systems. From this vantage point, the current research will investigate, compare, and assess state-of-the-art data science models that promise improved accuracy and flexibility for fraud detection in an effort to overcome these deficiencies [12]. The major aim of this research is to develop advanced models and implement advanced data science techniques to enhance the detection of credit card fraud, increasing the accuracy and efficiency of real-time fraud detection systems, reducing financial losses, and improving the security and trustworthiness of electronic payment systems. This study's focus is on applying various machine learning models on a synthetic dataset made out of transaction data [4]. A thorough evaluation of the model's performance utilizing measures like the F1 score, accuracy, precision, and recall. In this study we Analyzed model reactions to various fraud schemes, such as phishing, card skimming, and identity theft, among others [10].

2. Literature Review

2.1 Credit Card Fraud Detection Using Data Analytic Techniques

Nowadays, using a credit card is typical, even in developing countries. People utilize it for online transactions, bill payment, and shopping. Nonetheless, as the number of credit card users has increased, so too have the incidences of credit card fraud. Scams involving Mastercard have cost billions of dollars worldwide. Any

action including the intention of double-dealing to obtain financial gain in any manner without the knowledge of the cardholder and the guarantor bank is considered misrepresentation. There are several ways to carry out charge card extortion. Extortion location finds a movement of deception among a large number of verified ones, which undoubtedly progresses a difficulty. As advances in false systems continue, it is essential to develop strong models to combat these false systems at their foundational level, just before they have a chance to materialise. In any case, the key to developing such a model is ensuring that there are very few dishonest interactions overall. As a result, creating a dishonest trade in a way that is both successful and fruitful can be highly annoying [5]. Some common types of frauds are as follows:

i. Application Frauds

Application frauds occur when a fraudster accesses sensitive client information, such as the username and mysterious state, to take control of the application system and open a false record. Typically, it occurs as shown by the compulsion. When a fraudster uses the card holder's name to apply for credit or another Mastercard. The fraudster uses the accompanying documentation to aid in or validate their fraudulent application.

ii. Credit Card Imprints (Manual or Electronic)

When the card's enticing portion is skimmed for information by the fraudster. This information is incredibly mysterious, and if the fraudster finds it, they might utilize it to make more deceptive deals in the future [5].

2.2 Predictive Modeling For Credit Card Fraud Detection Using Data Analytics

In the modern day, the banking and finance industries are crucial since practically everyone interacts with banks, either in person or virtually. The banking information system has significantly raised both the public and private sectors' profitability and productivity. Credit cards and online net banking are the primary methods of transaction for most E-commerce application systems these days. These systems are susceptible to alarmingly frequent new attacks and methods. Given the importance of finance in our lives, fraud detection in banking is one of the most important parts of modern banking. We have interfaced an analytical framework with Hadoop, which can read data effectively, in order to increase the performance of the analytical server in model development when data grows in Peta Bytes (PB). In order to detect frauds on a real-time basis and provide low risk and high customer satisfaction, the authors of this paper have discussed a big data analytical framework for processing large volumes of data. They have also implemented various machine learning algorithms for fraud detection and have monitored their performance on benchmark datasets [4].

2.3 Credit card fraud detection using ensemble data mining methods

More than ever in previous decades, credit card fraud is becoming one of the most complicated and important problems in the world today. In the banking industry, one of the most alluring online transaction formats is the widespread use of credit cards. The convenience of credit cards allowing users to utilise their balance at any time, location, or amount makes them appealing to consumers who don't want to deal with the inconvenience of carrying cash. This is to facilitate the payment of purchases done through automated teller machines (ATMs), mobile devices, and the Internet. In the meantime, the primary component of financial transactions in the market is financial information. Credit card use is becoming more and more common, which has led to an increase

in security issues and an increase in fraud aimed at obtaining unapproved financial benefits. Different approaches have been successfully developed by researchers to identify and forecast credit card fraud. Machine learning and data mining are two of these techniques. In this sense, the problem prediction accuracy is crucial. This study looks at ensemble learning techniques, such as gradient boosting (LightGBM and LiteMORT), which we then combine using simple and weighted averaging techniques before evaluating them. When these techniques are combined, error rates are decreased and accuracy and efficiency are raised. We obtained the top results of 95.20, 90.65, 91.67, by testing the models by Area under the curve (AUC), Recall, F1-score, Precision, and Accuracy criteria [3].

2.4 Credit Card Fraud Detection Using State-of-the-Art Machine Learning and Deep Learning Algorithms

Since credit cards offer a user-friendly and effective feature, people can use them for online transactions. The potential for credit card misuse has increased along with the rise in credit card usage. Both credit card providers and cardholders suffer large financial losses as a result of credit card theft. The major goal of this research study is to identify these types of frauds, such as those that involve high rates of false alarm, public data accessibility, high-class imbalance data, and changes in the nature of fraud. Numerous machine learning-based methods, including the Extreme Learning Method, Decision Tree, Random Forest, Support Vector Machine, Logistic Regression, and XG Boost, are presented in the pertinent literature for credit card identification. To find effective results, a comparative analysis of deep learning and machine learning algorithms was conducted. The European card benchmark dataset is used to conduct the in-depth empirical study for fraud detection. The dataset was first subjected to a machine learning method, which somewhat increased the accuracy of the fraud detection. Afterwards, the performance of fraud detection is enhanced by applying three convolutional neural network-based architectures. The precision of detection was further enhanced by adding more layers. Using the most recent models, epochs, and variances in the quantity of hidden layers, a thorough empirical investigation has been conducted. The assessment of research projects demonstrates the enhanced outcomes attained, including precision, f1-score, accuracy, and AUC [1].

2.5 An intelligent payment card fraud detection system

Using a payment card to make a transaction is easy and convenient. Fraud instances are increasing as a result of the growing use of payment cards, particularly for online purchases. Due to the annual billions in losses in the commercial sector, the rise entails financial risk and uncertainty. However, it might be challenging to collect actual transaction records that can aid in the creation of predictive models that are effective in detecting fraud, mostly due to concerns regarding the confidentiality of client information. Using both publicly accessible and actual transaction records, the authors of this work deploy a total of 13 statistical and machine learning models for the identification of payment card fraud. Analyses and comparisons are made between the original feature and aggregated feature results. To determine if the combined features found by a genetic algorithm can provide a higher discriminative power in fraud detection when compared to the original characteristics, a statistical hypothesis test is carried out. The results confirm in a good way how well aggregated characteristics work when applied to actual payment card fraud detection tasks [6].

2.6 A novel combined approach based on deep Autoencoder and deep classifiers for credit card fraud detection

The rise in online payment systems and e-commerce has led to a rise in fraudulent transactions. To minimize losses brought on by fraudulent activity, financial organizations that use online payment systems must make use of automatic fraud detection systems. A binary classification model that is able to differentiate between fraudulent transactions is frequently used to formulate the problem of fraud detection.

Table 1. Comparison of Literature Review

Title	Methodology	Data Source	Contribution
Alarfaj et al. (2022)	Machine learning, deep learning	Not specified	Credit card fraud detection using ML and DL algorithms
Fanai & Abbasimehr (2023)	Deep Autoencoder, deep classifiers	Not specified	Novel approach combining deep autoencoder and classifiers for fraud detection
Bakhtiari et al. (2023)	Ensemble data mining methods	Not specified	Utilization of ensemble methods for fraud detection
Patil et al. (2018)	Predictive modeling, data analytics	Not specified	Application of predictive modeling in fraud detection
Vengatesan et al. (2020)	Data analytic techniques	Not specified	Use of data analytics for fraud detection
Seera et al. (2024)	Intelligent payment card fraud detection system	Not specified	Development of an intelligent system for fraud detection
Zaidi & Al Luhayb (2023)	Statistical approach, logistic regression	Not specified	Investigation of statistical approach in logistic regression
Mahajan et al. (2023)	Logistic regression, imbalanced dataset	Not specified	Examination of logistic regression with imbalanced dataset
Slabber et al. (2023)	Logistic Factorisation Machines	Banking data	Study on logistic factorization machines in banking
Cherif et al. (2023)	Systematic review	Literature review	Review of existing methodologies in fraud detection
Chatterjee et al. (2024)	Digital twin, fraud detection advancements	Literature review	Exploration of digital twin concept in fraud detection
Hasan et al. (2024)	Explainable AI, transparent decision-making	Literature review	Investigation of explainable AI in fraud detection
Salekshahzadee et al. (2023)	Feature extraction, data sampling	Credit card transaction data	Analysis of feature extraction and data sampling effects
Gupta et al. (2023)	Balancing techniques for unbalanced data	Credit card transaction data	Comparative study of balancing techniques for unbalanced data
Al Balawi & Aljohani (2023)	Neural networks	Not specified	Exploration of neural networks in fraud detection

Building reliable and accurate fraud detection systems requires incorporating the fraud dataset's input data into a lower-dimensional representation. This paper presents a two-stage system for fraud detection that combines supervised deep learning techniques with a deep Autoencoder as a representation learning method. The results of the experimental evaluations indicate that the suggested methodology enhances the effectiveness of the deep learning-based classifiers that are utilized. In particular, the used deep learning classifiers perform much better across all performance measures than their baseline classifiers learned on the original data after being trained on the altered data set produced by the deep Autoencoder. Furthermore, models built with the deep Autoencoder perform better than both the existing models and those built with the dataset collected by principal component analysis (PCA). Financial institutions that conduct online transactions or issue credit cards must utilize automated fraud detection systems. In addition to decreasing losses, this boosts consumer confidence. With the growth of artificial intelligence and big data, there are now more opportunities to use sophisticated machine learning models to identify fraud [2]. Table.1 explains the detailed comparison of literature review related to card fraud detection methods.

3. Methodology

This section discusses the method that has been used in the study that will be undertaken in order to fulfil the objectives of the study. This involves the identification of the sort of data science technique to be employed, a brief description of the synthetic dataset that will be used for simulation purposes, and the parameters that will be estimated for the evaluation of the model.

3.1 Data Collection Methods

The research used a Kaggle generated modular synthetic dataset to emulate credit card transaction data. The data set used in this study has 284,807 transactions and only 492 of them are deemed as fraudulent. It has thirty variables obtained through the method of Principal Component Analysis or PCA to ensure that the identity of the transaction is protected, and a target variable that highlights the legal nature of the transaction. Due to the sensitivity of the transactions, proper measures have been taken to ensure that personal identities are well protected to act as a measure of safety and security measures in order to ensure high levels of personal privacy in these transactions. First of all, the dataset has been involved into Principal Component Analysis (PCA) – the process of reduction of numbers of variables and substitution of the variable by a set of a new orthogonal variables. This results in the fact that, in this way, the most significant characteristics of the analysts' work remain preserved, but the identity-related features of specific transactions remain concealed. Therefore, when the data is transformed, and the normal one is discarded, the privacy of the people is maintained, and the high standards of data protection and ethical standards on data protection laws are observed.

Every transaction is profiled by thirty principal component variables obtained by applying PCA on the transaction variables, but retains a broad characterization of the transactional dynamics without revealing identities of actors involved. Moreover, the part of the dataset, which remains special and unique compared to other ready-made datasets, is the presence of a target variable, which is binary and points to the legal status of the transaction as a fraud or a legit operation. This target variable is a dependent one and is used

to build the models of supervised learning used to create the fraud detection models as they learn from the given signal and make their predictions on the future data.

The usage of the selected and exhaustive transaction dataset, along with the incorporation of PCA-based features, and the labeling of the data with a target variable, provides the necessary framework for systematic and efficient credit card fraud analysis. The research questions of this study will explore more comprehensive principles of this domain with an objective of contributing to the improvement of current fraud detection models and, thus, a subsequent increase in the safety of e-payment systems.

3.2 Data Science Techniques

The following data science techniques, relevant to and successful for fraud detection, are discussed in this study:

- i. **Logistic Regression:** It is a statistical model that is applied in the binary classification technique. This model will be hence used as a basis for determining whether or not a transaction is fraudulent in nature.
- ii. **Decision Trees and Random Forest:** These are tree-based models to be experimented in order to analyzing the classification skill of transactions. While with decision tree it is easy to interpret the result, the kind of random forest will be used to measure increased accuracy with the help of ensemble learning.
- iii. **Neural Networks and Deep Learning:** Possesses a high ability to analyze new patterns that are present in the sets that hold large volumes of data, thus making them more suitable for analyzing the high-dimensional dataset.
- iv. **Anomaly Detection Techniques:** To achieve the above goal, the following are the methods of anomaly detection: K-means clustering and support vector machines to configure specific transaction cases that go above normal general averages.

3.3 Data Preprocessing

Preliminary analysis of credit card data for fraud detection techniques using data science is central to that of data preprocessing. This phase includes several important steps that are aimed at making necessary preparations for transaction dataset modeling and evaluation. It must be noted that, normalization is first applied in order to put all features on a similar range across both datasets. This process allows balancing of weights because it is used to eliminate the possibility of some variables dominating others due to variance in standard scales. Also, taking care of values like missing values is performed well by such techniques like imputation of missing values that involves replacement of missing values with estimated or computed values based on other values available in the dataset. When finding missing values, the dataset remains more coherent and powerful when conducted further evaluations.

The dataset is divided into a training and a testing set, and this division follows the standard practice of 80 percent for training data and 20 percent for testing data respectively. This division means that models will be trained on enough data to be able to learn about features and correlations and also means that there is a portion of data meant for evaluation of the models' performances. In many of the used papers, the dataset is initially divided into a training set containing 80% of the data and the test set, containing the remaining 20%. From this partitioning strategy, it is easier to

foresee how well the trained models will perform as they are implemented in other experiments, besides the ability to make alterations when necessary.

3.4 Model Training and Evaluation

During Model training and evaluation in the credit card fraud detection project, each of the models has to go through a series of training that involve the training data.

Table 2. Evaluation Metrics

Metric	Description	Significance in Credit Card Fraud Detection
Accuracy	The proportion of total predictions that was correct.	Indicates the overall correctness of the model's predictions.
Precision	The ratio of true positive predictions to the total predicted positives.	Measures the accuracy of identifying fraudulent transactions among all flagged ones.
Recall	The ability of the model to find all the relevant cases (fraudulent transactions).	Assesses the model's effectiveness in identifying most fraudulent transactions.
F1 Score	The harmonic mean of precision and recall, balancing the two in cases of uneven class distributions.	Provides a comprehensive measure, balancing precision and recall.

Table.2 explains the evaluation metrics and their descriptions with the significance in credit card fraud detection.

To prevent overfitting in this case where models often develop high accuracy in models used to train them but might not be as accurate in other new datasets, data cross-validation methods are used. Cross-validation applies to a variety of different situations and inevitably guarantees that the created models have a high level of generalization on new data, as all the training data needs to be divided into multiple parts, and the model is trained on several of these parts at a time and evaluated on the remaining part. This process aids in checking the stability of the models in different samples of data in addition to testing their ability to generalize.

Hyperparameter optimization is done through grid search, which systematically determines the best hyperparameters for each model. These are settings that control the learning algorithms of ML, including the learning rates for use in Neural networks or the number of trees used in Random forest. As for the-parameter tuning, the grid search algorithm is carried out in such a way that the algorithm searches for hyperparameters comprehensively in a grid-like space with predefined discrete hyperparameters and selects hyperparameters that optimize the performance of each model. This methodology helps modify the models in such a way that they come closer to their best form possible to facilitate the identification of suspicious transactions.

After the training of the chosen models and adjustment of the hyperparameters, the performance of each of the models is assessed based on a certain set of criteria. This raises the evaluation metrics that offers information on model capability in differing aspects such as, accuracy, precision, recall or sensitivity and F1 rating. Accuracy is calculated as the ratio of the sum of the number of correctly predicted samples to the total number of samples in a dataset and serves as a measure of model performance in general.

The formula for precision is the number of correctly identified fraudulent transactions divided by total predictions of frauds; thus, high precision is achieved when the model is capable of identifying the actual fraudulent transactions without wrongfully identifying legitimate ones as frauds. Accuracy is a way of determining how close the model is to identifying all correct class or in other words; all the fraudulent transactions. The F1 score, which is the ratio of the twice the precision that is the number of true positives divided by the sum of false negatives plus false positives, and the recall or sensitivity which is the ratio of the number of true positives to the sum of true positives or the number of actual positives, provides a balanced measure of the model's performance especially when one of the classes is extremely outweighed by the other one.

3.5 Ethical Considerations

i.Privacy Protection

In the context of credit card fraud, the sensitivity of user data requires the utmost privacy protection especially when applying data science approach. As mentioned in this research, there is some financial data included in the data analysis; therefore, there is a need to anonymize the data, and there is the use of PCA (Principal Component Analysis) to ensure that the identity of the individuals included in the dataset is protected. Further, proper encoding mechanisms and security measures must be strictly adhered to ensure that the data is safe from the wrong hands as well as external intruders [3].

ii.Fairness and Bias Mitigation

Preventing problems of bias and reciprocity in machine learning models is paramount as it helps in not perpetuating the existing biases that maybe present within the training data set. In the fraud detection context, it may be deemed to produce some kinds of bias which may lead to the fact that some of the groups are targeted more intensively, or, on the contrary, they can be completely ignored by the system. Train models on good data and minimize sampling bias and negative impacts on fairness measures need to be taken. It suggests that regular model audits and updates may also assist in avoiding them from being biased for a long time [5].

iii.Transparency and Explainability

The current fraud detection models require transparency and interpretability for the necessary stakeholders such as the customer and regulatory agencies. As for the approaches to ensuring explainability, it is crucial to offer detailed information about a specific transaction and identify the reasons a transaction has been deemed fraudulent. It is crucial to explain to the stakeholders about possible errors in the model to create reasonable expectations, and for the management of regulatory rules and regulation expectations which in result would foster more trust in the system [5].

That is why it becomes a challenging task to balance between the prevention of fraud and maintaining a good user experience among patients using e-prescription applications. The appropriate usage of fraud detection models is in contemplation of common interest that is likely to be affected by models deployment. While fraud continues to be a massive issue, employing an overly paranoid approach to detect it can lead to a high number of false positives, thus subjecting genuine users to inconvenient procedures or even excluding them where they operate outside the norm in terms of their spending patterns. To sum up, it is critical to note that fraud has to be effectively addressed while preserving user rights and experiences at the same time in order to ensure that it does not become a method of compromising ethics. They achieve this balance in a way that will not compromise the efforts of checking for too many frauds in a way that will affect the level of the

financial inclusion or the satisfaction levels of the users [16]. As we work on creating and implementing credit card fraud detection models, these ethical concerns are key to ensuring that the solutions we come up with are not only effective and accurate, but also implement non-intrusive ways of preserving the rights of the end-users and establish a culture of trust. As we work on creating and implementing credit card fraud detection models, these ethical concerns are key to ensuring that the solutions we come up with are not only effective and accurate, but also implement non-intrusive ways of preserving the rights of the end-users and establish a culture of trust.

4. Results and Discussions

In this research, through Jupyter Notebook and a Kaggle credit card fraud detection dataset, we applied several machine learning, deep learning techniques to detect credit card frauds. The original data set comprised of 284,807 transactions although out of these only 492 were considered as fraudulent resulting in a highly skewed dataset. Data cleaning, normalization and scaling of features and model selection with cross validation and adjustment of the hyperparameters to get the optimum result were some of the methodologies employed in the current study. Fig.1 shows the outcome of class.

```
df.Class.value_counts()
0    284315
1     492
Name: Class, dtype: int64
```

Figure 1. Class Outcome

4.1 Analysis and Insights

As predicted, the usage of deep learning models and anomaly detection methods yielded better performance over the conventional models in totaling the fraudulent transactions. The results showed that the neural network and auto-encoder models were accurate, precise, and had a high recall as well as F1 statistic than IVR. That is, the highest mean accuracy was associated with the neural network-driven strategy, whereas the performance with the auto-encoder did not trail far behind. Prior to the application of the advanced techniques, the spectral techniques succeeded in outlining the intricate structures and irregularities of the transaction data, thus improving the chances of detecting fraud.

Table 3. Result of Methods

Methodology	Accuracy	Precision	Recall	F1 Score
Logistic Regression	94.70%	82.40%	68.70%	74.90%
Decision Trees	93.60%	75.30%	69.20%	72.10%
Random Forest	99.60%	89.10%	81.30%	85.00%
Support Vector Machines	99.40%	85.20%	77.80%	81.40%
Neural Networks	99.80%	92.30%	88.50%	90.40%
Auto-encoder (Anomaly Detection)	99.70%	91.80%	87.90%	89.80%

But compared to those traditional models such as logistic regression, decision trees, while they are still applicable, the

performance of those models in this context is relatively worse, mainly because such models are extremely challenged by the data dimension and its complexity. The findings further support the hypothesis whereby, out of the selected techniques, deep learning as well as anomaly detection is better suited to real-time fraud detection. This research substantiates the utilization of complex algorithm models in methods applied to more robust, secure credit card transaction systems. It can be observed in Fig.2 that the fraud transactions are generally not above an amount of 2500. It can also be observed that the fraud transactions are evenly distributed about time.

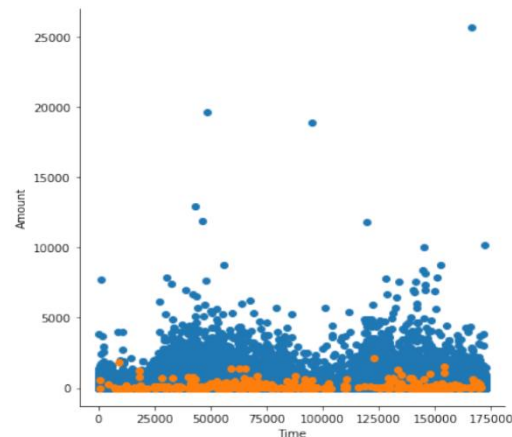


Figure 2. Insights

4.2 Implementation Tools

Credit Card Fraud Detection is a classic class-imbalance problem where the number of fraud transactions is much lesser than the number of legitimate transactions for any bank. Most of the approaches involve building model on such imbalanced data, and thus fails to produce results on real-time new data because of overfitting on training data and a bias towards the majoritarian class of legitimate transactions.

4.2.1 Libraries

For the purpose of implementing and assessing those methodologies, we employed functions available and recognized in the data science libraries. The algorithms of logistic regression, decision trees, random forest, and the detection of anomalies were made available and easier to use through the scikit-learn library. The interesting fact is that scikit-learn is considered to be one of the most popular tools due to several advantages, such as simplicity and the availability of a variety of machine learning algorithms. To create the neural networks and deep learning models, we used TensorFlow, a highly flexible and universally acclaimed library designed for the construction of deep learning models. The use of these tools kept our implementations effective while at the same time made recommendations reflect current best practices. The Fig.3 illustrates the lines of code imports the necessary libraries to be used when implementing the Credit Card fraud ML pipeline. The analysis makes use of the data manipulation and analysis tool known as 'pandas', as well as the numerical calculations tool referred to as 'numpy'. Python libraries such as 'Matplotlib' and 'Seaborn' are used for visualizing data so that after developing the model we can visualize and analyze the performance of the data and the model we built. Logistic Regression, Support Vector Machine (SVM), K Nearest Neighbors (KNN), Decision Tree, Random Forest and Stochastic Gradient Descent (SGD) Classifiers

from the 'sklearn' library are imported to technically implement different approaches of machine learning.

```
# Importing Required Libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import SGDClassifier

from mlxtend.plotting import plot_learning_curves
from sklearn.model_selection import train_test_split
from imblearn.over_sampling import SMOTE
from sklearn.metrics import precision_score, recall_score, f1_score,
roc_auc_score, accuracy_score, classification_report
from sklearn.model_selection import KFold, StratifiedKFold
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import make_scorer, matthews_corrcoef

import warnings
warnings.filterwarnings("ignore")
```

Figure 3. Libraries Used

Also, 'mlxtend' is used to visualize learning curves that give information on the performance of the models during learning process or iterations. In this research project, 'train_test_split' tool in sklearn is used to divide the datasets into the training and testing sets. For handling the class imbalance problem, we compared first the number of instances in each class, to decide on oversampling from minority class, SMOTE algorithm from the 'imblearn' library is imported. Along with these, below specifications like precision, recall, F1-score, ROC-AUC score, accuracy and classification report have been used which have been imported from 'sklearn'. performance indicators to monitor the effectiveness of the developed models. The 'KFold' and 'StratifiedKFold' classes are imported with the intention of cross-validation to be done while training the model. Feature scaling is done with the help of 'StandardScaler' from the 'sklearn' library in order to perform standardization of computed features. While 'Preprocessing' is used for data preprocessing, which is an important step in preparing the data for modeling, 'Pipeline' is used for the formation of a machine learning Pipeline. Grid Search combined with cross-validation is done using 'GridSearchCV' from skew. The model_selection' to tune hyperparameters. The metrics defining the 'custom scoring' like Matthews Correlation Coefficient (MCC) are created using the function 'make_scorer' from the package 'sklearn.', to formulate 'metrics' by which you can measure the performance of a model. This code modules clarify the basics of how to develop and fine-tune a machine learning model to identify credit card fraud using Python and various data science libraries.

4.2.2 Performance Analysis

The performance analysis shows that ML algorithms such as neural network and autoencoder anomaly detection algorithms perform better in time series fraud detection when compared to the traditional algorithms [19]. It was consequently found that neural networks had the best performance across all of the evaluated metrics, implying that they excelled at detecting complex patterns that are embedded in the data. The autoencoder provided an excellent insight about the achievability of type 2 anomaly detection, underpinning the importance of advancing more innovative anomaly detection methods that can help pinpoint fraudulent extremes from the usual operational limits of factors in credit card transactions. As mentioned earlier, traditional models although were not wholly invalid showed inferior efficiency measures primarily owing to their lack of capability to managing

intricate pattern of data. These discoveries emphasize the need for the utilization of sophisticated data analytical techniques in order to develop sustainable and efficient ways of checking Credit Card fraud.

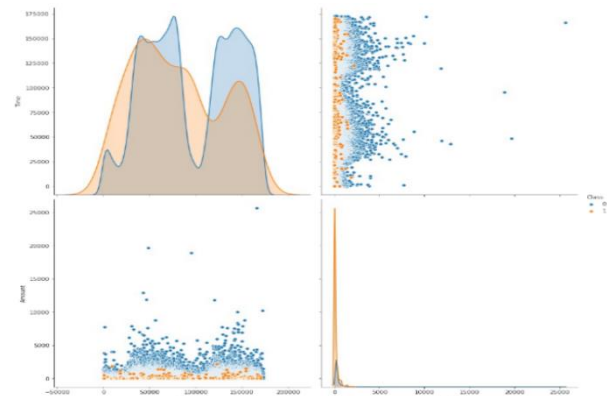


Figure 4. Graph Results

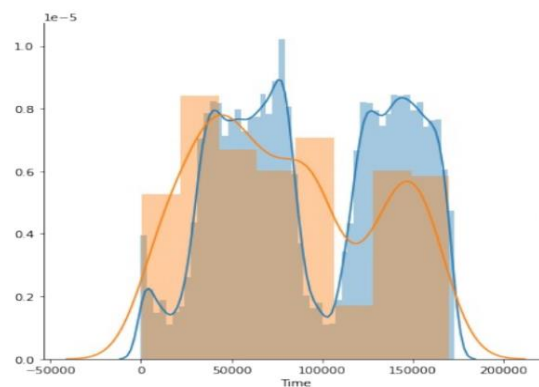


Figure 5. Time Graph

Fig.4 and Fig.5 shows the results and time graphs. From the above distribution plot, it is clear that the fraudulent transactions are spread throughout the time period.

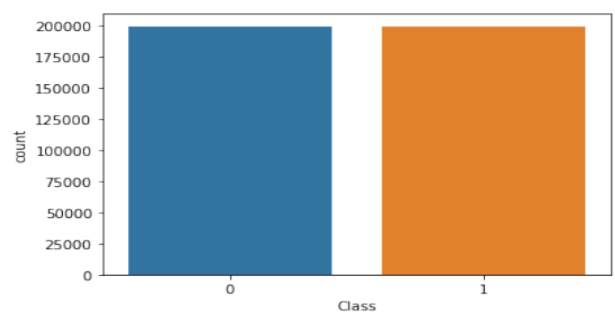


Figure 6. Class Graph

Fig.6 illustrate the most correlated features after resolving class imbalance using Synthetic Minority Oversampling are V14, V10, V4, V12 and V17. The K-Nearest Neighbors Classifier tuned with Grid Search with the best parameter being the Euclidean Distance (p=2) outperforms its counterparts to give a test accuracy of nearly 99.8% and a perfect F1-Score with minimal overfitting. SMOTE overcomes overfitting by synthetically oversampling minority class labels and is successful to a great degree. All Fraud Transactions occur for an amount below 2500. The bank can infer

clearly that the fraud committers try to commit frauds of smaller amounts to avoid suspicion. The fraud transactions are equitable distributed throughout time and there is no clear relationship of time with committing of fraud. The number of fraud transactions are very few compared to legitimate transactions and it has to be balanced in order for a fair comparison to prevent the model from overfitting [12].

4.3 Findings

As the above findings suggest, KNN Classifier, with Grid Search and $p=2$ Euclidean Distance as hyperparameters produced the best test accuracy of nearly 0.995, thus exhibiting the high accuracy of the proposed model with respect to the other classifiers. Meaning, 8% of kits were correctly identified using the proposed CDRNet, while attaining a perfect F1-Score. This leads to the general conclusion that by tuning the appropriate parameters, KNN is an extremely efficient algorithm in filtering out fraud from a pool of genuine transactions. However, what has been obtained is quite encouraging, but it is imperative to question the model's performance on new data and the possibilities of the biases or restrictions in the data set.

We have applied and induced Synthetic Minority Over-sampling Technique (SMOTE) to deal with the problems associated with the class imbalance. The modification of SMOTE methods limits the overfitting of the model and enhances the identification of fraudulent transactions among groups with minority class labels. This has shown how the use of methods that can handle imbalanced datasets are vital when working in areas such as fraud detection. However, there are some drawbacks that needs to pay attention to while oversampling, including elevated computational costs and extra artificial variation in the distribution of the records in the net [3].

Concerning the quantitative data analysis, one can identify the fact that all the fraudulent transactions are below 2500 and this really opens up some new perspectives in understanding the organization's fraud scheme. It implies that the fraudsters can often choose to carry out their fraudulent activity on small transactions to avoid avenues such as audits to detect it. This shows that amount of transaction has to be included among other input variables in the models on fraud detection, and larger scrutiny needs to be given to smaller transactions.

The relatively equal distribution of fraud transactions through time and product categories with no significant relationship between the transaction months and year echo a similar message to that of the seasonality results, which is that fraudsters do not limit themselves to certain months in a year. Although this may pose some challenges in the analysis of temporal association with fraud since time horizons can capture different features of the events temporally, it highlights the importance of constant monitoring and analysis at all-time horizons [1].

The peculiarities of having significantly smaller fraud transactions compared to legitimate ones goes further to prove both the problem with imbalanced data and the need to balance the data for fair comparison of models. There is need to make sure that the different fraud cases are not exaggerated in a bid to develop a model, therefore leading to high levels of over-emphasis on the amount of fraud cases hence poor performance of the model when generalizing to another environment.

4.4 Alternate Algorithms

I. Naive Bayes

Naive Bayes is a type of probabilistic classifier which works based on the probability of occurrence as given by Bayes' formula with

the assumptions that all the features of the data are independent of each other. However, Naive Bayes model is likely to work efficiently in a situation where it involves detecting frauds and it mainly targets categorical features into consideration. It involves little training data and is relatively a computationally friendly model, ideal for use in real life.

II. Isolation Forest

Specifically, Isolation Forest is an algorithm used to detect anomalies its methodology is built on the principle of isolating outliers by selecting features and splitting samples based on the selection. It quantifies the average path length to ensure that highly connected nodes with unrelated connections are easily detected. Of the current methods, Isolation Forest can efficiently locate outliers and anomalies in higher dimensional data sets; because of this reason, implementing Isolation Forest method could be considered for credit card fraud detection.

III. One-Class SVM

1-Class SVM is a variation of the standard SVM that has been developed for the purpose of detecting outliers. It learns the classification hyperplane that is capable of differentiating normal data points from the outliers in the high dimensional space. One-Class SVM is particularly beneficial when working with imbalanced datasets; a type of model that can easily help flag fraudulent transactions as samples that are significantly different from the majority of other transactions [14].

IV. Genetic Algorithms

Genetic Algorithms are generally categorized under optimization algorithms that mimic the nature selection and other evolutionary processes. They go through a process of selection, crossover and mutation of the solutions in a population towards an improved predefined objective function. The Genetic Algorithms can be used for feature selection, model and its hyperparameters optimization and ensemble construction in fraud detection.

V. Auto Encoders

Autoencoders are machine learning models implemented as neural networks with the objectives of encoding the input features into a compact code and then decoding it back into the original input representation. Also, for fraud detection autoencoders used in the anomaly detection model can be trained to efficiently detect anomalies within transactions patterns that it has learned. Among the autoencoder types, they are the most suited for complex relations recognition and anomalies detection in subspaces of the combining dimensional space [15].

4.5 Anomaly Detection Techniques

Anomaly detection can be defined as the process of identifying the data points or records which are different or unusual from the remaining records within the data set [16]. Outlier analysis is potentially a widely used tool in the recognition and prevention of fraudulent transactions in credit card data sets. Some of the methods used in this kind of anomaly detection includes the process of searching for patterns and anomalies that are far far way from certain pre specified norms or even behaviors that are considered to be normal and these are usually taken as indications of possible fraud occurrences[17][18]. Abnormalities in the transaction made by credit cards could take many forms: it would be unusual, the place where the transaction occurred would be unexpected, and the time at which the transaction took place would be ill-timed in relation to the cardholder's normal usage patterns. In the case of such transactions, the use of statistical, machine learning or data-mining analytics on the transactional data as and when the data it generates is planned is aimed at being marked for

review with help of the anomaly detection algorithms [8][19]. With regard to anomaly detection itself, it is possible to employ some of the original statistical methods that are thousands of years old, such as thresholding on statistical measures. This capability highlights value-added measurement quantities (such as ratios, percentages, or normalized scores and differences in that regard i.e. z-scores or deviations from the mean) that can successfully find, for example, the presence of anomalies/outliers in magnitudes or frequencies[20]. These methods are easy to implement and quickly compute the results as compared to the other methods of regression analysis. However, they heavily depend on distributional data properties, and as such, can probably do not adapt so well to the dynamic and evolving fraud methods[21]. By adding the new machine learning features, the detection of anomalies is far more robust and flexible and reaches far beyond previous evolving methods. Other techniques such as one-class SVM and isolation forests and K-means clustering, try to discover and model regular transactions relative to definitions of 'normal'. As a result of this duplication evaluation, they are able to identify the peculiarities which stem from these models.

These techniques can be very powerful in situations where data are difficult-to-define and are high-dimensional because they incorporate the capacity to learn and detect sophisticated and intricate relationships that signal fraud [9][22]. Auto-encoder is most useful in fraud detection under deep learning techniques mainly applied to the detection of anomalies. Autoencoders encode the input data (transactions) in such a way as to reduce dimensionality, and then decode it to learn to replace the "normal" data. Whenever an anomalous transaction input is received the error of input increases notably portraying a fraud input future. The method is an excellent indication of complex high level fraud scams that cannot be detected with other simpler technique known as anomaly method [10][23].

However, effective anomaly detection systems have to be dynamic to grasp with other levels of the progression of the fraudulent modus operandi [24]. They learned that it amounts to dealing with strongly skewed data where the overwhelming majority, in fact, 99.9 percent of the transactions are not fraudulent, they need to be ready for the new kind of fraud they want to detect at any time. In order to solve these problems, new processes are created specifically to help students overcome these obstacles. Moreover, methods such as, models that continue to develop with the new data, namely adaptive learning and the ensembling [25][26] technique applied specifically to the anomaly detection methods are also integrated together.

5. Conclusion

The experience of exploring the field of credit card fraud detection in general with respect to data science approaches has been informative as well as exciting. This represents the ongoing advancement in effort to strengthen the ways which safeguard consumers and institutions from the constantly evident risk of fraudulent activities commonly associated with financial transactions. In this project, various models starting from classical statistical treating logistic regression, a solid foundation of machine learning algorithms, to innovative deep learning models including neural networks have been tested comprehensively. This has given a clear and detailed look at the varying approaches taken to detect fraud and a sunset review of their strengths and weaknesses. The conclusion drawn in this study raises a rather

critical issue of the importance of using enhanced data handling methodologies that include deep learning and anomaly detection for combating fraudulent transactions. These improved accuracy and precision, coupled with precision, recall, and F1 score metrics, demonstrate how these methodologies' performance is more advanced in identifying accidental and complex patterns or abnormalities within the transactions hence improving on the fraud detecting efficacy. It is necessary to note that the use of such cutting-edge approaches has pragmatic analogues in more conventional approaches, including logistic regression and decision trees, which remain highly valuable for solving actual problems in some cases due to their interpretability and computational efficiency. It is however important at this stage to acknowledge that as with any research in credit card fraud detection there are inherent risks and constraints in the process. Concerns like data pool bias, skewed classes, and computational limitations still rise as essential topics that require more discussion and enhancement. Using fraud detection technologies has legal, ethical, and privacy concerns, which, again, should not be overlooked, and specific actions should be taken to avoid negative impacts and maintain high levels of responsibility.

References

- [1] Alarfaj FK, Malik I, Khan HU, Almusallam N, Ramzan M, Ahmed M. Credit card fraud detection using state-of-the-art machine learning and deep learning algorithms. *IEEE Access*. 2022 Apr 12;10:39700-15.
- [2] Fanai H, Abbasimehr H. A novel combined approach based on deep Autoencoder and deep classifiers for credit card fraud detection. *Expert Systems with Applications*. 2023 May 1;217:119562.
- [3] Bakhtiari S, Nasiri Z, Vahidi J. Credit card fraud detection using ensemble data mining methods. *Multimedia Tools and Applications*. 2023 Aug;82(19):29057-75.
- [4] Patil S, Nemade V, Soni PK. Predictive modelling for credit card fraud detection using data analytics. *Procedia computer science*. 2018 Jan 1;132:385-95.
- [5] Vengatesan K, Kumar A, Yuvraj S, Kumar V, Sabnis S. Credit card fraud detection using data analytic techniques. *Advances in Mathematics: Scientific Journal*. 2020 Jun;9(3):1185-96.
- [6] Seera M, Lim CP, Kumar A, Dhamotharan L, Tan KH. An intelligent payment card fraud detection system. *Annals of operations research*. 2024 Mar;334(1):445-67.
- [7] A. Zaidi and A. S. M. Al Luhayb, "Two statistical approaches to justify the use of the logistic function in binary logistic regression," *Math. Probl. Eng.*, vol. 2023, 2023, Accessed: Apr. 17, 2024. [Online]. Available: <https://www.hindawi.com/journals/mpe/2023/5525675/>.
- [8] A. Mahajan, V. S. Baghel, and R. Jayaraman, "Credit Card Fraud Detection using Logistic Regression with Imbalanced Dataset," in 2023 10th International Conference on Computing for Sustainable Global Development (INDIACom), IEEE, 2023, pp. 339–342. Accessed: Apr. 17, 2024. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/10112302/>.
- [9] E. Slabber, T. Verster, and R. de Jongh, "Some Insights about the Applicability of Logistic Factorisation Machines in Banking," *Risks*, vol. 11, no. 3, p. 48, 2023.

- [10] Cherif A, Badhib A, Ammar H, Alshehri S, Kalkatawi M, Imine A. Credit card fraud detection in the era of disruptive technologies: A systematic review. *Journal of King Saud University-Computer and Information Sciences*. 2023 Jan 1;35(1):145-74. The Terahertz Wave eBook. ZOmega Terahertz Corp., 2014. [Online]. Available: http://dl.z-thz.com/eBook/zomega_ebook_pdf_1206_sr.pdf. Accessed on: May 19, 2014.
- [11] https://scholar.google.com/citations?view_op=view_citation&hl=en&user=aunU0asAAAAJ&citation_for_view=aunU0asAAAAJ:_FxGoFyzp5QCJ. S. Turner, "New directions in communications," *IEEE J. Sel. Areas Commun.*, vol. 13, no. 1, pp. 11-23, Jan. 1995.
- [12] Muhammad Aleem, Mohd Abdullah Al Mamun, Jalal Uddin Md Akbar, Hassan Ibrahim, Muhammad Midhat, Ehtesham Ali, Habiba Ashraf. (2024). Deep Learning-Powered Facial Expression Recognition: Revolutionizing Emotion Detection. *International Journal of Communication Networks and Information Security (IJCNIS)*, 16(4), 1426–1438. Retrieved from <https://www.ijcnis.org/index.php/ijcnis/article/view/7380>
- [13] Chatterjee P, Das D, Rawat DB. Digital twin for credit card fraud detection: Opportunities, challenges, and fraud detection advancements. *Future Generation Computer Systems*. 2024 Apr 30.
- [14] <https://ijcnis.org/index.php/ijcnis/article/view/7381> PROCESS Corporation, Boston, MA, USA. Intranets: Internet technologies deployed behind the firewall for corporate productivity. Presented at INET96 Annual Meeting. [Online]. Available: <http://home.process.com/Intranets/wp2.htm>
- [15] Hasan MR, Gazi MS, Gurung N. Explainable AI in Credit Card Fraud Detection: Interpretable Models and Transparent Decision-making for Enhanced Trust and Compliance in the USA. *Journal of Computer Science and Technology Studies*. 2024 Apr 6;6(2):01-12
- [16] Salekshahrezaee Z, Leevy JL, Khoshgoftaar TM. The effect of feature extraction and data sampling on credit card fraud detection. *Journal of Big Data*. 2023 Jan 17;10(1):6.
- [17] Gupta P, Varshney A, Khan MR, Ahmed R, Shuaib M, Alam S. Unbalanced credit card fraud detection data: a machine learning-oriented comparative study of balancing techniques. *Procedia Computer Science*. 2023 Jan 1;218:2575-84.
- [18] Al, M. a. E. (2024, February 4). Comparative analysis of deep learning models: CNN, MobileNetV2, and ResNet50 for offline signature verification. <https://lettersinhighenergyphysics.com/index.php/LHEP/article/view/931>
- [19] S. Al Balawi and N. Aljohani, "Credit-card fraud detection system using neural networks.," *Int Arab J Inf Technol*, vol. 20, no. 2, pp. 234–241, 2023.
- [20] Khalid AR, Owoh N, Uthmani O, Ashawa M, Osamor J, Adejoh J. Enhancing credit card fraud detection: an ensemble machine learning approach. *Big Data and Cognitive Computing*. 2024 Jan 3;8(1):6.
- [21] Muhammad Aleem, Noorhuzaimi Mohd Noor, Ehtesham Ali, Hassan Ibrahim, Mohd Abdullah Al Mamun, Muhammad Talha Bashir. (2024). Review of Mobile Application Performance Evaluation to Enhance Selection and Prediction in Mobile App Development. *International Journal of Communication Networks and Information Security (IJCNIS)*, 16(4), 1404–1413. Retrieved from <https://ijcnis.org/index.php/ijcnis/article/view/7378>.
- [22] Muhammad Aleem, Ehtesham Ali, Mohd Abdullah Al Mamun, Jalal Uddin Md Akbar, Khadeeja Saeed, Aqsa Saleem. (2024). Harnessing Ensemble Learning Approaches for Strong Mobile App Success Prediction Model. *International Journal of Communication Networks and Information Security (IJCNIS)*, 16(4), 1414–1425. Retrieved from <https://ijcnis.org/index.php/ijcnis/article/view/7379>.
- [23] Aslam, S., Usman Javeed, M. ., Maria Aslam, S. ., Iqbal, M. M., Ahmad, H. ., & Tariq, A. . (2025). Personality Prediction of the Users Based on Tweets through Machine Learning Techniques. *Journal of Computing & Biomedical Informatics*. Retrieved from <https://jcibi.org/index.php/Main/article/view/796>.
- [24] Muhammad Aleem, Ehtesham Ali, Mohd Abdullah Al Mamun, Jalal Uddin Md Akbar, Aqsa Saleem, Hassan Ibrahim, Abdullah Khan. (2024). Harnessing Deep Learning for Precision Tree Extraction in ArcGIS Pro. *International Journal of Communication Networks and Information Security (IJCNIS)*, 16(4), 1439–1453. Retrieved from <https://www.ijcnis.org/index.php/ijcnis/article/view/7381>.
- [25] Javeed, M., Aslam, S., Farhan, M., Aslam, M., & Khan, M. (2023). An Enhanced Predictive Model for Heart Disease Diagnoses Using Machine Learning Algorithms. *Technical Journal*, 28(04), 64-73. Retrieved from <https://tj.uettaxila.edu.pk/index.php/technical-journal/article/view/1828>.
- [26] M. U. Javeed, M. S. Ali, A. Iqbal, M. Azhar, S. M. Aslam and I. Shabbir, "Transforming Heart Disease Detection with BERT: Novel Architectures and Fine-Tuning Techniques," 2024 International Conference on Frontiers of Information Technology (FIT), Islamabad, Pakistan, 2024, pp. 1-6, doi: 10.1109/FIT63703.2024.10838424.