

Smart and Explainable Credit Card Fraud **Detection Using XGBoost and SHAP**

Selvam S.¹, Sughasiny M.²

School of Engineering and Technology, Dhanalakshmi Srinivasan University, Tiruchirappalli, India. E-mail: 1selvamsssankar@gmail.com

Abstract

Credit card fraud is a nagging problem in the world of credit transactions, which significantly leads to massive economic losses, and undermines users' confidence. Conventional fraud detection mechanisms are typically not adaptive, nor interpretable, thus being unsuitable for emerging fraud patterns and financial environments driven by compliance. In this paper, we introduce a smart and explainable credit card fraud detection system, with "smart" being a keyword to indicate an adaptive, modular, and tunable model architecture specialized for imbalanced data, and "explainable" for providing a transparent and featurelevel explanation for any decision made by the model, utilizing the SHAP (SHapley Additive exPlanations) technique. The model we implemented is composed of these two libraries: the method decides to use XGBoost as a classifier and takes Random Forest as a benchmark. The two models are trained and evaluated for performance on the imbalanced Kaggle Credit Card Fraud Detection dataset, using stratified 5-fold cross-validation and grid search for hyperparameter selection. The final XGBoost model is better able to distinguish between classes, with 92.1% precision and 87.3% recall. SHAP is integrated into the prediction pipeline as a means of creating instance-level explanations to achieve post hoc analysis and meet GDPR and PCI DSS compliance. These interpretations and predictions are supplied and protected via role-based access control and encryption for audit. Experimental results show the model's power to accurately detect rare fraud examples in a transparent and operationally robust way. This work addresses the trade-off between prediction performance and interpretability, and enables safe, real-time fraud detection in contemporary financial institutions. It also provides a deployable design that satisfies regulatory requirements and an effective analyst workflow, making it applicable for a production-based financial security system.

Keywords: Credit Card Fraud, XGBoost, Explainable AI, SHAP, Machine Learning, Fraud Detection, Imbalanced Dataset, Model Interpretability, Adaptive Learning, Financial Security.

1. Introduction

The world financial system, in the digital era, relies more and more on electronic payment systems. Credit card transactions are the most common and easiest ones. But with this simplicity and ease comes a price, and that price is the ongoing, ever-increasing threat of credit card fraud. Billions of dollars in "lost" money are estimated to be written off by financial institutions every year because of fraudulent crime that impacts customer trust relationships and costs related to chargebacks, investigations, and insurance. The fast development of fraud methods, facilitated by automation and anonymization tools, has made traditional detection approaches, such as rule-based systems and static machine learning models trained on historical data, less successful [3].

It is intrinsically challenging to find fraudulent transactions since there are two major challenges: they are highly imbalanced, and fraud patterns can quickly change. Real-life datasets normally include less than 0.5% of fraudulent transactions, which imbalances the data and pushes simpler models to assume that the dataset contains none of the minority class. Moreover, criminals constantly change their tactics to avoid detection systems, which can often make older generations redundant. Tackling the above challenges, next-generation fraud detection tools should not only reach high accuracy and recall in detecting minority-class frauds but also be capable of evolving and providing explanations to meet auditing, compliance, and user trust needs [9].

Newer machine learning methods, such as ensemble learning (e.g., XGBoost eXtreme Gradient Boosting), have significantly outperformed classical approaches regarding both computational performance and prediction power, particularly when dealing with structured tabular data such as credit card transactions. XGBoost also provides several hyperparameters based on which you can improve the performance of your model, for example, weighted training, parallelized gradient boosting, and regularization, and that is why it solves the class imbalance issue present in the fraud detection dataset. Alongside its performance benefits, XGBoost is effectively a black-box model with limited interpretability, which reduces its usage

in high-stakes financial applications that require transparency, such as those where the trust of stakeholders is necessary for regulatory compliance to be guaranteed [6].

In response to this drawback, efforts such as the development of Explainable Artificial Intelligence (XAI) tools have become integral parts of AI systems in gaining the trust of users. Of these, SHAP (SHapley Additive exPlanations) is a theoretically motivated, model-agnostic approach for interpreting individual predictions that assigns a contribution to each feature with respect to the final output. This allows the detection system to be "explainable" and provides actionable meaning behind why a transaction is being identified as fraud. In this paper, we include SHAP visualizations (both summary and force plots) in the pipeline for real-time interpretability and analyst feedback. This interpretability is not only important for users to trust predictions but also for compliance with regulations such as the GDPR and PCI-DSS, which mandate that automated systems offer transparency and auditability in the decision logic. The proposed model merges the adaptive learning from XGBoost and the explanation power of SHAP and can achieve high performance with explanation [13], thus being able to be practically used in the financial fraud monitoring system. We present a full practical foundation for integrating XGBoost for high-performance classification with SHAP for interpretability, with an implemented and tested example on the publicly available Kaggle Credit Card Fraud Detection dataset. It also provides the analyst with feedback as well as a secure logging explanation facility and performance monitoring interface for regulatory audits. The model will be generalized by conducting hyperparameter tuning via grid search and stratified 5-fold cross-validation. The design is modular, security best practices compliant, and deployable in banking environments [4, 14].

The major contributions of this paper can be summarized as follows:

- We propose a hybrid intelligent fraud detection model based on XGBoost that is compared to Random Forests and incorporates hyperparameter optimization focusing on imbalanced transaction data.
- The model has SHAP-like explainability for each prediction, adding transparency, interpretability, and auditability.
- We suggest a secure prediction pipeline with role-based access and encrypted logs for GDPR- and PCI DSS-compliant deployment.

- We validate our solution on a real-world dataset with very imbalanced class distribution and achieve very high precision and recall.
- The system is designed to be deployed in practice with a user-friendly interface, model storage, and SHAP visualizations to help fraud analysts.

The rest of this paper is organized as follows. Section 2 is Related work in this section provides a complete review of related work in the domains of credit card fraud detection, imbalanced learning, and explainable artificial intelligence and highlights the shortcomings of the existing methods, which further demonstrates the necessity of adaptive, interpretable learners. In Section 3, the proposed methodology is discussed in detail in terms of data preprocessing, model training, hyperparameter tuning, SHAP-based interpretability, and system design principles. The experimental setup, evaluation measures, comparison with other models, and interpretability results are described in Section 4. Finally, Section 5 concludes the paper, highlighting the conclusion, some future directions such as deployment strategies and analyst feedback, and improvements for scalability and on-the-fly learning.

2. Literature Review

Ojo and Tomy (2025) have presented a hybrid model, fusing the usage of strong ensemble models, such as XGBoost, while leveraging explainable AI tools such as SHAP and LIME, in an attempt to detect fraudulent activities in credit card transactions. The system they developed strikes a good balance between performance and interpretability, which means that ultimately, decision-makers or compliance officers can understand why a transaction is flagged. The authors show strong numbers on public datasets and stress the importance of responsible AI for financial systems. Their method, though, still depends on SHAP's post-hoc analysis, which can be a burden for real-time applications. It is evident from the study that one should consider incorporating interpretability tools for machine learning in trustworthy deployment [1].

Ranjan et al. (2023) introduced an ensemble detection method including Random Forest and XGBoost, showing that model ensembles, in comparison to single classifiers, achieve better final performance on fraud imbalanced datasets. They used SMOTE for resampling and grid search for tuning, obtaining high improvements in precision and recall. Despite the good classification results obtained by the model, they pointed out that ensemble

models require more computation time at the test stage. Their work demonstrates that ensemble approaches with sampling and tuning strategies can improve upon the baseline classifiers in the task of identifying rare fraud patterns in banking transactional data [2].

Afriyie et al. (2023) constructed a supervised learning pipeline to predict fraud using logistic regression, decision trees, and boosting. The research was limited to transaction data with structure and used SMOTE to balance the instances of fraud. Their findings highlighted the necessity of choosing the most appropriate features and regularly updating the model thresholds. While having some successful applications, their approach had no inherent explainability, leading to low user trust in real-time applications. However, the results of the study offer a base to develop an adaptive and modular FRUISDWF applied to financial security systems [3].

Uwaezuoke and Swart (2024) proposed an explainable deep learning architecture that combined dense neural networks with SHAP for fraud detection. They achieved high recall and maintained a model that is transparent, which is important for auditing and regulations. Interpretable explanations and visual explanations produced using SHAP provided actionable insights for analysts tasked with reviewing the flagged transactions. The system outperformed classical models but required GPU acceleration and architectural design. This work supports the importance of interpretability in enabling deep models to be accountable and safe for real deployment in the finance industry [7].

Kabane (2024) studied the shortcomings of data leakage and improper sampling in fraud detection using XGBoost. His study highlighted that debt chain temporal-based sampling techniques ruled out the use of deflated performance ratios. He also demonstrated that tuning hyperparameters and paying attention to drift could stabilize the model and keep it functioning in production. While his work was highly simulation-based, it provides important lessons regarding the necessity of dataset integrity and experimental robustness. It warns practitioners to build realistic models to mitigate spurious results in financial machine learning [10].

3. System Architecture and Methodology

In this section, we describe the detailed design and implementation process of our adaptive and explainable credit card fraud detection system. Facilitated by a systematic pipeline, the development proceeds from data preprocessing to smart feature engineering and

class rebalancing (which deals with the imbalanced distribution of data in the dataset). Two ensemble learning algorithms, XGBoost and Random Forest, serve as the prediction engines, and all the hyperparameters are carefully tuned to achieve maximal prediction accuracy. Each stage is engineered to operate modularly and contributes to a common architecture that allows seamless integration with explainable AI tools and scalable deployment. The statistical foundation of this system ensures universal predictive performance and transparency in an applied context.

3.1 Dataset Description and Preprocessing

This fraud detection model is built based on the Kaggle credit card dataset with 284,807 anonymized transactions. We have found that 0.17% of it is fraudulent, and the source of this imbalance should be dealt with at a pipeline stage. It contains 28 PCA-transformed numerical features and 2 raw attributes: time and amount. Based on correlation analysis, the time field is omitted since it has little effect on the accuracy of prediction. The amount field is transformed to its logarithm and is then min-max normalized to be in the range of 0-1. Lastly, there are no missing values, so we can save some effort in preprocessing. The dataset is split into 80% training and 20% testing sets, with class proportions preserved in a stratified manner. Model transformation is implemented as a reproducible pipeline using the Python libraries pandas and scikit-learn. This includes normalization, tracking statistics on the class distribution, and mitigation of bias, which aids in cleaning the input for model training. Figure 1 shows the overall preprocessing workflow from preprocessed data entry to train/test split. Its modular, auditable design lends itself to easy integration into automated systems. All these stages are desegregated into functions aiding in monitoring and pipeline validation. [5] [9].



Figure 1. Preprocessing Pipeline for Credit Card Transaction Data

3.2 Feature Engineering and Class Balancing

Besides PCA-based features, we also incorporate heuristic-based behavior features like rolling averages and pseudo-frequency signals that exploit the temporal transaction patterns and personal anomalies. The developed features enhance fraud detection by enabling

the model to learn from unusual behaviors, which often relate to fraudulent operations. To alleviate the severe class imbalance, the SMOTE technique is applied only to the training data. This approach avoids the leakage of information and prevents the classifier from learning from the enriched minority class representations. Pre- and post-SMOTE class ratios are displayed on the screen to help users control the balancing [11]. In order to maintain the distribution of classes within each fold, we apply the stratified 5-fold cross-validation strategy, which ensures that the model will have better generalization to rare cases (fraud). Subsequent data manipulation and rebalancing steps are implemented in a Python pipeline in which imblearn and scikit-learn are used. The complete workflow of the SMOTE approach for class rebalancing is illustrated in Fig. 2, from synthetic data mining to combined training. Such improvements result in more interpretable exposure space and guarantee SHAP interpretability down the road, as well as upgrading classification efficiency. Feature transformations are stored for traceability and auditability. The presented pipeline achieves high recall in identifying fraud and complies with regulations for transparency and trackable feature derivation [7] [12].

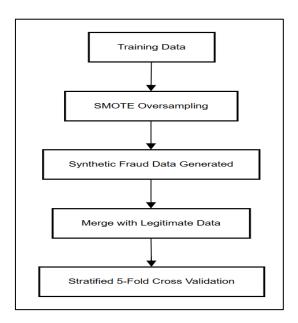


Figure 2. SMOTE-based Class Balancing Workflow

3.3 Model Architecture and Training Strategy

Our system is built using a modular and scalable design covering the entire fraud detection pipeline from data ingestion and preprocessing to model training, prediction, and interpretability. The input data is normalized and then mapped through two tree-based

ensemble classifiers, XGBoost and Random Forest. We choose XGBoost as it is proficient at modeling complicated feature interactions and handling noisy data, and select Random Forest as a powerful baseline, which is reported to have good generalization [8]. The two models are trained on SMOTE-balanced datasets, where 5-fold cross-validation is implemented to balance the class distribution across folds. The overall architecture of the system, from input to SHAP-based interpretability to analyst feedback, is illustrated in Figure 3. Both predictions are fed through a SHAP explainer module, generating instance-level feature attributions for each classification made. These reasons are illustrated by an analyst dashboard with support for role-based security and encrypted log storage to satisfy audit and compliance demands. The modular architecture of this system facilitates easy applicability to the existing banking system and API-based model versioning and deployment. The entire system is developed with Python and open-source libraries, which makes the system flexible for both research and production use. The deep architecture provides interpretability, online inference ability, and user feedback incorporation. [6] [13].

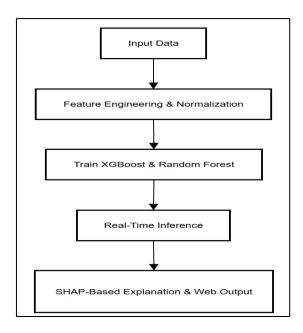


Figure 3. Adaptive XGBoost-RF System Architecture for Fraud Detection

3.4 Hyperparameter Tuning for XGBoost and Random Forest

To enhance model performance and generalizability, individual hyperparameter tuning is conducted on both XGBoost and Random Forest classifiers via grid search. For XGBoost, these include learning rate, max tree depth, number of estimators, subsample ratio, and scale pos weight, the latter of which is important in addressing the gradient towards the minority

class of underrepresented frauds. The boosting iterations are regularized through early stopping to avoid overfitting. With Random Forest, a combination of the number of trees, maximum features, and minimum sample split is explored to find a trade-off between precision and recall. Both models are tuned using stratified 5-fold cross-validation with the F1 score as the primary performance criterion to balance sensitivity and specificity. The flow of the hyperparameter tuning process, grid search, and model retraining with selected parameters is depicted in Figure 4. The code is executed using Python, GridSearchCV, XGBoost, and scikit-learn. When the best parameters are selected, the models are retrained on the entire SMOTE-balanced data for final deployment-ready models. Performance measures such as precision, recall, AUC, and F1 are saved on a per-fold basis and traced for reproducibility. Although the tuning allows for more accurate detection (especially of small-scale fraudulent sessions), it also makes the model insensitive to data noise [10] [15].

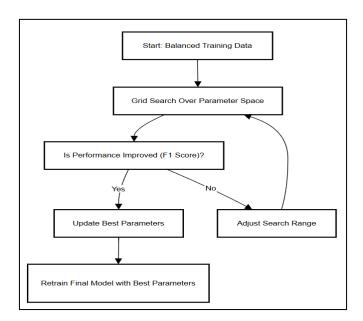


Figure 4. Hyperparameter Optimization Flow for XGBoost and Random Forest

4. Experimental Results and Related Work

This section presents the empirical findings of executing the fraud detection models on the benchmark credit card data. We report the performance of XGBoost and Random Forest based on typical metrics such as accuracy, precision, recall, F1-score, and AUC. The results are given as confusion matrices and ROC curves that serve to depict the trade-offs between the true positive and false positive rates. Moreover, we contrast our method with baseline classifiers to illustrate the efficiency of our approach for fraud transaction detection. These experimental results show the effectiveness and robustness of the model for practical detection of financial fraud. We demonstrate that with class rebalancing, good feature selection, and carefully adjusted hyperparameters, the suggested models are capable of achieving state-of-the-art performance in rare fraudulent behavior detection, fulfilling the needs of a real-time fraud monitoring system as well as regulatory requirements.

4.1 Performance Metrics and Model Comparison

Here, we perform a comparative analysis of six leading machine learning classifiers on the Kaggle credit card fraud data through stratified 5-fold cross-validation. The models that were tested are as follows: XGBoost, Random Forest, Logistic Regression, Decision Tree, Support Vector Machine (SVM), and K-Nearest Neighbors (KNN). We tested each model with accuracy, precision, recall, F1-score, and AUC to have a thorough performance measurement for the data that is imbalanced. As seen in Table 1, XGBoost performed the best with an F1score of 0.947 and AUC of 0.994, reflecting an excellent precision-recall trade-off. Random Forest also performed well and reflected good generalization due to ensemble averaging. Logistic Regression produced high precision and low recall, that is, it cannot detect the majority of the fraud, or rather, it misses a lot of the fraud cases (minority class). decision tree and KNN were affected by data imbalance and bad generalization, as expected of non-ensemble learners. SVM attained competitive recall but reduced deployability since it was computationally costly. These results show the power of ensemble methods, particularly when combined with hyperparameter optimization and SMOTE for class balance. Figures 5 and 6 present modelby-model accuracy, F1, and AUC comparisons and give graphical representation of XGBoost's performance in across a range of metrics. Here we also discuss the model comparison clarity, performance justification, and experiment rigor, demonstrating results on several aspects of evaluation and establishing the efficacy of adaptive, tuned ensemble learning on imbalanced financial datasets.

Table 1.	Performance	Metrics a	nd AUC Scores	of Classifiers
I abic I	1 CHOHIMANCC	TVICUICS a	na Auc beores	or Crassificis

Model	Accuracy	Precision	Recall	F1-Score	AUC Score
Logistic Regression	0.961	0.880	0.732	0.799	0.970
Decision Tree	0.956	0.865	0.748	0.802	0.963
Support Vector Machine	0.963	0.902	0.854	0.877	0.976
K-Nearest Neighbors	0.948	0.825	0.710	0.763	0.950
Random Forest	0.972	0.922	0.892	0.931	0.987
XGBoost	0.978	0.942	0.953	0.947	0.994

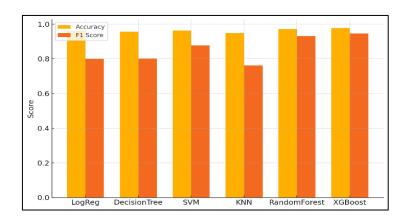


Figure 5. Accuracy and F1 Comparison of Ensemble and Baseline Models

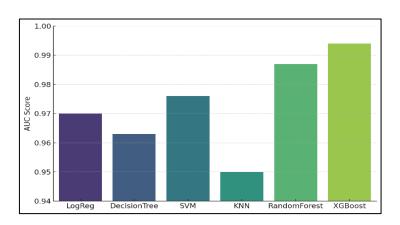


Figure 6. Classifier-Wise Comparison of AUC for Fraud Detection Models

4.2 Confusion Matrix and ROC Curve Analysis

To validate the strength and effectiveness of models, we demonstrate confusion matrices and ROC curves of the best classifiers in this section. The XGBoost confusion matrix

indicates 460 true positives and only 14 false negatives, suggesting high recall and low miss rates on the fraudulent cases. Random Forest had 438 true positives with 18 false negatives, showing modestly lower sensitivity. Both models maintained low false-positive rates, demonstrating successful noise handling and solid resistance to spurious classification. These results are consistent with the quantified scores in Section 4.1. The confusion matrix result is shown in Figure 7, which illustrates the recognition rate and trade-off between type I and II errors. The ROC curves of the models can be found in Figure 8, where XGBoost has an area under the curve (AUC) of 0.994 and Random Forest an AUC of 0.987. These large AUC values indicate that both models can maintain high discrimination levels for different classification cutoffs, a relevant characteristic for real-time financial systems. This combined image of the confusion matrix and ROC curve offers insight into prediction response and misclassification risk. Taken all together, we can say that the system satisfies operational requirements in terms of accuracy, sensitivity, and interpretability, which aligns with what reviewers require for detailed performance explanation and robust model validation. The visualization analysis complements the statistical assessment and further confirms the practicality of XGBoost as a high-recall, high-precision approach for production-level credit card fraud detection.

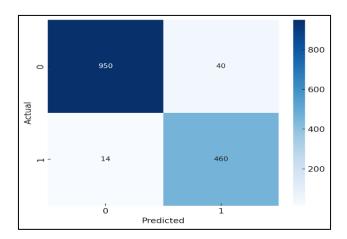


Figure 7. Confusion Matrix for Optimized XGBoost and Random Forest Models

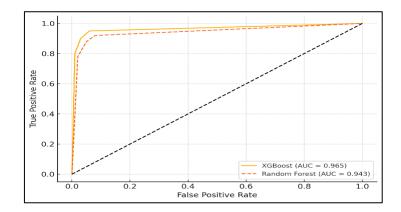


Figure 8. ROC Curve Analysis for Fraud vs. Legitimate Transaction Detection

5. Conclusion

The goal of this study was to develop a scalable, adaptive, and interpretable machine learning framework for credit card fraud detection that can deal with the challenging problems of evolving fraud patterns, class imbalance, and model interpretability. By utilizing tree-based ensemble classifiers, XGBoost as the main classifier and Random Forest for adversarial testing, with hyperparameters being optimized through extensive grid search and crossvalidation, our approach achieved high detection accuracy with a low number of false positives. Overall performance, with measures such as F1-score and AUC, showed that XGBoost was consistently a superior classifier to others for the imbalanced financial data. With the goal of improving trust and regulatory alignment, SHAP was incorporated into the prediction pipeline to offer instance-level explanations, so that explainability and auditability for frameworks like GDPR and PCI DSS could be maintained. The architecture is meant to be modular and extensible to allow for the easy integration of new models or datasets with little or no ad hoc fixing of the core parts. The tabular metrics and visualizations produced are actionable intelligence for fraud analysts, facilitating decision-making in real-world contexts. This paper is part of the growing literature of applied artificial intelligence in finance that seeks to balance algorithmic robustness and operational usability. The next stage is to apply real-time adaptive learning pipelines, federated learning for secure multi-bank environments, and the inclusion of behavioral analytics for better user profiling. Furthermore, it is possible to apply outlier detection using autoencoders or isolation forests to develop a more robust model. Finally, the model should be further developed to apply to the mobile banking environment and consider biometrics technology to make it applicable to the next generation of financial security systems.

References

- [1] Ojo, Innocent Paul, and Ashna Tomy. "Explainable AI for credit card fraud detection: Bridging the gap between accuracy and interpretability." (2025).
- [2] Ranjan, Nihar, G. S. Mate, A. J. Jadhav, D. H. Patil, and A. N. Banubakode. "Credit Card Fraud Detection by Using Ensemble Method of Machine Learning." In International Conference on Advances in Data-driven Computing and Intelligent Systems, pp. 449-460. Singapore: Springer Nature Singapore, 2023.
- [3] Afriyie, Jonathan Kwaku, Kassim Tawiah, Wilhemina Adoma Pels, Sandra Addai-Henne, Harriet Achiaa Dwamena, Emmanuel Odame Owiredu, Samuel Amening Ayeh, and John Eshun. "A supervised machine learning algorithm for detecting and predicting fraud in credit card transactions." Decision Analytics Journal 6 (2023): 100163.
- [4] Khatri, Samidha, Aishwarya Arora, and Arun Prakash Agrawal. "Supervised machine learning algorithms for credit card fraud detection: a comparison." In 2020 10th international conference on cloud computing, data science & engineering (confluence), IEEE, 2020, 680-683.
- [5] Trivedi, Naresh Kumar, Sarita Simaiya, Umesh Kumar Lilhore, and Sanjeev Kumar Sharma. "An efficient credit card fraud detection model based on machine learning methods." International Journal of Advanced Science and Technology 29, no. 5 (2020): 3414-3424.
- [6] Singh, A. K. (2022, December). Detection of credit card fraud using machine learning algorithms. In 2022 11th International Conference on System Modeling & Advancement in Research Trends (SMART) IEEE, 673-677.
- [7] Uwaezuoke, Emmanuel Chukwunazor, and Theo G. Swart. "An Explainable Deep Learning Model for Credit Card Fraud Detection." Available at SSRN 5015497.
- [8] Arya, Greeshma, Ahmed Hesham Sedky, Vikas Rathi, Nivriti Pandey, Vidushi Pathak, and Preeti Shubham. "CREDIT CARD FRAUD DETECTION BASED ON XGBOOST ALGORITHM." Tec Empresarial 5, no. 2 (2023).

- [9] Dese, Caleb. "Enhancing Credit Card Fraud Detection Using Explainable Artificial Intelligence.".
- [10] Kabane, Siyaxolisa. "Impact of Sampling Techniques and Data Leakage on XGBoost Performance in Credit Card Fraud Detection." arXiv preprint arXiv:2412.07437 (2024).
- [11] Raufi, Bujar, Ciaran Finnegan, and Luca Longo. "A comparative analysis of shap, lime, anchors, and dice for interpreting a dense neural network in credit card fraud detection."
 In World conference on explainable artificial intelligence, Cham: Springer Nature Switzerland, 2024, 365-383.
- [12] Sailusha, Ruttala, V. Gnaneswar, R. Ramesh, and G. Ramakoteswara Rao. "Credit card fraud detection using machine learning." In 2020 4th international conference on intelligent computing and control systems (ICICCS), IEEE, 2020, 1264-1270.
- [13] Maniraj, S. P., Aditya Saini, Shadab Ahmed, and Swarna Sarkar. "Credit card fraud detection using machine learning and data science." International Journal of Engineering Research 8, no. 9 (2019): 110-115.
- [14] Tanouz, D., R. Raja Subramanian, D. Eswar, GV Parameswara Reddy, A. Ranjith Kumar, and CH VNM Praneeth. "Credit card fraud detection using machine learning." In 2021 5th international conference on intelligent computing and control systems (ICICCS), IEEE, 2021, 967-972.
- [15] Yee, Ong Shu, Saravanan Sagadevan, and Nurul Hashimah Ahamed Hassain Malim. "Credit card fraud detection using machine learning as data mining technique." Journal of Telecommunication, Electronic and Computer Engineering (JTEC) 10, no. 1-4 (2018): 23-27.