# Credit Card Fraud Detection Using Machine Learning

*Team - 8 Information:* Ushaswini Punna, Sai Snehitha Ramisetty, Bhavana Chinnamgari, Uday Rohith

## Introduction

Credit card fraud is a significant issue faced by financial institutions and customers worldwide. Developing a fraud detection system has immediate real-world applications and implications, making it inherently interesting and relevant. However, it involves tackling hurdles like imbalanced datasets and evolving fraud tactics. Our project aims to address these challenges by leveraging machine learning techniques to enhance security and minimize financial losses. By utilizing advanced algorithms and PCA-based features, our system aims to accurately identify fraudulent transactions, benefiting all stakeholders. Success will be measured by reduced fraud complaints and chargebacks, alongside metrics like AUPRC. While there are risks of model inaccuracies, the payoff in terms of financial savings and security enhancements is significant. Costs are minimal, leveraging freely available tools and datasets, with a flexible development timeline. Validation will be conducted against industry benchmarks to ensure effectiveness.

## Problem Definition

**Formal Problem Definition:** The problem of credit card fraud detection involves identifying fraudulent transactions from a large set of transactions where the number of fraudulent cases is significantly smaller compared to non-fraudulent ones. The goal is to develop a system that can accurately detect fraudulent transactions in real-time, minimizing financial losses and enhancing security for both financial institutions and customers.

**Jargon-Free Version (Heilmeier Question #1)**

How can we create a system that quickly and accurately spots fake credit card transactions among millions of real ones, helping banks and customers save money and stay secure?

## Literature Survey

[1] Utilizing deep learning techniques, specifically OSCNN, for credit card fraud detection by integrating oversampling preprocessing and CNN architecture. Relevant for understanding the application of deep learning in fraud detection, providing insights into model architecture and preprocessing techniques applicable to the project. The paper may lack detailed analysis of specific challenges encountered in credit card fraud detection, such as handling imbalanced datasets or adapting to evolving fraud tactics.

[2] This paper evaluates anomaly detection techniques for credit card transactions and suggests Isolation Forest as the most effective for real-time fraud detection. Valuable insights into credit card fraud detection methods; Isolation Forest's accuracy and efficiency make it a promising solution for real-time fraud detection. Scalability with large data volumes, need for continuous monitoring and updates to adapt to evolving fraud tactics. Further exploration could enhance performance and applicability.

[3] The paper analyzes credit card fraud using machine learning, aiming to identify fraud types and propose mitigation strategies. Insights into machine learning for credit card fraud analysis; discusses fraud detection techniques and proposes strategies for enhanced detection. Potential shortcomings: data quality, model interpretability, scalability. Further improvements needed in these areas for better fraud analysis.

[4] Introduces a federated learning model for credit card fraud detection, tackling data privacy and class imbalance issues. Insights on implementing federated learning, preserving data privacy, and selecting effective resampling techniques for fraud detection improvement. Include complexity in implementation and computational overhead for privacy-preserving techniques. Research could focus on novel resampling methods and platform-specific optimization for improved performance.

[5] Utilizes ensemble learning with various supervised algorithms to detect credit card fraud, integrating data-level techniques for handling imbalanced data. Relevant for credit card fraud detection by enhancing accuracy and minimizing misclassification through ensemble learning and data-level techniques. Shortcoming may include requiring precise tuning of the ensemble model and parameters for optimal performance.

[6] Enhances credit card fraud detection with large-scale data mining techniques, addressing scalability, efficiency, skewed data, and variable error costs. Relevant for improving fraud detection capabilities by applying distributed data mining, potentially reducing losses from fraudulent activities. Challenges include real-world implementation complexity and adapting to evolving fraud patterns. Future research could focus on practical implementation and benchmarking against existing systems.

[7][8] Surveys credit card fraud detection methods, emphasizing efficiency amid rising fraud rates. Covers AI, data mining, fuzzy logic, ML, sequence alignment, and genetic programming. Provides an overview of fraud detection methods, aiding in understanding and evaluating approaches for the project. Offers insights into emerging techniques. Potential shortcomings: Lack of detailed analysis and empirical evaluation for each method. Might not address challenges like imbalanced data or evolving fraud tactics.

[9] Proposes using a Hidden Markov Model (HMM) for credit card fraud detection by modeling transaction processing sequences. Introduces a unique approach to complement existing techniques, broadening fraud detection capabilities. Challenges include capturing complex fraud patterns and requiring high-quality training data. Future research could enhance HMMs with additional features and improve scalability for real- world use.

[10] Emphasizes the importance of fraud detection methods in combating rising digital fraud, focusing on statistics and machine learning. Offers an overview of statistical and machine learning techniques in fraud detection, providing insights into their usefulness. Challenges include adapting to evolving fraud tactics and handling complex datasets, suggesting room for improvement.

[11] Explores performance measures for plastic card fraud detection tools, introducing metrics focused on minimizing costs and visualizing algorithm performance. Offers insights into assessing fraud detection tools for plastic card transactions, aiding in metric selection and system optimization. May lack practical implementation details; improvements could focus on addressing implementation challenges.

[12] Reviews statistical methods in consumer credit scoring, vital for categorizing credit applicants amidst rising demand. Offers insights into statistical techniques relevant to developing a fraud detection system, aiding in understanding credit scoring for identifying fraud. Possible limitations in accessing public literature due to commercial confidentiality. Future research could validate techniques using publicly available datasets for credit scoring and fraud detection.
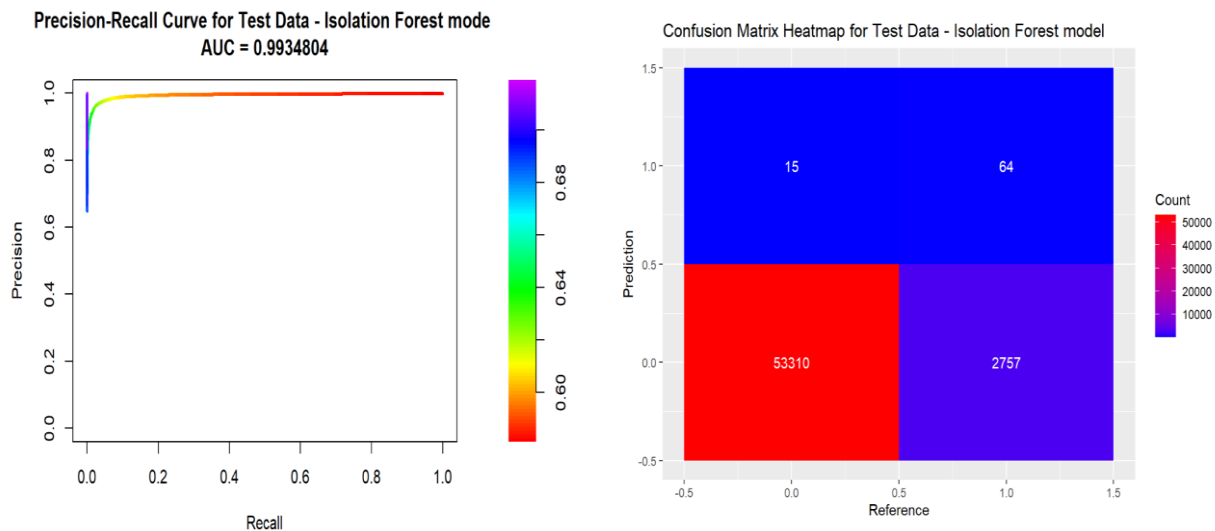
## Proposed Methods:

**Intuition:** By integrating various advanced machine learning algorithms and anomaly detection techniques, our proposed method aims to enhance the accuracy and efficiency of fraud detection. Leveraging PCA-based features and ensemble learning methods, our system is expected to outperform state-of-the-art techniques in identifying fraudulent transactions, addressing issues like imbalanced datasets and evolving fraud tactics effectively. We have performed the following steps to train and evaluate all the models for credit card fraud detection:

1. Train-test split: The data is divided into training (80%) and test (20%) sets with reproducibility ensured by setting a seed.
2. Train model: An model is trained on the training data excluding the target variable.
3. Predict anomalies: Anomaly scores are calculated for both training and test datasets. Anomalies are identified using the 95th percentile threshold of the training scores.
4. Evaluate the model: Confusion matrices for both datasets are created and printed to assess
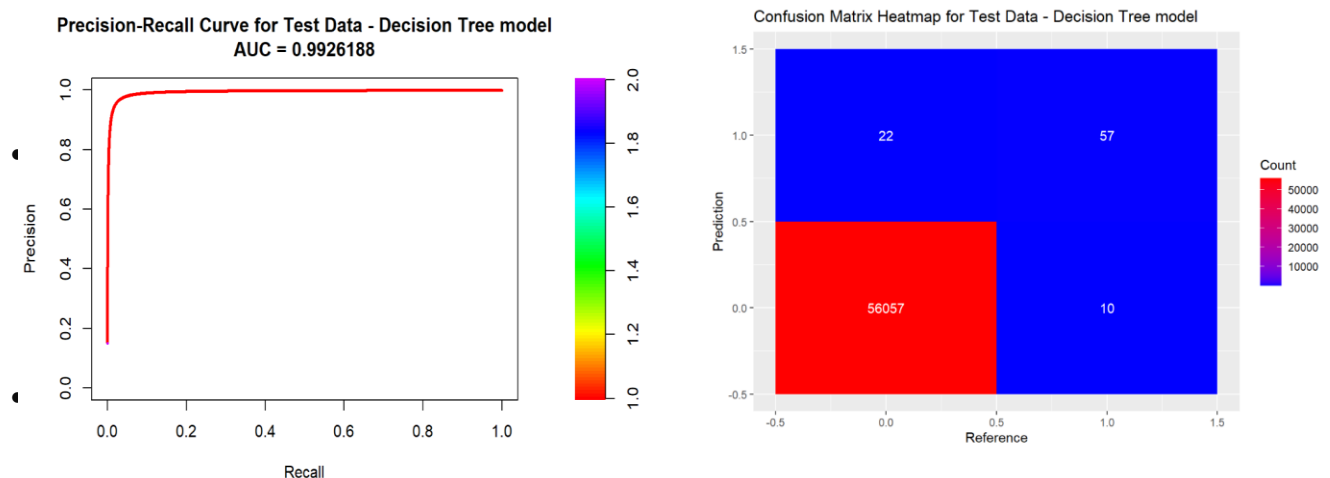
model performance.

5. Identify fraudulent transactions: Predictions are added to the test data, and fraudulent transactions are displayed.
6. ROC and Precision-Recall Curves: These curves are plotted for the test data to visualize model performance.
7. Confusion matrix heatmap: A heatmap is generated to visually represent the model's classification accuracy on the test data.

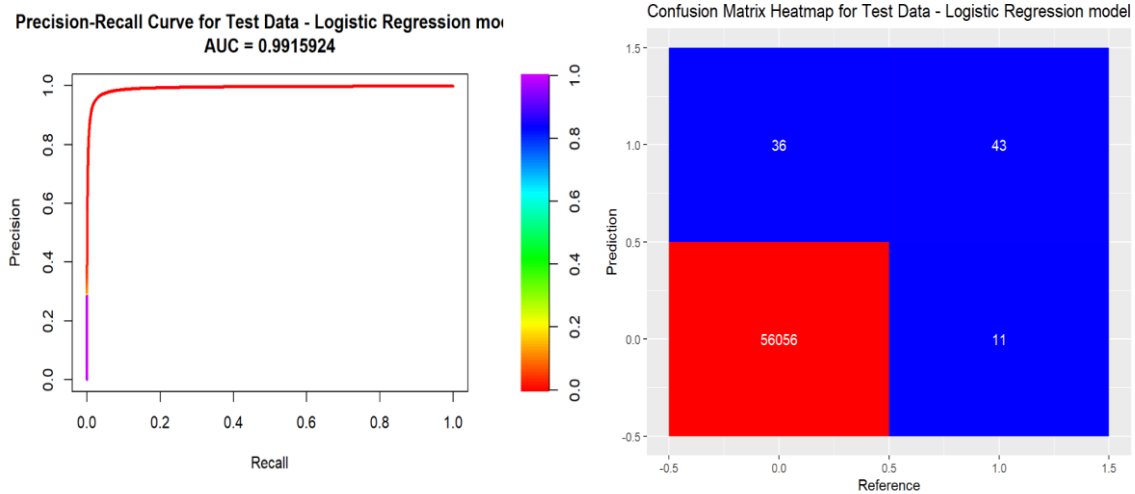**Description of Algorithms:**

- **Isolation Random Forest:** The Isolation Forest model is an unsupervised machine learning algorithm used for anomaly detection. It operates on the principle that anomalies are few and different, making them easier to isolate compared to normal data points. The model constructs a set of random binary trees, where each split isolates a data point. Anomalies, being distinct, tend to be isolated quickly, resulting in shorter path lengths within these trees. The average path length of a data point across many trees is used to determine its anomaly score, with shorter average path lengths indicating higher likelihood of being an anomaly. This method is efficient, scales well to large datasets, and does not require labeled data.
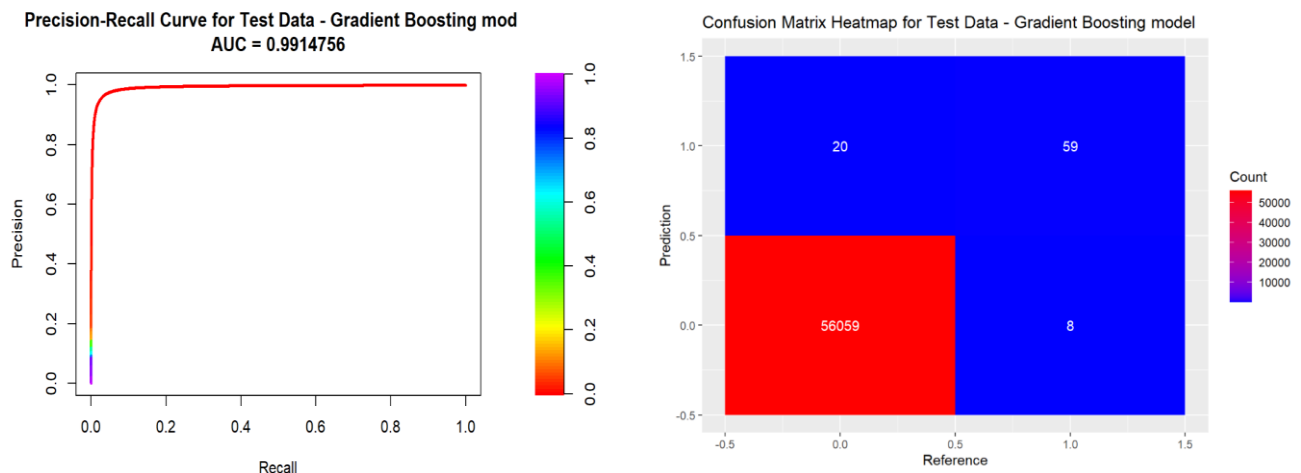


- **Decision Tree:** Decision Tree is a non-parametric supervised learning method used for both classification and regression tasks. It partitions the feature space into segments, making decisions based on simple rules inferred from the data. Each internal node represents a feature, each branch represents a decision based on that feature, and each leaf node represents the outcome. We employed the rpart package in R to build a Decision Tree model. The rpart function constructs the tree using the Classification and Regression Trees (CART) algorithm. The resulting model was visualized using the prp function to illustrate the decision-making process.

**ogistic Regression Model:** Logistic regression is a statistical method used for binary classification tasks, predicting the probability of one of two possible outcomes. The model equation combines input features linearly, and parameters are estimated using Maximum Likelihood Estimation (MLE). The code trains a logistic regression model on training data and predicts probabilities for the test data, classifying transactions based on a 0.5 threshold. It evaluates the model using a confusion matrix, displays predicted fraudulent transactions, and visualizes performance with ROC and Precision-Recall curves. Finally, it generates a heatmap of the confusion matrix to show classification results visually.



- **XGBoost Model:** XGBoost, short for Extreme Gradient Boosting, is a highly efficient and scalable implementation of the gradient boosting framework. It sequentially builds an ensemble of decision trees, correcting errors made by previous trees to improve overall model accuracy. With built-in regularization to prevent overfitting, XGBoost excels in handling large datasets and supports various objective functions and hyperparameters for fine-tuning. We had set the seed for reproducibility, then train the model with specified parameters on training features and labels. Predict class probabilities on test data, convert probabilities to binary classes, and evaluate using a confusion matrix. Add predictions to the test data and identify predicted fraudulent transactions. Generate and plot ROC and Precision-Recall curves to assess performance, and create a heatmap of the confusion matrix for visual evaluation of classification results.

## Experiments & Evaluation:

**Description of Testbed:**

Our experiments are designed to evaluate the effectiveness of different machine learning algorithms and anomaly detection techniques in identifying fraudulent credit card transactions. The key questions include:

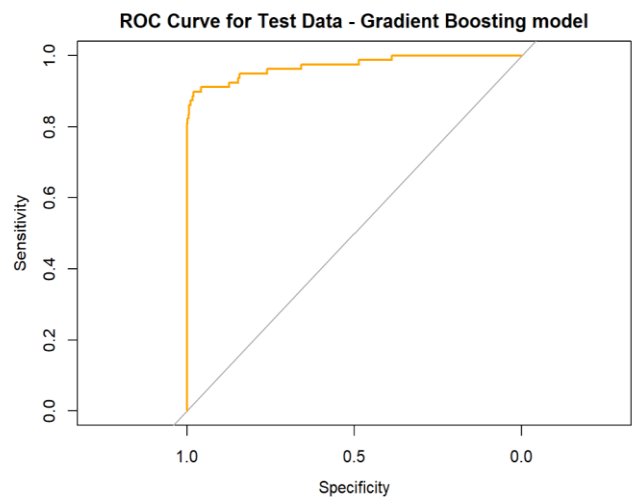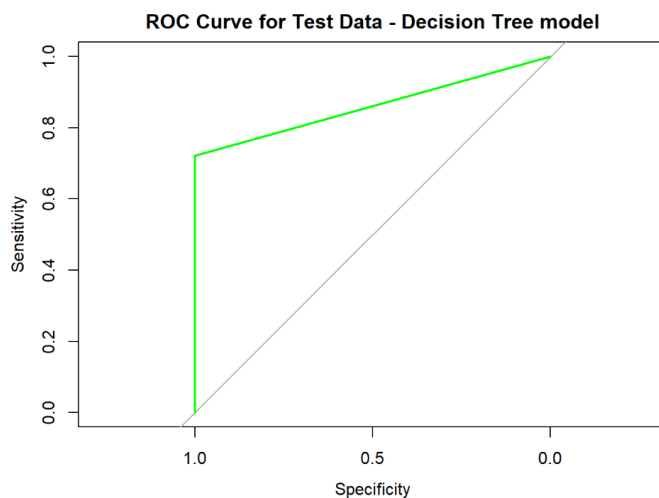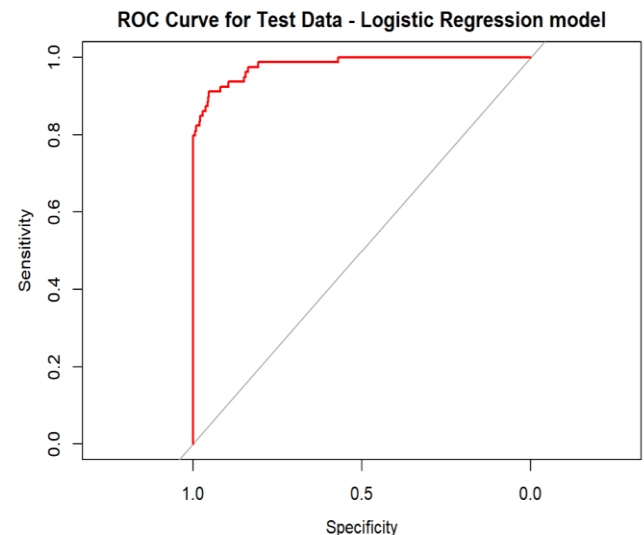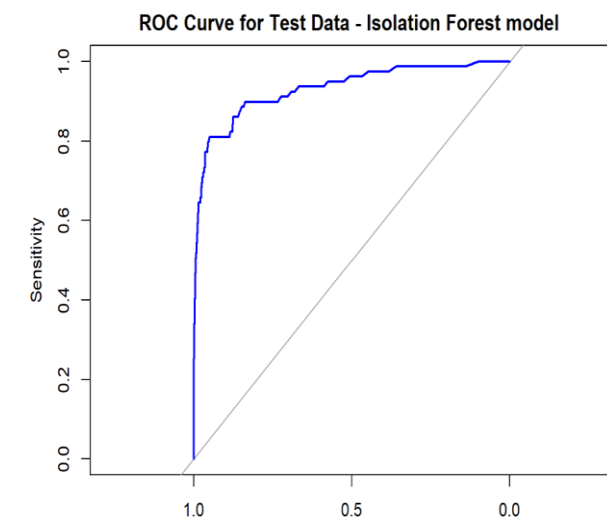How accurately can each algorithm detect fraudulent transactions?

How well do the models handle imbalanced datasets?

What are the performance metrics (e.g., AUPRC, accuracy, precision, recall) for each model?

**Description of Data:**

The dataset used in this project comprises credit card transactions made by European cardholders in September 2013, totaling 284,807 transactions. Of these, only 492 are fraudulent, highlighting the highly imbalanced nature of the data, with fraudulent transactions constituting a mere 0.172%. The features include the transaction time, amount, and anonymized variables derived from a PCA transformation, providing a robust foundation for detecting fraud amidst a predominantly non-fraudulent dataset.

## Observations & Results:



ROC Curve for Test Data - Isolation Forest model



ROC Curve for Test Data - Logistic Regression model



ROC Curve for Test Data - Decision Tree model



ROC Curve for Test Data - Gradient Boosting model

```
## Fraudulent transactions predicted by the model:
```

```
head(fraudulent_transactions)
```

```
##          Time        V1        V2          V5          V6          V9
## 221 -1.993313 -1.306267  1.303301 -0.2571324  1.89916501 0.614042788
## 222 -1.993313 -1.306267  1.303301 -0.2571324  1.89916501 0.614042788
## 224 -1.993313 -1.305774  1.304315 -0.2597587  1.90091651 0.613876435
## 226 -1.993271 -1.448397  2.999577  1.0106762 -1.26802022 3.948326704
## 418 -1.990008 -0.544400  1.154011 -0.6655057 -0.05184682 0.005856417
## 469 -1.989124 -1.877608 -2.889081  0.8656660 -0.94494998 1.529419992
##            V10         V11        V12        V14        V15        V16
## 221  0.13744134 -1.19998956  1.3145462 -0.2986622 -2.2044032 -0.8863493
## 222  0.13744134 -1.19998956  1.3145462 -0.2986622 -2.2044032 -0.8863493
## 224  0.13709520 -1.19987504  1.3142392 -0.2988880 -2.2043625 -0.8860504
## 226  5.79713772  2.41778766 -0.7287479 -6.8586143  1.9326695 -0.7473123
## 418 -0.01176964  0.07079537  0.9467465  0.3822799  0.8726434  0.2748692
## 469 -1.80702723 -0.72658489  0.4668575 -0.4622043  0.1340417 -1.2339759
```

*(Each team member has contributed an equitable amount of effort.)*

## Comparative Analysis

The following table summarizes the performance of all models:

```
                  Model  Accuracy    AUC_ROC     AUC_PR
1       Isolation Forest 0.9506287 0.9313276 0.9934804
2    Logistic Regression 0.9991629 0.9803463 0.9915924
3          Decision Tree 0.9994301 0.8606703 0.9926188
4                XGBoost 0.9995013 0.9714112 0.9914756
```

## CONCLUSION

The development of a fraud detection system using machine learning techniques represents a critical step in addressing the pervasive issue of credit card fraud. By overcoming challenges such as imbalanced datasets and evolving fraud tactics, our project aims to enhance security and reduce financial losses for both financial institutions and customers. Leveraging advanced algorithms and PCA-based features, our system promises to accurately identify fraudulent transactions, ultimately benefiting all stakeholders involved.

This report outlines our credit card fraud detection project, employing diverse machine learning techniques to detect fraudulent transactions. Models investigated encompass Isolation Forest, Logistic Regression, Decision Tree, and XGBoost. Key findings indicate that XGBoost achieves the highest accuracy (0.9995), Logistic Regression excels in AUC-ROC (0.9803), and Isolation Forest leads in AUC-PR (0.9935).

While there are risks associated with model inaccuracies, the potential payoff in terms of financial savings and enhanced security far outweighs these concerns. With minimal costs and a flexible development timeline, our project is poised to make a significant impact in combating credit card fraud. Future research could focus on integrating additional features, optimizing models for real-world implementation, and exploring novel resampling methods and platform-specific optimizations to enhance performance and applicability.

## REFERENCES

[1] Aya Abd El Naby; Ezz El-Din Hemdan; Ayman El-Sayed *"Deep Learning Approach for Credit Card Fraud Detection."* International Conference on Electronic Engineering (ICEEM) 2021 https://ieeexplore.ieee.org/document/9480639

[2] Soumaya Ounacer, Houda Jihal, Soufiane Ardchir & Mohamed Azzouazi , "*Anomaly Detection in Credit Card Transactions*", Advanced Intelligent Systems for Sustainable Development (AI2SD'2019) https://link.springer.com/chapter/10.1007/978-3-030-36674-2_14

[3] N. R. Shetty, N. H. Prasad, H. C. Nagaraj, "*Credit Card Fraud Analysis Using Machine Learning*", Advances in Communication and Applications Proceedings of ERCICA 2023, Volume 2 https://link.springer.com/chapter/10.1007/978-981-99-7633-1_21

[4] John MacIntyre, "*Federated learning model for credit card fraud detection with data balancing techniques*", Neural Computing and Applications https://link.springer.com/article/10.1007/s00521-023-09410-2

[5] Raunak Chhabra, Shailza Goswami & Ranjeet Kumar Ranjan, *"A voting ensemble machine learning based credit card fraud detection using highly imbalance data"*, Springer - Multimedia Tools and Applications https://link.springer.com/article/10.1007/s11042-023-17766-9

[6] P.K. Chan; W. Fan; A.L. Prodromidis, "*Distributed data mining in credit card fraud detection"*, IEEE Intelligent Systems and their Applications ( Volume: 14, Issue: 6, Nov.-Dec. 1999) https://ieeexplore.ieee.org/abstract/document/809570

[7] S. Benson Edwin Raj; A. Annie Portia, *"Analysis on credit card fraud detection methods*", 2011 International Conference on Computer, Communication and Electrical Technology (ICCCET). https://ieeexplore.ieee.org/abstract/document/5762457

[8] Ghosh; Reilly, "*Credit card fraud detection with a neural-network",* 1994 Proceedings of the Twenty- Seventh Hawaii International Conference on System Sciences https://ieeexplore.ieee.org/abstract/document/323314

[9] Abhinav Srivastava; Amlan Kundu; Shamik Sural; Arun Majumdar, *"Credit Card Fraud Detection Using Hidden Markov Model",* IEEE Transactions on Dependable and Secure Computing ( Volume: 5, Issue: 1, Jan.-March 2008) https://ieeexplore.ieee.org/abstract/document/4358713

[10] Richard J. Bolton and David J. Hand, *"Statistical Fraud Detection: A Review",* Institute of Mathematical Statistics Vol. 17, No. 3 (Aug., 2002), https://www.jstor.org/stable/3182781?searchText=credit+card+fraud+detection&searchUri=%2Faction%2FdoBasicSearch%3FQuery%3Dcredit%2Bcard%2Bfraud%2Bdetection%26so%3Drel&ab_segments=0%2Fbasic_search_gsv2%2Fcontrol&refreqid=fastly-default%3Ad61965cfa5e5d59ba67fd37aa4c5cc39

[11]     D. J. Hand, C. Whitrow, N. M. Adams, P. Juszczak and D. Weston, *"Performance Criteria for Plastic Card Fraud Detection Tools"*, The Journal of the Operational Research Society Vol. 59, No. 7 (Jul., 2008), pp. 956-962

https://www.jstor.org/stable/20202156?searchText=credit+card+fraud+detection&searchUri=%2Faction

%2FdoBasicSearch%3FQuery%3Dcredit%2Bcard%2Bfraud%2Bdetection%26so%3Drel&ab_segments

=0%2Fbasic_search_gsv2%2Fcontrol&refreqid=fastly-default%3Ad61965cfa5e5d59ba67fd37aa4c5cc39

[12]     D. J. Hand, and W. E. Henley, *"Statistical Classification Methods in Consumer Credit Scoring: A Review"*, Journal of the Royal Statistical Society. Series A (Statistics in Society)Vol. 160, No. 3 (1997)

https://www.jstor.org/stable/2983268?searchText=credit+card+fraud+detection&searchUri=%2Faction%

2FdoBasicSearch%3FQuery%3Dcredit%2Bcard%2Bfraud%2Bdetection%26so%3Drel&ab_segments=0

%2Fbasic_search_gsv2%2Fcontrol&refreqid=fastly-default%3Ad61965cfa5e5d59ba67fd37aa4c5cc39