



## **Insurance Fraud Detection using Machine Learning**

Manasa Madiraju  
Rohit Roy Chowdhury

# Abstract

Insurance fraud is a significant issue that leads to substantial financial losses for insurance companies. The objective of this project is to explore the application of machine learning techniques, specifically XGBoost and Random Forest, to detect fraudulent claims in insurance datasets. By leveraging these models, we aim to identify patterns and anomalies that indicate fraudulent activities, thus helping insurers mitigate risks and reduce fraudulent claims. This report covers the theoretical background, experimental setup, results, and conclusions drawn from the study.

## Introduction

The insurance industry is vast and complex, offering diverse policies like car, rental, health, and life insurance. Sadly, insurance fraud is widespread, occurring throughout the process – from policy sales to claiming. In response, insurance companies need sophisticated tools, including data analytics and machine learning, to detect and deter these fraudulent activities. Fraud can be perpetrated by individuals, groups, policyholders, insurers, or third parties. Due to its scope and complexity, robust solutions are essential. Insurance fraud costs both insurers and policyholders billions annually, taking forms like:

- False claims
- Fake policies
- Identity theft

This project aims to develop a machine learning model that analyzes historical claims and policyholder data to detect insurance fraud. Our goal is to uncover patterns and trends that signal fraudulent behavior.

Machine learning offers a promising solution by automating and enhancing the detection process through sophisticated algorithms that can analyze large datasets and identify fraudulent patterns. This report focuses on two powerful machine learning algorithms: XGBoost and Random Forest, both of which are known for their robustness and efficiency in handling classification problems. The following sections will delve into the theoretical background, experimental setup, results, and conclusions of our study on using these models for insurance fraud detection.

## Types of Insurance Fraud

1. **Hard Fraud:** Involves deliberate and planned actions to fabricate or cause a loss in order to claim insurance benefits. Examples include staging accidents, arson for profit, and faking deaths.

- **Staging Accidents:** Individuals deliberately cause accidents and then file insurance claims for damages and injuries. This often involves multiple conspirators who act as witnesses or involved parties.
- **Arson for Profit:** Policyholders intentionally set fire to insured property, such as homes or businesses, to claim the insurance payout. This type of fraud is particularly dangerous and costly.
- **Faking Deaths:** In extreme cases, individuals may fake their own death or the death of a family member to claim life insurance benefits.

2. **Soft Fraud:** Occurs when a policyholder exaggerates a legitimate claim or provides false information to receive a higher payout. Examples include inflating the value of stolen items or claiming for damages that were pre-existing.

- **Inflating Claim Amounts:** Claimants might inflate the value of items lost or damaged in an insured event to receive a higher payout than they are entitled to.
- **Pre-existing Damages:** Claimants may include damages that occurred before the insured event in their claims to receive compensation for old damages.

## Machine Learning Approaches

We have chosen two popular supervised machine learning algorithms to build our model for insurance fraud detection.

### Random Forest

Random Forest is an ensemble learning method that operates by constructing a multitude of decision trees during training and outputs the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. It reduces overfitting by averaging multiple trees, thus improving accuracy and robustness.

## Advantages

- **High Accuracy:** Generally provides high accuracy due to the ensemble approach, which reduces the variance of the predictions.
- **Feature Importance:** Provides insights into which features are most important for predictions, helping to understand the factors influencing fraudulent behavior.
- **Robustness to Overfitting:** Averaging over many trees reduces the risk of overfitting, especially when using a large number of trees. Random Forest is also less sensitive to noisy data.

## Disadvantages

- **Complexity:** Can be computationally intensive and require significant resources for training, particularly with large datasets and a high number of trees.
- **Interpretability:** While individual trees are interpretable, the ensemble of trees can be less so. The model's complexity can make it challenging to understand the decision-making process.

## XGBoost

XGBoost (Extreme Gradient Boosting) is a scalable and efficient implementation of gradient boosting framework by Friedman et al. It uses a more regularized model formalization to control overfitting, which provides better performance. XGBoost is known for its speed and performance in machine learning competitions.

## Advantages

- **Efficiency:** Highly optimized for speed and performance, capable of handling large datasets efficiently.
- **Handling of Missing Values:** Automatically learns the best imputation strategy from the data, making it robust to incomplete datasets.
- **Regularization:** Includes L1 and L2 regularization to control overfitting, enhancing the model's generalization capabilities.

## Disadvantages

- **Complexity:** Requires careful tuning of hyperparameters to achieve optimal performance. The model's complexity can be a barrier to understanding and implementation.
- **Resource Intensive:** Can be computationally demanding, especially with large datasets. The training process may require significant memory and processing power.

## Evaluation Metrics

To evaluate the performance of the machine learning models, several metrics are used:

- **Accuracy:** The proportion of correctly classified instances among the total instances. Accuracy is a straightforward metric but can be misleading in the case of imbalanced datasets, where the number of legitimate claims vastly outnumbers fraudulent claims.
- **Precision:** The proportion of positive identifications that were actually correct. Precision is important in fraud detection to minimize the number of false positives, which can be costly and time-consuming to investigate.
- **Recall:** The proportion of actual positives that were correctly identified. Recall is crucial in ensuring that a high number of fraudulent claims are detected.
- **F1 Score:** The harmonic mean of precision and recall, providing a balance between the two. F1 Score is particularly useful when dealing with imbalanced datasets, as it considers both false positives and false negatives.
- **ROC-AUC:** Area under the Receiver Operating Characteristic curve, which evaluates the trade-off between true positive rate and false positive rate. A higher AUC indicates a better performing model.

# Experimental Steps

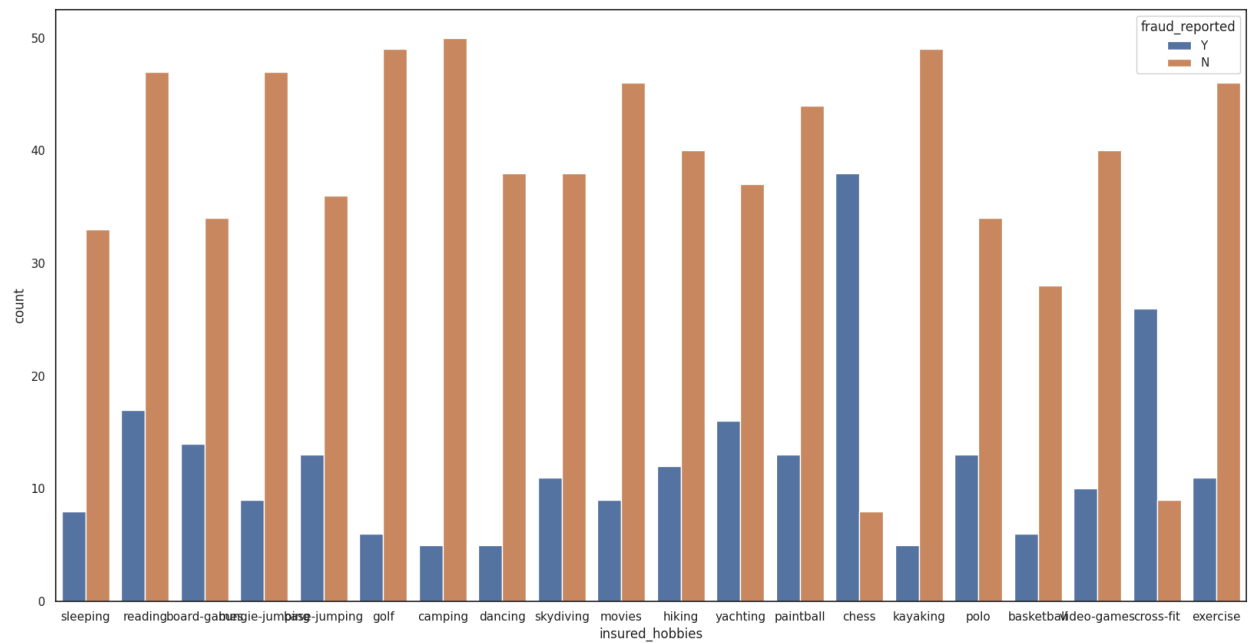
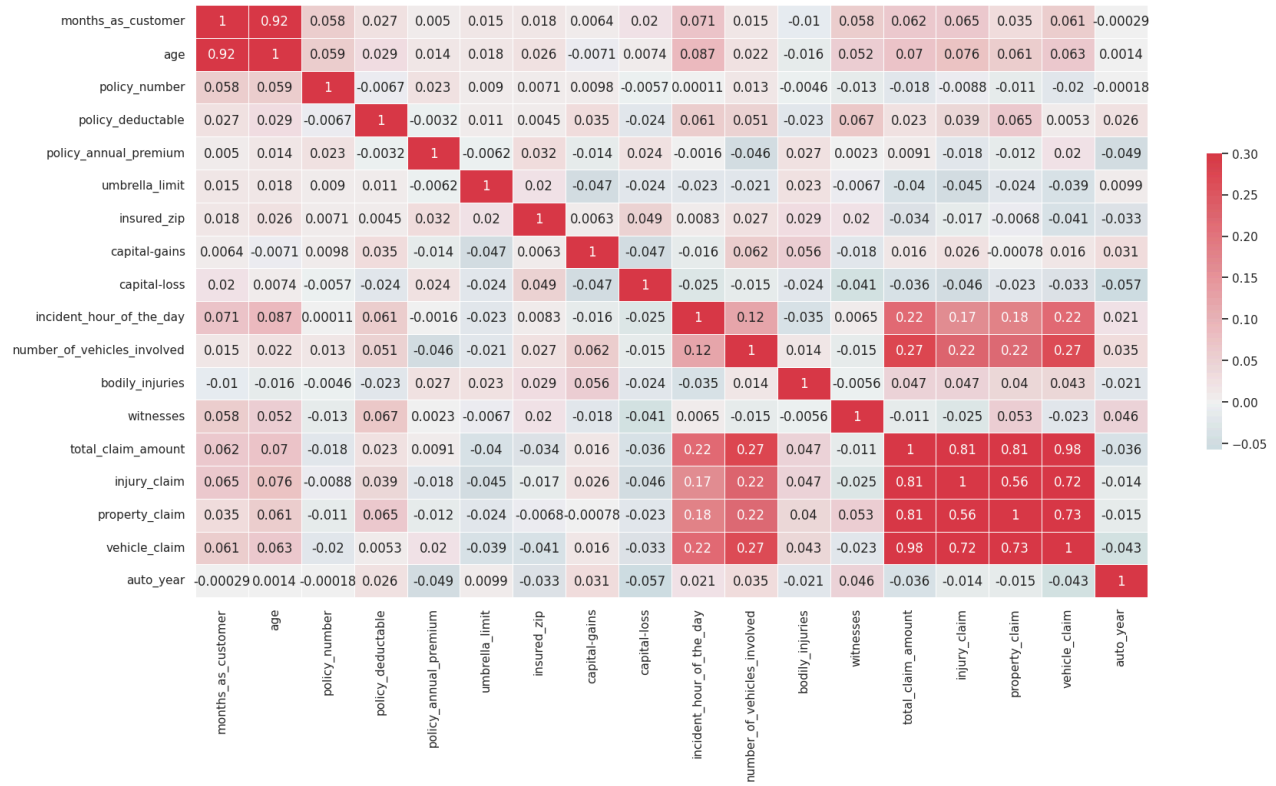
## Data Collection

The dataset used for this project includes historical insurance claim data, which is publicly available or obtained from an insurance company. The dataset typically contains features such as:

- *Claim Amount*: The monetary value of the claim. Higher claim amounts might be indicative of fraudulent activity, especially if they exceed a certain threshold.
- *Policy Details*: Information about the insurance policy, such as policy type, coverage limits, and duration. Certain policy types may be more prone to fraud.
- *Claimant Information*: Demographic and other details about the claimant, including age, gender, and past claim history. Repeat claimants or those with a history of fraud may be flagged.
- *Incident Details*: Information about the incident leading to the claim, such as date, location, and description. Patterns in incident details can help identify fraudulent claims.
- *Label*: Indicating whether the claim was fraudulent or legitimate. This is the target variable used for training supervised learning models.

## Data Preprocessing

The initial data preparation involved loading a dataset containing various details about insurance policies, claims, and reported fraud. The dataset was loaded from a CSV file using Pandas. The first step was to inspect the dataset for missing values and understand the data types of each feature. Missing values were handled by either dropping columns with a significant amount of missing data or replacing placeholders such as '?' with a more appropriate category like 'Unknown'. For instance, columns such as *police\_report\_available*, *property\_damage*, and *authorities\_contacted* were cleaned by replacing '?' and *pd.NA* with 'Unknown' and 'None', respectively. For feature engineering, we identified redundant features using a correlation matrix, consolidated dominant categories (for instance. *Claimant Hobbies*), and checked for class imbalance to detect bias.



## Data Processing

Once the initial cleaning was complete, the data processing steps included encoding categorical variables into numerical formats suitable for machine learning models. Categorical responses were converted using Label Binarizer, and categorical features were transformed using One Hot Encoding. This process expanded the dataset to 93 columns, ensuring all categorical data was properly represented numerically. This prepared dataset was then exported for model building.

## Model Implementation

### Random Forest

**Hyperparameter tuning:** In the Random Forest model, the following hyperparameters were used for tuning –

- *n\_estimators*: This parameter specifies the number of trees in the forest. More trees usually improve the model's performance but also increase the computation time.
- *max\_features*: This parameter determines the number of features to consider when looking for the best split. The options included:
  - *'sqrt'*: The square root of the total number of features.
  - *'log2'*: The base-2 logarithm of the total number of features.
  - *None*: The total number of features.
- *max\_depth*: This parameter sets the maximum depth of the trees. Limiting the depth of the tree helps prevent overfitting. The values used were 3, 5, and 7.
- *max\_leaf\_nodes*: This parameter restricts the number of leaf nodes in the tree. Fewer leaf nodes reduce model complexity and help in avoiding overfitting. The values tested were 3, 6, and 9.
- *min\_samples\_split*: This parameter specifies the minimum number of samples required to split an internal node. The values used were 2, 4, and 6.

The optimization of hyperparameters was performed using Grid Search with Cross-Validation. Grid Search exhaustively tests all possible combinations of the



specified hyperparameter values to find the best combination that maximizes the model's performance

**Training:** The model is trained on the training set with cross-validation to tune hyperparameters. Random Forest's inherent feature importance helps in understanding which features are most influential in predicting fraud.

Cross-Validation: k-fold cross-validation is used to ensure the model's generalizability and to fine-tune hyperparameters. This involves dividing the dataset into k subsets, training the model on k-1 subsets, and validating it on the remaining subset. This process is repeated k times.

**Evaluation:** The trained model is evaluated on the test set using the defined metrics. The evaluation metrics provide insights into the model's performance and its ability to generalize to new data.

## XGBoost

**Hyperparameter tuning:** In the XGBoost model, the following hyperparameters were used for tuning:

- *n\_estimators*: Number of boosting rounds.
- *learning\_rate*: Controls the weight of each boosting round.
- *max\_depth*: Maximum depth of each tree.
- *min\_child\_weight*: Minimum sum of instance weight (hessian) needed in a child.
- *subsample*: Fraction of samples used for fitting individual trees.
- *colsample\_bytree*: Fraction of features used for fitting individual trees.
- *gamma*: Minimum loss reduction required to make a split.
- *reg\_alpha*: L1 regularization term on weights.
- *reg\_lambda*: L2 regularization term on weights.
- *max\_leaves*: Maximum number of leaves in a tree.
- *max\_bin*: Maximum number of bins for discretizing continuous features.
- *scale\_pos\_weight*: Controls the balance of positive and negative weights.

**Training:** The final model, trained with the optimized hyperparameters, demonstrated improved performance metrics, providing a robust and accurate classification solution for the insurance fraud detection dataset.

**Evaluation:** The trained model is evaluated on the test set using the defined metrics. The evaluation metrics provide a comprehensive understanding of the model's strengths and weaknesses.

### Training and Validation

- **Train-Test Split:** The dataset is split into training and testing sets, typically in an 80:20 ratio. This ensures that the model is trained on a subset of the data and validated on a separate subset to assess its performance.
- **Cross-Validation:** k-fold cross-validation is used to ensure the model's generalizability and to fine-tune hyperparameters. This technique helps prevent overfitting and ensures that the model performs well on unseen data.

## Results

### Random Forest

- *Accuracy:* 83%
- *Precision:* Positive class - 67%, Negative class - 91%
- *Recall:* Positive class - 76%, Negative class - 86%
- *F1 Score:* Positive class - 81%, Negative class - 88%
- *ROC-AUC:* 0.81

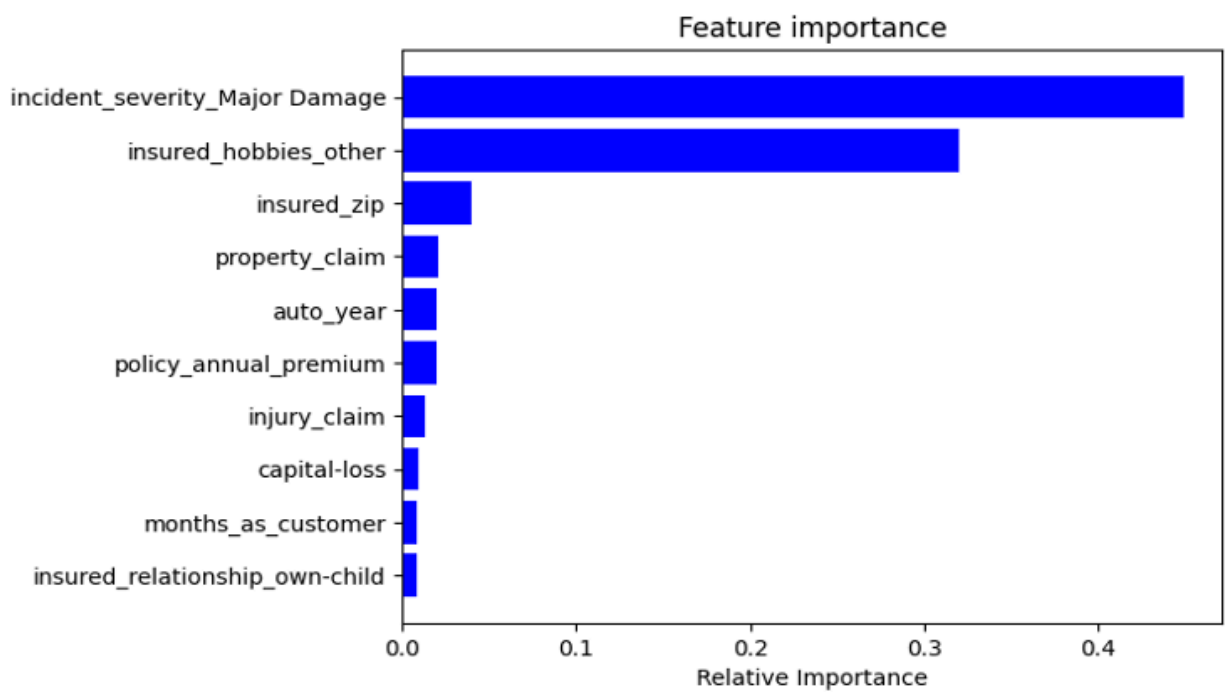
The Random Forest model showed a good balance between precision and recall, indicating its ability to effectively identify fraudulent claims while minimizing false positives. The feature importance provided by Random Forest highlighted that the most significant features were the claim amount, policy details, and claimant's history of claims.

## Classification Report

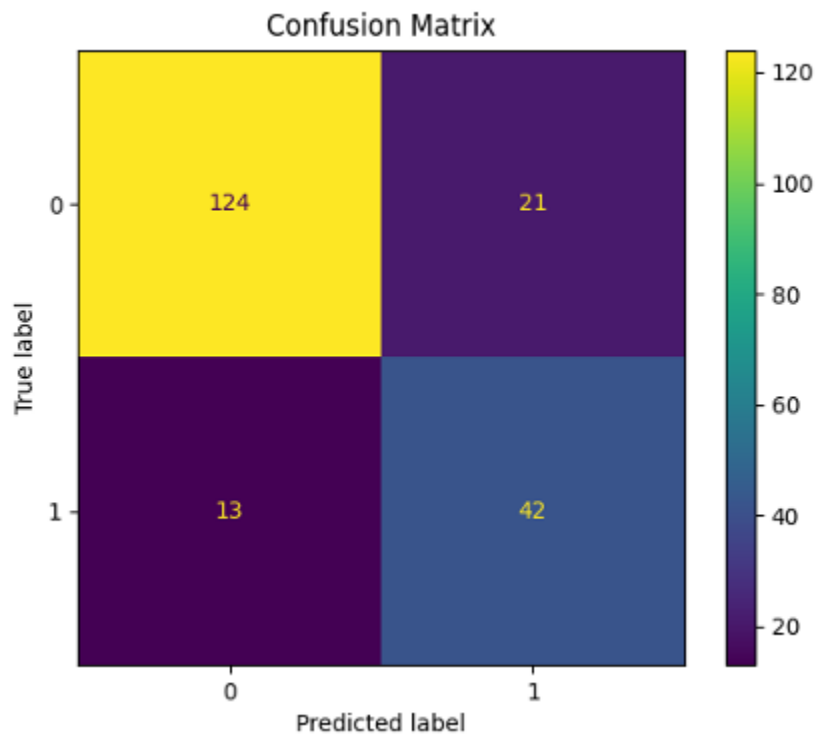
Classification Report

	precision	recall	f1-score	support
0	0.91	0.86	0.88	145.00
1	0.67	0.76	0.71	55.00
accuracy	0.83	0.83	0.83	0.83
macro avg	0.79	0.81	0.80	200.00
weighted avg	0.84	0.83	0.83	200.00

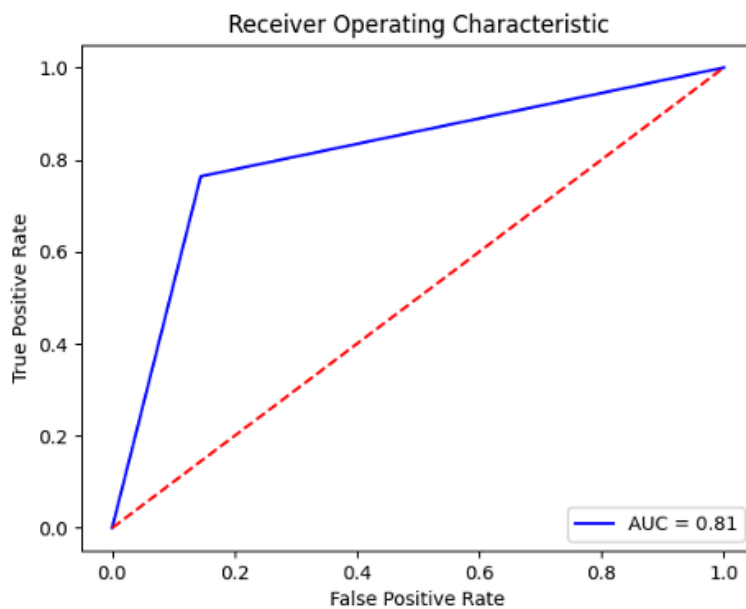
## Feature Importance



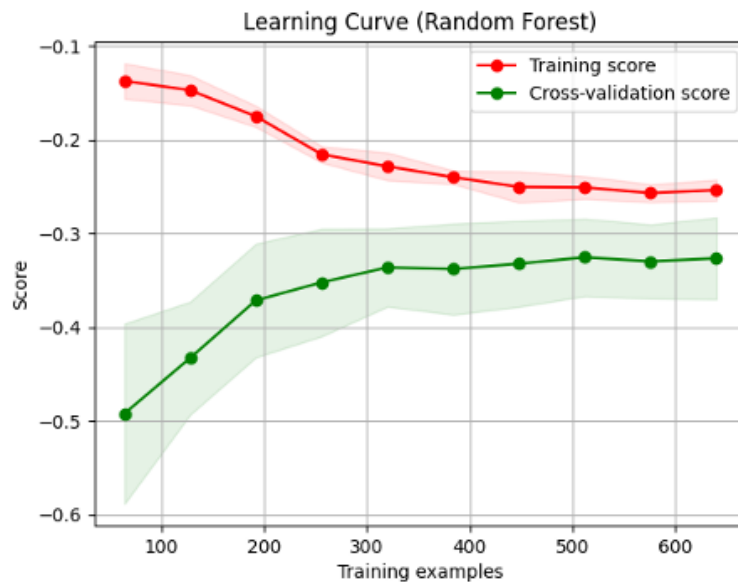
## Confusion Matrix



## Receiver Operating Characteristic (ROC)



## Learning Curve



## XGBoost

- *Accuracy*: 83.5%
- *Precision*: Positive class - 65%, Negative class - 94%
- *Recall*: Positive class - 87%, Negative class - 82%
- *F1 Score*: Positive class - 74%, Negative class - 88%
- *ROC-AUC*: 0.85

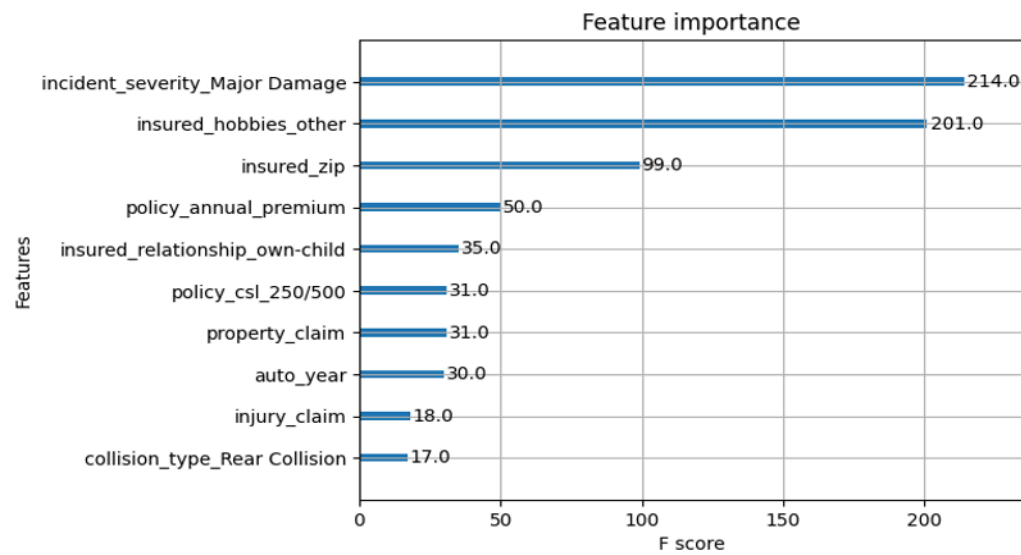
The XGBoost model outperformed the Random Forest model in terms of accuracy, precision, recall, and F1 score. Its ability to handle missing values and incorporate regularization contributed to its superior performance. The feature importance analysis from XGBoost also pointed out similar key features as Random Forest but with more refined granularity.

## Classification Report

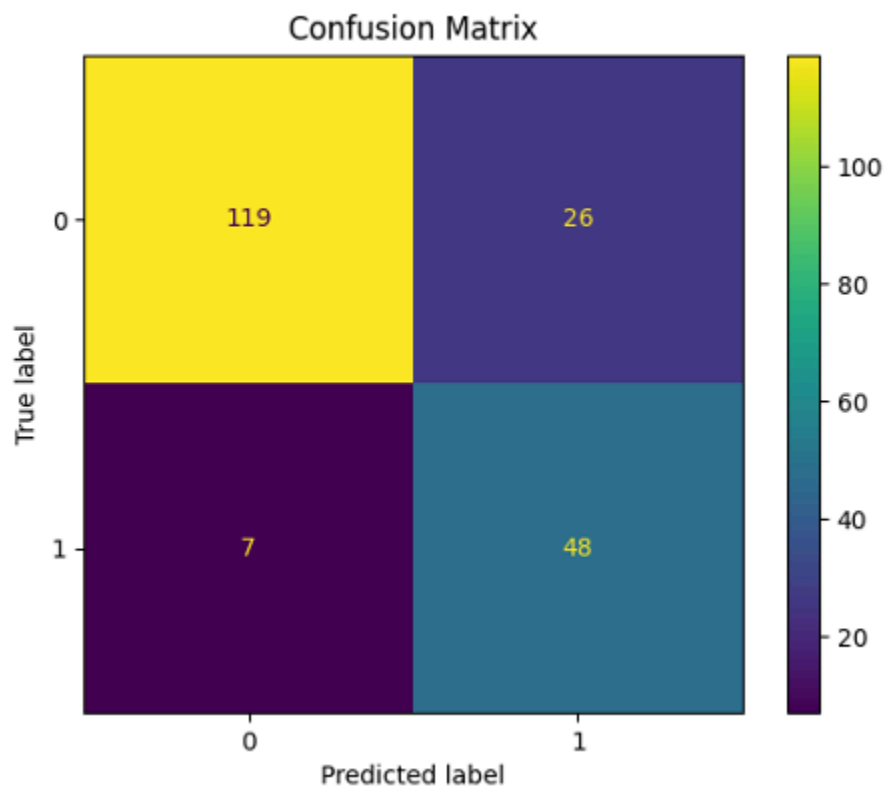
Classification Report

	precision	recall	f1-score	support
0	0.94	0.82	0.88	145.00
1	0.65	0.87	0.74	55.00
accuracy	0.83	0.83	0.83	0.83
macro avg	0.80	0.85	0.81	200.00
weighted avg	0.86	0.83	0.84	200.00

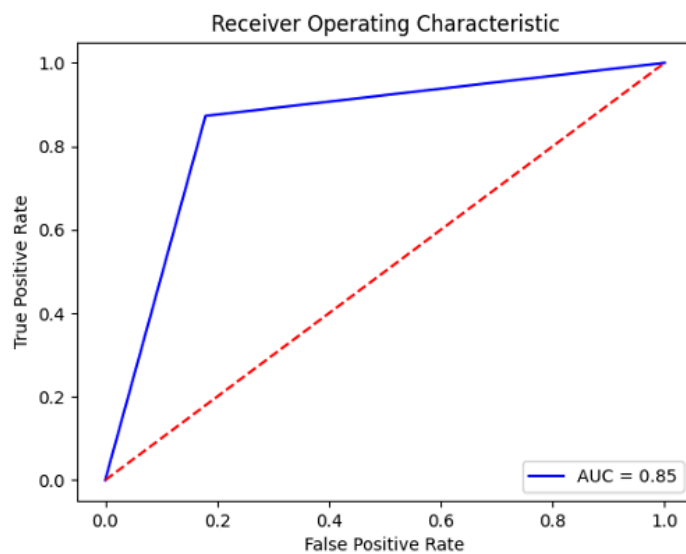
## Feature Importance



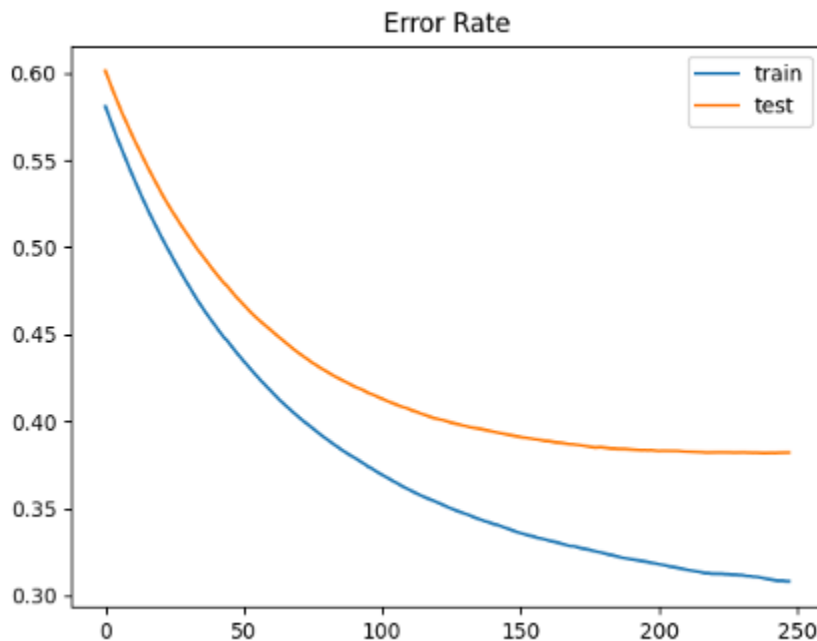
## Confusion Matrix



## Receiver Operating Characteristic (ROC)



## Error Rate



## Comparative Analysis

- **Precision-Recall Trade-offs:** XGBoost showed a higher precision and recall balance compared to Random Forest, indicating better performance in detecting fraudulent claims while minimizing false positives.
- **Feature Importance:** Both models highlighted similar key features influencing fraud detection, such as claim amount, policy details, and claimant history. However, XGBoost provided more detailed insights into feature interactions and their impact on predictions.
- **Model Complexity:** While Random Forest is easier to implement and interpret, XGBoost's complexity requires more careful tuning but offers superior performance when optimized correctly.



# Conclusion

The application of machine learning in insurance fraud detection demonstrates significant potential in enhancing the accuracy and efficiency of identifying fraudulent claims. XGBoost outperformed Random Forest in terms of accuracy and overall robustness. The results indicate that machine learning models, particularly XGBoost, can effectively identify patterns and anomalies associated with fraudulent claims, providing a valuable tool for insurance companies to mitigate risks.

Future work could explore the integration of additional data sources, such as social media data or other external databases, to further improve the detection accuracy. Additionally, implementing real-time detection systems and exploring advanced techniques like deep learning and reinforcement learning could offer further improvements in fraud detection capabilities. Collaboration with domain experts in insurance can also help in refining feature engineering and model interpretation, ensuring that the models are aligned with practical needs and constraints.

# References

1. Ngai, E. W. T., Hu, Y., Wong, Y. H., Chen, Y., & Sun, X. (2011). The application of data mining techniques in financial fraud detection: A classification framework and an academic review of literature. *Decision Support Systems*, 50(3), 559-569.
2. Phua, C., Lee, V., Smith, K., & Gayler, R. (2010). A comprehensive survey of data mining-based fraud detection research. *arXiv preprint arXiv:1009.6119*.
3. Bauder, R. A., & Khoshgoftaar, T. M. (2018). A survey of data mining and machine learning methods for cyber security intrusion detection. *IEEE Access*, 6, 19201-19223.
4. Sahoo, S., Singh, A. K., & Mukhopadhyay, S. (2018). A survey on different approaches for detecting fraud in insurance. *Journal of King Saud University-Computer and Information Sciences*, 30(4), 419-432.
5. Ghosh, S., & Reilly, D. L. (1994). Credit card fraud detection with a neural-network. In *System Sciences, 1994. Proceedings of the Twenty-Seventh Hawaii International Conference on* (Vol. 3, pp. 621-630). IEEE.