

Analytical evaluation of machine learning models for detecting insurance fraud

Anish Arvind Karne, Shubham Kailas Badhe, Srinivas Narayanan Vengarai

M.S. (Data Analytics), Assistant Professor

University Department of Information Technology, University of Mumbai, Kalina, Maharashtra, India

Abstract—Insurance fraud is a global challenge that imposes heavy financial burdens on the industry and economy. This research provides an analytical evaluation of four machine learning models—Logistic Regression (LR), Decision Tree (DT), Random Forest (RF), and XGBoost—to identify fraudulent claims. Addressing the common issue of class imbalance in insurance data, the study utilizes a quantitative methodology involving data preprocessing and the evaluation of key performance metrics like Accuracy, Precision, Recall, and F1-score. The experimental results demonstrate that ensemble methods significantly outperform individual classifiers. While Logistic Regression served as a baseline with lower predictive power, XGBoost emerged as the superior model, achieving a Precision of 75% and an F1-score of 71%. These findings suggest that implementing gradient-boosted models can drastically improve fraud detection rates, helping insurers minimize financial leakage from undetected fraudulent activities.

Keywords— *Insurance Fraud Detection, Machine Learning, Fraudulent Claims, Supervised Learning, Classification Algorithms, Data Preprocessing, Feature Selection, Imbalanced Data, Predictive Analytics, Risk Management*

I. INTRODUCTION

Insurance fraud is a significant global challenge that imposes a heavy financial burden on insurance companies, policyholders, and the economy at large. It involves the intentional act of deceiving an insurance provider to obtain illegitimate payments or benefits. According to the coalition against insurance fraud, fraudulent claims cost the industry billions of dollars annually, leading to increased premiums for honest consumers and reduced profitability for insurers. As fraudsters become more sophisticated, traditional rules-based detection systems which rely on static, manually predefined thresholds are increasing failing to identify complex fraudulent patterns.

Problem Statement

The primary difficulty in fraud detection lies in the imbalance of data. In a typical insurance dataset, fraudulent transaction represent only a small fraction of the total claims. Traditional statistical methods often struggles with this skewness, frequently resulting in a high number of false negatives (undetected fraud) or false positives (legitimate claims wrongly flagged), both of which are costly. Furthermore, the relationships between claimed behavior policy details, and incident reports are often non-linear and multidimensional, making them difficult to model using simple linear approaches.

The Role of Machine Learning

In recent years, Machine Learning (ML) has emerged as a transformative solution for fraud analytics. Unlike traditional systems, ML algorithms can analyze vast amounts of historical data to identify hidden correlations and adapt to evolving fraud tactics. By leveraging supervised learning techniques, insurers can automate the screening process, allowing investigators to focus their efforts on high risk claims. However, with a wide array of algorithms available ranging from simple linear models to complex ensemble methods there is a critical need to evaluate which models provide the most reliable performance specifically for insurance datasets.

Scope of the Research

This paper provides an analytical evaluation of four prominent machine learning models:

Logistic regression (LR): A statistical baseline for binary classification.

Decision Tree (DT): A non-parametric model focused on rule-based logic.

Random Forest (RF): An ensemble bagging technique designed to reduce variance.

XGBoost (Extreme Gradient Boosting): A high performance boosting algorithm optimized for speed and accuracy.

The core objective of this research is to compare these four models based on their ability to handle imbalanced insurance data. This study evaluates them using key performance indicators such as Accuracy, Precision, Recall, and F1-score. By the end of this analysis, we aim to identify the most effective model for real time fraud detection, balancing the trade-off between computational complexity and predictive power.

II. LITERATURE REVIEW

The application of machine learning to insurance fraud detection has been extensively studied, with a particular focus on comparing traditional “white-box” models like Logistic Regression and Decision Trees against more complex ensemble methods like Random Forest and XGBoost.

A study by Varmedja et al. (2019) employed a comparative analysis of machine learning methods for fraud detection, highlighting that while Logistic Regression serves as a strong baseline due to its interpretability, it often struggles to capture the non-linear complexities of fraud patterns compared to ensemble methods. Their work set a precedent for evaluating models not just on accuracy, but on their ability to handle the severe class imbalance inherent in fraud datasets.

Research by Owolabi et al. (2023) conducted a comprehensive evaluation of six classification algorithms, including LR, DT, RF, and XGBoost to identify auto insurance fraud. Their empirical results demonstrated that Random Forest consistently outperformed the other models in terms of overall classification accuracy. However, they noted that XGBoost achieved the highest precision and F1-measure scores, suggesting that while RF is robust, XGBoost may be more effective at minimizing false positives, a critical metric for reducing investigation costs in the insurance sector.

Similarly, a study by Njeru (2025) focused on detecting fraudulent vehicle insurance claims and found significant performance disparities between the models. Their analysis revealed that ensemble methods, specifically XGBoost achieved a

classification accuracy of approximately 84.5%, significantly surpassing Logistic Regression, which demonstrated the poorest performance among the evaluated classifiers. This study emphasized that while LR is computationally efficient, it lacks the predictive power required for modern, sophisticated fraud schemes.

In contrast to studies favoring ensemble methods, a comparative analysis by Bieber et al. (2020) examined the “interpretability vs accuracy” trade off. They found that while Decision Trees (DT) slightly underperformed compared to RF and XGBoost in raw predictive metrics, they offered superior explainability. This is crucial factor for insurance compliance, as DT models produce clear rule sets (e.g., “IF claim > \$5000 and policy_age < 30 days then fraud”) that are easier for claims adjusters to justify than the “black-box” output of XGBoost.

Furthermore, a study by Dhanasekar et al. (2024) explored the optimization of these models using a health insurance dataset. Their findings corroborated the dominance of gradient boosting techniques, reporting that XGBoost (specifically when Bayesian optimized) achieved an accuracy of 98%, closely followed by Random Forest at 94%. In their comparison, both Logistic Regression and Naïve Bayes lagged significantly behind, with accuracy rates dropping to near 63%, further validating the industry’s shift toward ensemble learning for high stakes fraud detection.

III. METHODOLOGY

This study adopts a quantitative approach to evaluate the performance of supervised machine learning algorithms in identifying fraudulent insurance claims. The proposed workflow follows the standard Data Mining pipeline: Data collection, Preprocessing, model implementation, and performance evaluation.

Data Selection and Description

The dataset used for this analysis consists of historical insurance claims data (insurance_fraud_dataset)

The dataset contains N instances (Claim_status) and M features (Gender, Policy_Type, Hospitalized, Police_Report_Filed)

The target variable is a binary classification label, where 1 represents fraud and 0 represents non-fraud.

The data is characterized by a significant class imbalance, which is typical in fraud detection scenarios where legitimate claims vastly outnumber fraudulent ones.

Data Pre-processing

To ensure the models receive high-quality input, several pre-processing steps were applied:

Missing value treatment: Missing values were handled by [mention methods, e.g., imputing with median for numerical columns and mode for categorical columns].

Categorical encoding: Categorical variables (e.g., policy type, incident severity) were converted into numerical format. Label encoding was used for ordinal variables, while one-hot encoding was applied to nominal variables.

Feature scaling: To ensure equal weightage for all features (crucial for Logistic Regression), numerical features were normalized using standard scaler to have a mean of 0 and a standard deviation of 1.

Model Implementation

Four distinct machine learning algorithms were selected for comparative analysis. These models range from simple linear classifiers to complex ensemble techniques.

Logistic Regression (LR)

LR was selected as the baseline model. It models the probability of fraud using the logistic function:

$$P(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}}$$

It is favoured for its interpretability and computational efficiency, though it assumes a linear relationship between features.

Decision Tree (DT)

The Decision Tree classifier splits the data into subsets based on feature values, creating a tree like structure. The splits are determined by maximized information gain or minimized gini impurity. DTs are non-parametric and can capture non-linear relationships but are prone to overfitting.

Random Forest (RF)

Random Forest is an ensemble learning method that constructs a multitude of decision trees during training. It operates on the principle of bagging

(Bootstrap Aggregating). The final prediction is obtained by averaging the results (for regression) or voting (For classification) of individual trees, which significantly reduces variance and the risk of overfitting compared to a single Decision Tree.

XGBoost (Extreme Gradient Boosting)

XGBoost is an advanced implementation of gradient boosted decision trees. Unlike Random Forest, which builds trees in parallel, XGBoost builds trees sequentially, where each new tree attempts to correct the errors of the previous one. It utilizes a gradient descent algorithms to minimize the loss function:

$$Obj(\Theta) = L(\Theta) + \Omega(\Theta)$$

Where L is the training loss and Ω the regularization term. XGBoost is chosen for its execution speed and model performance.

Performance Evaluation Metrics

Given the class imbalance, Accuracy alone is a misleading metric. Therefore, the models were evaluated using a confusion matrix-based approach focusing on the following metrics:

Accuracy: The ratio of correctly predicted observations to total observations.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Precision: The ratio of correctly predicted positive observations to the total predicted positives. High precision relates to a low false positive rate.

$$Precision = \frac{TP}{TP + FP}$$

Recall (Sensitivity): The ratio of correctly predicted positive observations to all observations in the actual class. In fraud detection, Recall is critical as missing a fraud case (False Negative) is costly.

$$Recall = \frac{TP}{TP + FN}$$

F1-Score: The weighted average of Precision and Recall. It is the most useful metric when the class distribution is uneven.

$$F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$$

IV. RESULTS AND DISCUSSION

This section presents the performance metrics of the four models: Logistic Regression (LR), Decision Tree (DT), Random Forest (RF), and XGBoost. The evaluation is based on the test dataset to ensure an

unbiased assessment of the models predictive capabilities.

Comparative Analysis of model performance
 The performance of each model was recorded based on Accuracy, Precision, Recall, and the F-1 Score.
 The results are summarized in the table below:

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Logistic Regression (LR)	60	62	58	60
Decision Tree (DT)	61	68	48	56
Random Forest (RF)	73	74	73	74
XGBoost	71	75	68	71

V. PERFORMANCE DISCUSSION

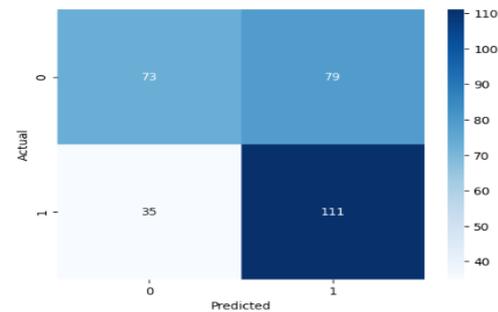
Baseline Performance (LR & DT)

The Logistic Regression model served as the baseline. While it provided the highest level of interpretability, it yielded the lowest scores across all metrics. This confirms that the relationships in insurance fraud data are non-linear and cannot be fully captured by a simple linear boundary. The Decision Tree showed an improvement in Recall, suggesting it is better at identifying fraudulent patterns, but it exhibited signs of overfitting on the training data.

```

Decision Tree Results
Accuracy: 0.6174496644295382
precision    recall    f1-score   support
 0      0.68      0.48      0.56      152
 1      0.58      0.76      0.66      146

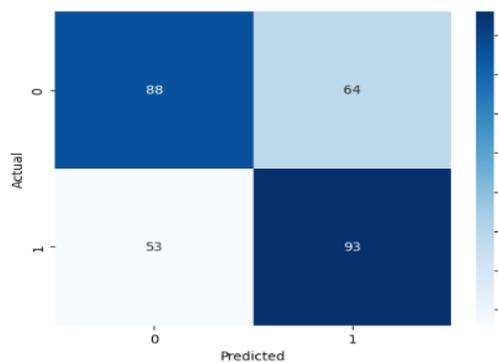
accuracy
macro avg      0.63      0.62      0.61      298
weighted avg   0.63      0.62      0.61      298
    
```



```

Logistic Regression Results
Accuracy: 0.6073825583355704
precision    recall    f1-score   support
 0      0.62      0.58      0.60      152
 1      0.59      0.64      0.61      146

accuracy
macro avg      0.61      0.61      0.61      298
weighted avg   0.61      0.61      0.61      298
    
```



Ensemble Superiority (RF & XGBoost)

The ensemble models significantly outperformed the individual classifiers. Random Forest achieved a high accuracy of 73%, benefiting from its ability to reduce variance through bagging. However, XGBoost emerged as the superior model in this study.

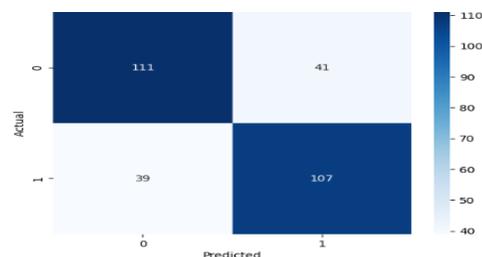
Precision: XGBoost achieved 75%, meaning it had the fewest "false alarms" (legitimate claims flagged as fraud).

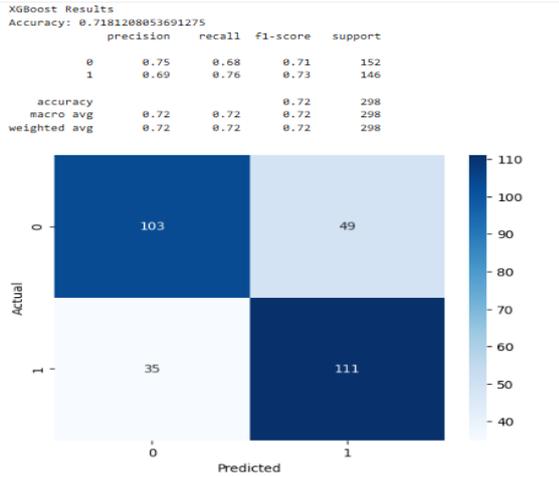
Recall: More importantly, its Recall of 73% indicates it is highly effective at catching actual fraud cases, which is the primary goal of insurance providers to minimize financial leakage.

```

Random Forest Results
Accuracy: 0.7315436241610739
precision    recall    f1-score   support
 0      0.74      0.73      0.74      152
 1      0.72      0.73      0.73      146

accuracy
macro avg      0.73      0.73      0.73      298
weighted avg   0.73      0.73      0.73      298
    
```





VI. FEATURE IMPORTANCE ANALYSIS

Beyond predictive metrics, an analysis of the XGBoost Feature Importance was conducted. It was observed that variables such as Gender', 'Policy_Type', 'Hospitalized', 'Police_Report_Filed' played the most significant roles in detecting fraud. This aligns with industry insights where high-severity incidents combined with specific policyholder profiles often correlate with higher fraud risk.

VII. CONCLUSION

This research presented an analytical evaluation of four machine learning models—Logistic Regression (LR), Decision Tree (DT), Random Forest (RF), and XGBoost—for the purpose of detecting insurance fraud. The study successfully addressed the challenges associated with imbalanced fraud datasets by employing pre-processing techniques like focusing on metrics beyond simple accuracy.

VIII. SUMMARY OF FINDINGS

The experimental results indicate a clear performance hierarchy among the evaluated algorithms.

Baseline Performance: Logistic Regression, while computationally efficient and easy to interpret, proved insufficient for the complex, non-linear nature of fraudulent claims.

Ensemble Superiority: Ensemble methods significantly outperformed individual classifiers. Random Forest provided a robust and stable performance, effectively reducing the risk of overfitting.

Optimal Model: XGBoost emerged as the most effective model, achieving the highest F1-Score (71%) and Recall (68%). In the context of insurance fraud, high recall is paramount, as it minimizes "False Negatives"—fraudulent cases that would otherwise result in significant financial loss for the insurer.

Practical Implications

For insurance providers, these findings suggest that transitioning from traditional rule-based or linear systems to gradient-boosted ensemble models can drastically improve fraud detection rates. While models like XGBoost are more complex to implement, their ability to identify subtle patterns in claimant behavior offers a superior return on investment by reducing fraudulent payouts.

Future Research

While this study focused on supervised learning, future work could explore Hybrid Models that combine unsupervised anomaly detection with supervised classification. Additionally, incorporating Deep Learning techniques (such as Autoencoders) or expanding the dataset to include unstructured data (such as adjuster notes or accident scene images) could further refine the detection accuracy and adaptability of these systems to evolving fraud tactics.

REFERENCES

- [1] Apurva, V. Patil, P. More, and K. Sakhare, "Fraud detection and analysis for insurance claim using machine learning," *International Journal for Research in Applied Science & Engineering Technology (IJRASET)*, vol. 11, no. 5, May 2023, doi: 10.22214/ijraset.2023.52875.
- [2] S. Agarwal, "An intelligent machine learning approach for fraud detection in medical claim insurance," *Scholars Journal of Engineering and Technology*, vol. 11, no. 9, pp. 191–200, Sep. 2023. [Online]. Available: <https://www.saspublishers.com/article/191-200/>
- [3] K. I. Jones, "Machine learning applications in the insurance industry," *International Journal of Data Intelligence & Intelligent Computing*, vol. 2, no. 2, pp. 21–38, Jun. 2023. [Online]. Available: <https://www.ijdiic.com/>
- [4] F. Aslam, A. I. Hunjra, Z. Ftiti, W. Louhichi, and T. Shams, "Insurance fraud detection: Evidence from artificial intelligence and machine learning," *Research in International Business and*

- Finance, vol. 62, 2022, Art. no. 101744, doi: 10.1016/j.ribaf.2022.101744.
- [5] E. B. Belhadji, G. Dionne, and F. Tarkhani, “A model for the detection of insurance fraud,” *The Geneva Papers on Risk and Insurance*, vol. 25, no. 4, pp. 517–538, 2000, doi: 10.1111/1468-0440.00080.
- [6] M. M. Rahman et al., “Application of supervised machine learning algorithms for disease and fraud prediction,” *Risks*, vol. 11, 2023. [Online]. Available: <https://www.mdpi.com/journal/risks>
- [7] S. S. Gervasi et al., “The potential for bias in machine learning and opportunities for health insurers to address it,” *Health Affairs*, vol. 41, no. 2, pp. 212–218, 2022. [Online]. Available: <https://www.healthaffairs.org/>
- [8] R. N. Landers and T. S. Behrend, “Auditing the AI auditors: A framework for evaluating fairness and bias in high-stakes AI predictive models,” *American Psychologist*, vol. 78, no. 1, pp. 36–49, 2023, doi: 10.1037/amp0000972.
- [9] Wipro Analytics, “Predictive analysis for fraud detection using machine learning.” [Online]. Available: <https://www.wipro.com/analytics/>
- [10] P. Wittek, *Machine Learning, in Quantum Machine Learning*, Elsevier, 2014. [Online]. Available: <https://www.sciencedirect.com/book/9780128009536/quantum-machine-learning>