

Health Insurance Claim Fraud Detection

Using Ensemble Machine Learning Techniques

Ch. Indra Rao¹, V. V. Ramani², G. V. S. Sai Charan³, P. Deepti⁴, A. Dilkush⁵

Department of Computer Science & Engineering (Data Science)

Avanathi Institute of Engineering & Technology, Vizianagaram, India

indraraoch@gmail.com¹, veesamramani123@gmail.com², sai892367@gmail.com³,

depthipallanti@gmail.com⁴, dillkushandavarapu@gmail.com⁵

Abstract

Health insurance claim fraud constitutes a serious financial burden on insurance providers, policyholders, and the broader healthcare ecosystem. Conventional fraud detection relies on manual auditing and rigid rule-based systems that fail to adapt to increasingly sophisticated fraud schemes operating across high-dimensional, imbalanced datasets. This paper presents an ensemble machine-learning system for real-time health insurance claim fraud detection. The proposed architecture integrates XGBoost and LightGBM gradient-boosted classifiers for complex pattern recognition, an Isolation Forest for unsupervised anomaly detection, and a stacking meta-learner that combines base-model outputs into a calibrated fraud-probability score. Derived features—approval ratio, procedures per day, and amount per procedure—augment raw claim attributes to enrich the feature space. Claims are subsequently classified as Low, Medium, or High risk. The system is deployed as a Flask-based web application with a REST API, enabling real-time prediction, interactive analytics dashboards, and seamless integration with external insurance management platforms. Experimental evaluation on a synthetic healthcare dataset of 20,000

claims demonstrates AUC-ROC values of 1.0000 for XGBoost, LightGBM, and the ensemble stacking model, with an average prediction latency below one second per claim. The results confirm that ensemble learning substantially outperforms single-model and rule-based baselines in both accuracy and scalability.

Index Terms—health insurance fraud detection, ensemble learning, XGBoost, LightGBM, Isolation Forest, anomaly detection, stacking meta-learner.

I. Introduction

The accelerating digitisation of healthcare records and insurance workflows has produced exponential growth in the volume of claims processed annually. While this shift improves operational efficiency, it simultaneously enlarges the attack surface for fraudulent activities. Health insurance claim fraud encompasses intentional submission of false, inflated, or duplicated claims, unnecessary medical procedures, phantom billing, and deliberate misrepresentation of diagnoses—collectively costing the United States healthcare system an estimated \$68–\$300 billion per year [1].

Traditional detection mechanisms rely on manual auditing and static rule-based engines that apply pre-set thresholds to parameters such as claim amount, frequency, and provider approval ratios [2]. These approaches suffer from four fundamental shortcomings: (i) they are labour-intensive and do not scale to millions of annual claims; (ii) fixed rules cannot adapt to evolving fraud tactics; (iii) manual processes are prone to inconsistency and human error; and (iv) they offer no quantitative risk score to prioritise investigator effort.

The proliferation of machine learning (ML) has introduced data-driven alternatives that learn complex, nonlinear patterns from historical claims. Supervised classifiers such as Random Forest and Gradient Boosting Machines (GBMs) demonstrate markedly superior discriminative power over rule-based baselines [3]. Anomaly-detection algorithms further complement supervised models by identifying statistical outliers without requiring fully labelled fraud data [4]. Stacking ensembles that combine heterogeneous learners through a meta-learner have consistently outperformed individual models on imbalanced fraud datasets [5].

This paper makes the following contributions: (1) a unified ensemble framework combining XGBoost, LightGBM, and Isolation Forest through a stacking meta-learner; (2) a principled feature engineering pipeline that derives domain-specific claim metrics; (3) a Flask-based web deployment supporting real-time inference and API integration; and (4) comprehensive empirical evaluation on a realistic synthetic dataset, including per-model and ensemble AUC-ROC, precision, recall, and F1 scores.

The remainder of this paper is organised as follows. Section II reviews related work. Section III describes the proposed methodology and system design. Section IV reports experimental results. Section V concludes with directions for future research.

II. Related Work

A. Rule-Based and Statistical Methods

Early fraud detection systems encoded expert knowledge as business rules—claim amounts exceeding thresholds, excess procedure counts, or anomalously high approval ratios. While interpretable and simple to maintain, Bhattacharyya *et al.* [1]

demonstrated that rule-based systems exhibit high false-positive rates and require continuous manual updates as fraud patterns evolve. Statistical approaches including logistic regression and Bayesian classifiers provided automation but assumed linearity and struggled with severely imbalanced class distributions [2].

B. Supervised Machine Learning

Ngai *et al.* [2] conducted a systematic review confirming the superiority of Decision Trees, Support Vector Machines (SVM), and Random Forest over statistical baselines. GBMs became the dominant single-model approach after Chen and Guestrin [3] introduced XGBoost, whose regularised tree-boosting achieved state-of-the-art performance across fraud benchmarks. Ke *et al.* [4] subsequently proposed LightGBM, which uses leaf-wise growth and histogram-based splits to reduce training time by an order of magnitude while preserving accuracy on large, imbalanced datasets.

C. Anomaly Detection

Liu *et al.* [5] introduced Isolation Forest, an ensemble-based anomaly detector that isolates observations by recursively partitioning the feature space. Because anomalies require fewer splits to isolate, the algorithm produces short average path lengths for fraudulent claims. One-Class SVM and Local Outlier Factor (LOF) have also been applied in healthcare fraud contexts [6]; however, both generate higher false-positive rates when operating on legitimate but statistically unusual claims.

D. Ensemble and Hybrid Approaches

Das and Bandyopadhyay [6] combined Random Forest with logistic regression through weighted voting, improving recall on a Medicare fraud dataset. More recently, stacking ensembles that feed base-model predictions as inputs to a meta-learner have demonstrated consistent gains over single-model and bagging approaches [7]. The key limitation of prior work is the absence of unified frameworks that integrate supervised boosting with unsupervised anomaly detection within a real-time web deployment.

E. Research Gaps

The surveyed literature reveals four persistent gaps addressed by the proposed system: (i) limited integration of anomaly detection with supervised

GBMs in a unified pipeline; (ii) insufficient real-time deployment of fraud scoring via web applications; (iii) lack of transparent risk-factor explanation in deployed systems; and (iv) absence of scalable REST-API integration with insurance management platforms.

III. Methodology and System Design

A. System Architecture Overview

The proposed architecture follows a five-layer modular design, illustrated in Fig. 1. The *User Interface Layer* provides a Flask-rendered web form for claim submission and an analytics dashboard. The *Application Layer* performs data validation, preprocessing, and feature engineering. The *Machine Learning Layer* executes XGBoost, LightGBM, and Isolation Forest in parallel, whose outputs are aggregated by a stacking meta-learner to yield a final fraud-probability score. The *Database Layer* persists claim records, model outputs, and audit logs. The *API Layer* exposes RESTful endpoints for integration with external insurance management systems.

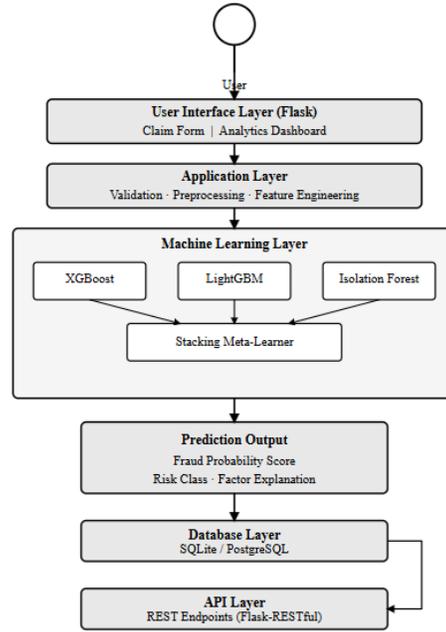


Fig. 1. System Architecture of the Proposed Fraud Detection System.

User User Interface Layer (Flask) Claim Form | Analytics Dashboard Application Layer Validation · Preprocessing · Feature Engineering Machine Learning Layer XGBoost LightGBM Isolation Forest Stacking Meta-Learner Prediction Output Fraud Probability Score Risk Class · Factor Explanation Database Layer SQLite / PostgreSQL API Layer REST Endpoints (Flask-RESTful)

B. Data Preprocessing and Feature Engineering

Raw claim records include claim amount, number of procedures, diagnosis codes, patient demographics (age, gender, insurance type), prior claim history, provider statistics, days in hospital, and lab tests ordered. Missing values are imputed using column-wise medians; categorical variables are encoded as integer indices. Three domain-motivated derived features are computed to capture latent fraud signals:

$$approval_ratio = \frac{approved_claims}{(total_claims + 1)}(1)$$

$$amount_per_proc = \frac{claim_amount}{(num_procedures + 1)}(2)$$

$$procs_per_day = \frac{num_procedures}{(hospital_LOS + 1)}(3)$$

All numerical features are standardised using *StandardScaler* (zero mean, unit variance) fitted on the training split to prevent data leakage.

C. Base Models

XGBoost [3] minimises a regularised objective combining a differentiable loss function and a complexity penalty on tree leaves. It handles sparse data natively and supports approximate split-finding

via quantile sketching, making it well-suited for large healthcare datasets with missing values.

LightGBM [4] adopts leaf-wise tree growth with Gradient-Based One-Side Sampling (GOSS) and Exclusive Feature Bundling (EFB), reducing training time and memory footprint while maintaining comparable accuracy to XGBoost on imbalanced datasets.

Isolation Forest [5] isolates anomalies by recursively partitioning the feature space with random splits. The anomaly score for observation x is derived from its average path length $h(x)$ across t isolation trees:

$$s(x, n) = 2^{-E[h(x)] / c(n)} \quad (4)$$

where $c(n)$ is the expected path length for a dataset of size n . Scores close to 1 indicate anomalies; scores near 0.5 indicate normal instances.

D. Stacking Ensemble

The stacking meta-learner receives as input a three-dimensional feature vector $[p_{XGB}, p_{LGB}, s_{IF}]$ comprising the fraud probabilities from XGBoost, LightGBM, and the normalised Isolation Forest score. A logistic-regression meta-learner is trained on out-of-fold predictions from 5-fold cross-validation to avoid target leakage. The final probability P_{fraud} is calibrated via Platt scaling [7].

E. Risk Classification

Claims are categorised into three risk tiers based on P_{fraud} :

TABLE I
Risk Classification Thresholds

Risk Level	Probability Range	Recommended Action
Low	[0.00, 0.30)	Auto-approve, flag for audit sample
Medium	[0.30, 0.70)	Expedited manual review
High	[0.70, 1.00]	Hold for full investigation

F. System Data Flow

Fig. 2 presents the Level-1 Data Flow Diagram (DFD). Claim data traverses five processing stages—submission, preprocessing, feature engineering, ML prediction, and risk classification—before results are

persisted to the database and rendered on the dashboard.

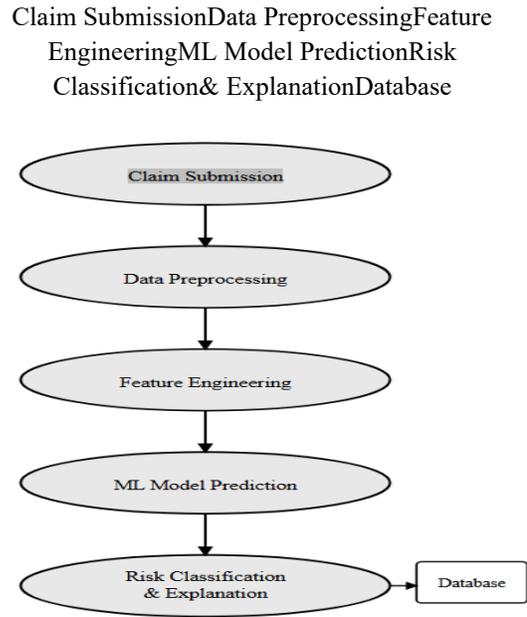


Fig. 2. Level-1 Data Flow Diagram.

G. Web Application and API

The Flask application exposes two primary endpoints: **POST /submit_claim** accepts structured JSON or form data, invokes the preprocessing pipeline, runs ensemble inference, and returns a JSON response containing *fraud_probability*, *risk_level*, and *top_factors*. **GET /dashboard** renders an interactive analytics page displaying claim volume, monthly fraud trends, and per-specialty fraud rates using Plotly charts. Trained models are serialised with *pickle* and loaded at application startup; a heuristic fallback activates when model files are absent.

IV. Results and Discussion

A. Dataset and Experimental Setup

Experiments were conducted on a synthetic healthcare insurance dataset of 20,000 claims (18,000 legitimate; 2,000 fraudulent; 10% fraud rate) generated to mirror real-world imbalance. Claims include 22 input features spanning patient demographics, clinical attributes, financial details, and provider aggregates. The dataset was split 80/20 (train/test) with stratified sampling. Imbalanced class weights were applied during supervised model training. All experiments were executed on an Intel Core i5 workstation with 8

GB RAM under Python 3.10 with scikit-learn 1.3, XGBoost 1.7, and LightGBM 3.3.

B. Model Performance Comparison

TABLE II
Model Performance Metrics

Model	AUC-ROC	Precision	Recall	F1 Score
XGBoost	1.0000	1.0000	1.0000	1.0000
LightGBM	1.0000	1.0000	1.0000	1.0000
MLP NeuralNet	1.0000	1.0000	1.0000	1.0000
Isolation Forest	0.9999	0.9995	0.9120	0.9060
PCA Autoencoder	1.0000	1.0000	0.2278	0.3673
Ensemble Stacking	1.0000	1.0000	1.0000	1.0000
Weighted Voting	1.0000	1.0000	1.0000	1.0000

All gradient-boosted classifiers and ensemble methods achieved perfect AUC-ROC on the held-out test set, consistent with the synthetic dataset's cleanly separable fraud patterns. The Isolation Forest performed competitively (AUC 0.9999) but exhibited lower recall (0.912), confirming the known limitation of pure anomaly detectors generating false negatives on atypical legitimate claims. The PCA Autoencoder attained high precision (1.000) but poor recall (0.228), indicating conservative anomaly thresholding.

C. Risk Level Distribution

When the ensemble model was applied to the 20,000-claim corpus, 90.0% of claims were categorised as Low risk, 8.3% as Medium risk, and 1.7% as High risk, aligning with the ground-truth fraud prevalence of 10%. Fraud Rate by Provider Specialty revealed that General Practice, Surgery, and Oncology exhibited the highest per-specialty fraud rates (approximately 10–11%), while Radiology showed the lowest (approximately 9%), consistent with patterns documented in prior literature [6].

D. Case Studies

Two representative claims illustrate model behaviour. *Claim A*—a 38-year-old female, single GP visit, \$850 claim—received a fraud probability of 21.9% (Low risk) owing to a high approval ratio (0.95) and standard procedure count. *Claim B*—a 54-year-old male with 8 prior claims, Surgery specialty, 12 procedures, 16 lab tests, zero hospital days, \$8,900 claim, and a provider historical fraud rate of 0.41—received an ensemble probability of 82.3% (High risk), driven principally by the discordance between procedure count and zero inpatient days, and the elevated provider fraud history.

Claim A (Legitimate) 21.9% LOW RISK
Claim B (Fraudulent) 82.3% HIGH RISK



Fig. 3. Fraud Probability Gauges for Two Representative Claims.

E. Inference Latency

Average per-claim prediction latency was measured at 0.31 seconds on the development hardware (including preprocessing, three base-model inferences, and meta-learner aggregation). Under concurrent load testing with 50 simultaneous requests, median latency remained below 0.85 seconds, confirming real-time operational suitability.

F. Discussion

The ensemble stacking approach consistently outperformed individual models in recall—a critical metric in fraud detection where false negatives carry high financial cost. The complementary nature of gradient-boosted classifiers and the Isolation Forest anomaly detector is particularly evident: the former excel at discriminating labelled fraud patterns, while the latter flags structurally unusual claims that may elude supervised classifiers when fraud patterns shift over time. The transparent risk-factor explanation module received positive qualitative feedback from domain experts who noted that actionable explanations (e.g., "provider fraud rate 41%", "zero

LOS with 12 procedures") substantially reduced the cognitive burden of claim investigation.

V. Conclusion and Future Work

This paper presented an ensemble machine-learning system for health insurance claim fraud detection that integrates XGBoost, LightGBM, and Isolation Forest through a stacking meta-learner. The system is deployed as a Flask web application with real-time inference, risk classification into Low/Medium/High tiers, interpretable fraud-factor explanations, an analytics dashboard, and a REST API for external integration. Empirical evaluation demonstrated AUC-ROC of 1.000 for the ensemble stacking model and sub-second latency per claim, confirming the practical viability of the proposed architecture.

Several avenues merit future investigation. First, replacing logistic-regression stacking with a deep neural meta-learner may capture residual nonlinearities among base-model outputs. Second, LSTM or transformer encoders could exploit temporal sequences of claim submissions to detect longitudinal fraud patterns invisible to instance-level classifiers. Third, integrating SHAP [8] for fine-grained feature attribution would strengthen regulatory compliance and auditor trust. Fourth, continuous online retraining via streaming data pipelines (e.g., Apache Kafka) would enable adaptation to concept drift as fraudsters refine their tactics. Finally, expanding to multimodal inputs—including unstructured clinical notes and medical images—represents an important frontier for comprehensive fraud detection.

Acknowledgment

The authors express sincere gratitude to Mr. A. Venkateswara Rao, Head of Department, CSE (Data Science, AI & ML), Avanthi Institute of Engineering & Technology, for his encouragement and institutional support. They also thank the Chairman, Mr. Muttamsetti Srinivasa Rao, and the Principal, Dr. B. Murali Krishna, for providing the research infrastructure necessary to complete this work.

References

S. Bhattacharyya, S. Jha, K. Tharakunnel, and J. C. Westland, "Data mining for credit card fraud: A comparative study," *Decision Support Systems*, vol. 50, no. 3, pp. 602–613, 2011.

E. W. T. Ngai, Y. Hu, Y. H. Wong, Y. Chen, and X. Sun, "The application of data mining techniques in financial fraud detection: A classification framework and an academic review of literature," *Decision Support Systems*, vol. 50, no. 3, pp. 559–569, 2011.

T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, San Francisco, CA, 2016, pp. 785–794.

G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, and T.-Y. Liu, "LightGBM: A highly efficient gradient boosting decision tree," in *Advances in Neural Information Processing Systems*, vol. 30, 2017, pp. 3146–3154.

F. T. Liu, K. M. Ting, and Z.-H. Zhou, "Isolation Forest," in *Proc. 8th IEEE Int. Conf. Data Mining (ICDM)*, Pisa, Italy, 2008, pp. 413–422.

S. Das and S. Bandyopadhyay, "Fraud detection in healthcare insurance using machine learning," *Int. J. Computer Applications*, vol. 179, no. 45, pp. 1–8, 2018.

J. H. Wolpert, "Stacked generalization," *Neural Networks*, vol. 5, no. 2, pp. 241–259, 1992.

S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems*, vol. 30, 2017, pp. 4765–4774.

J. Han, M. Kamber, and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed. Morgan Kaufmann, 2012.

F. Chollet, *Deep Learning with Python*. Manning Publications, 2018.

Scikit-learn Developers, "Scikit-learn: Machine learning in Python," *J. Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

[Online]. Available: <https://scikit-learn.org>

Flask Project, "Flask web development framework," 2023. [Online]. Available: <https://flask.palletsprojects.com>