

# FRAUD DETECTION AND ANALYSIS FOR INSURANCE CLAIM USING MACHINE LEARNING

<sup>1</sup>K JAYA KRISHNA, <sup>2</sup>SHAIK GOUSE LAJAM

<sup>1</sup>Associate Professor, Department of Master of Computer Applications, QIS College of Engineering & Technology, Ongole, Andhra Pradesh, India

<sup>2</sup>PG Scholar, Department of Master of Computer Applications, QIS College of Engineering & Technology, Ongole, Andhra Pradesh, India

# ABSTRACT

Insurance fraud is a significant problem for insurance companies, costing them billions of dollars each year. Traditional methods of fraud detection are often ineffective, as they rely on manual review of claims, which is time-consuming and error-prone. Machine learning (ML) offers a promising approach to fraud detection, as it can automate the process of identifying fraudulent claims. This paper presents a novel ML-based approach to fraud detection in insurance claims. The proposed approach uses a supervised learning algorithm to classify claims as either fraudulent or legitimate. The algorithm is trained on a dataset of historical insurance claims, which includes both fraudulent and legitimate claims. The features used to train the algorithm include a

variety of claim-related variables, such as the type of claim, the amount of the claim, and the policyholder's history. The proposed approach was evaluated on a dataset of realworld insurance claims. The results showed that the approach was able to achieve high accuracy in identifying fraudulent claims. The approach was also able to identify fraudulent claims that were not detected by traditional methods. The proposed approach has the potential to significantly improve the efficiency and effectiveness of fraud detection in insurance claims. The approach can be used to automate the process of identifying fraudulent claims, which can free up insurance investigators to focus on more complex cases. The approach can also be used to identify fraudulent claims that are not detected by traditional methods, which can help to reduce the cost of insurance fraud.

**Index :** insurance, claims, fraud, detections, supervised machine learning , approach

#### I. INTRODUCTION

Insurance fraud is a claim made for getting improper money and not actual amount of money from insurance company or any other underwriter. Motor and insurance area unit two outstanding segments that have seen spurt in fraud.Frauds is classified from a supply or nature purpose of read. Sources is client, negotiator or internal with the latter two being a lot of essential from control framework purpose of reads. Frauds cowl vary of improper activities that a private might commit so as to attain the favorable outcome from an underwriter. Frauds is classified into nature wise, for example, application, inflation, identity, fabrication, contrived, evoked accidents etc. This could vary from staging incident, misrepresenting matters as well as pertinent members and therefore reason behind finally the extent of injury occurred. Probable things might embrace packing up for a state of affairs that beneath wasn't lined the insurance. Misrepresenting the context of an event. This might embrace transferring blames to the incidents wherever the insured set is

accountable, failure to require approved the security measures. Increased impact of the incident .Inflated measure of the loss occurred through the addition of not much relatedlosses or/and attributing inflated price to the increased losses[1][2][3].

## **1.1 PROBLEM STSTEMENT**

The traditional method for the detecting frauds depends on the event of heuristics around fraud indicators. Supported these, the selection on fraud created is said to occur in either of situations like, in certain things the principles are shown if the case should be interrogated for extra examination. In numerous cases, an inventory would be prepared with scores for various indicators of the occurred fraud. The factors for deciding measures and additionally the thresholds are tested statistically and periodically recalibrated. Associate aggregation and then price of the claim would verify necessity of case to be sent for extra examination. The challenge with above strategies is that they deliberately believe on manual mediation which might end in the next restrictions:

1. Inability to perceive the context-specific relationships between the parameters (geography, client section, insurance sales

process) which may not mirror the typical picture.

2. Constrained to control with the restricted set of notable parameters supported the heuristic knowledge – whereas being aware that a number of the opposite attributes might conjointly influence the decisions.

3. Reconstruction of the given model is that the hand operated exercise that need to be conducted sporadically to react dynamic behavior. Also to make sure that the model gives feedback from the examinations. The flexibility to manage this standardization is tougher.

4. Incidence of occurrence of fraud is low generally but 1percent of claims area unit classified.

5. Consultations with business specialists point out that there is not a typical model to determine the model exactly similar to the context

# A. Motivation

Ideally, businesses ought to obtain the responses to prevent fraud from happening or if that is out of the question, to watch it before important damage is finished at intervals the strategy. In most of the companies, fraud is understood entirely once it happens. Measures are then enforced to forestall it from happening over again. At intervals the given time that they can't resist at different time intervals, but Fraud detection is that the most effective suited issue for removing it from the atmosphere and preventing from continuance once more.

## **B.** Significance of the Problem

Knowing a risk is that the beginning in bar, associated intensive assessment offers the lightness that want. This is typically usually performed exploitation varied techniques, like interviews, surveys, focus teams, feedback conducted anonymously, detailed study of record and analysis to spot traffic pumpers, service users, and subscription scam which are different fraudulent case. The association of Certified Fraud Examiners offers a detailed guide to follow. This can be usually alleged to be a preventive methodology, fraud analysis and detection is associate certain consequence of risk associate intensive evaluation. Recognize and classify threats to fraud in knowledge technology and telecommunications sector stereotypically yield the shape of the chances like:

Records showing associate degree inflated rates in calls at associate degree surreal time of day to associate degree uncertain location or far-famed fraud location. Unusual Dialing patterns showing one variety being referred to as additional of times by external numbers than job out.

Increased calls created in an exceedingly day than the minute's allotted per day, that might indicate an account has been hacked or shared

## **C. Major Contribution**

To compare machine learning algorithms: LR, XGB, DT, RF and SVM.

To construct a model that predict transactions could be fraudulent with high accuracy.

To detect if an insurance claim is fraudulent or not.

To analyze the performance of fraud detection algorithm.

#### LITERATURE SURVEY

Insurance fraud is a significant financial burden for insurance companies, leading to increased premiums for honest customers. Machine learning (ML) offers a powerful approach to automate claim analysis and detect fraudulent activity.

Here's a breakdown of key points from the literature:

**Problem:** Fraudulent insurance claims pose a financial threat to insurance companies,

impacting premiums for honest customers [2, 4].

**Machine Learning Applications:** ML algorithms excel at analyzing large datasets to identify patterns and anomalies indicative of fraud [2, 4]. These algorithms can analyze various claim data points, including policyholder details, claim history, and claim specifics, to flag suspicious cases [2].

**Benefits:** Utilizing ML for fraud detection offers several advantages:

Accuracy: ML models can learn from historical data to identify subtle patterns associated with fraudulent claims, potentially leading to higher accuracy than traditional methods [2, 4].

**Efficiency:** Automating fraud detection with ML streamlines the process, reducing investigation time and costs [2, 4].

**Predictive Modeling:** ML can be used to develop predictive models that assign a "fraud probability score" to each claim, allowing for prioritization of investigations [2].

**Common ML Algorithms:** Several machine learning algorithms are employed for fraud detection in insurance claims. Some frequently mentioned examples include:

**Decision Trees:** These algorithms create a tree-like structure to classify claims based on a series of decision rules learned from the data [2].

**Support Vector Machines (SVMs):** SVMs create a hyperplane to separate legitimate claims from fraudulent ones in a high-dimensional space [2].

**Random Forests:** These involve building an ensemble of decision trees, improving overall accuracy and reducing the risk of overfitting [2].

**XGBoost:** This powerful tree boosting algorithm is known for its effectiveness in various classification tasks, including fraud detection [4].

**Data Preprocessing:** A crucial aspect of using ML for fraud detection is data preprocessing. This involves cleaning, transforming, and engineering features from the raw claim data to ensure the model can learn effectively [2, 4].

# **Further Exploration:**

Research comparing the performance of different ML algorithms for insurance fraud detection. Explore specific types of insurance fraud (e.g., auto, health) and how ML is applied to detect them. Investigate the challenges of data privacy and bias in the context of using ML for insurance claim analysis. This survey provides a starting point for your research into fraud detection and analysis for insurance claims using machine learning. By delving deeper into the references provided and exploring the broader literature, you can gain a comprehensive understanding of this evolving field.

## **III. PROBLEM STATEMENT**

Fraud detection and analysis for insurance claims using machine learning is a popular application in the insurance industry. Machine learning algorithms are used to identify patterns in insurance claims data that can potentially indicate fraudulent activities. Here are some common components of a system designed for fraud detection and analysis in insurance claims using machine learning:

1. Data Collection: Gathering various types of data related to insurance claims, such as claimant information, policy details, claim history, and external data sources like historical fraud data or social media data.

2. Data Preprocessing: Cleaning, transforming, and preparing the data for analysis. This step may involve handling missing values, normalizing data, and encoding categorical variables. 3. Feature Engineering: Creating new features or selecting relevant features that can help improve the performance of machine learning models in detecting fraudulent claims.

4. Model Development: Developing machine learning models, such as supervised learning models (e.g., logistic regression, random forest, or neural networks) to predict the likelihood of a claim being fraudulent.

5. Model Training: Training the machine learning models using historical insurance claims data labeled as fraudulent or nonfraudulent.

6. Model Evaluation: Evaluating the performance of the trained models using metrics like precision, recall, F1 score, and ROC-AUC to assess how well the models are able to detect fraudulent claims.

7. Deployment: Implementing the trained machine learning models into a production system that can analyze new insurance claims in real-time and flag potential cases of fraud for further investigation.

8. Monitoring and Updates: Continuously monitoring the performance of the deployed models and updating them with new data to improve their accuracy and effectiveness in detecting insurance fraud. Overall, a robust fraud detection and analysis system for insurance claims using machine learning involves a combination of data processing, modeling, evaluation, and deployment to effectively identify and prevent fraudulent activities in the insurance industry.

## **3.1 Existing System Disadvantages:**

While machine learning-based systems have proven to be effective in fraud detection and analysis for insurance claims, there are some disadvantages and challenges associated with these systems. Here are some common disadvantages of using machine learning in this context:

1. \*\*Data Quality and Imbalance\*\*: Machine learning models require highquality and balanced data for training. Insurance fraud datasets may be imbalanced, with a small number of fraudulent claims compared to legitimate claims, leading to biased models. In addition, incomplete or erroneous data can negatively impact the performance of the models.

2. \*\*Interpretability\*\*: Some machine learning models, such as deep neural networks, are complex and difficult to interpret. Understanding how these models make decisions can be a challenge, especially in regulated industries like insurance where explainability is crucial. 3. **\*\***Overfitting**\*\***: Machine learning models may overfit the training data, capturing noise instead of actual patterns in the data. This can lead to poor generalization on unseen data, reducing the effectiveness of fraud detection in real-world scenarios.

4. \*\*Scalability\*\*: As the volume of insurance claims data grows, scalability can become a concern for machine learning systems. Training and deploying models that can handle large amounts of data efficiently may require significant computational resources.

5. \*\*Adversarial Attacks\*\*: Fraudsters may actively try to deceive machine learning models by manipulating the input data in a way that the model misclassifies fraudulent claims as legitimate. Adversarial attacks pose a significant challenge in maintaining the security and effectiveness of fraud detection systems.

6. \*\*Model Maintenance and Updates\*\*: Machine learning models require regular updates and maintenance to adapt to evolving fraud patterns and changing regulations. Without proper monitoring and updates, the performance of the models may degrade over time.

7. \*\*Ethical and Legal Implications\*\*: Automated decision-making in insurance

fraud detection using machine learning raises ethical concerns related to privacy, fairness, and transparency. Ensuring that the models comply with ethical standards and regulatory requirements is essential but can be challenging. Addressing these disadvantages requires careful consideration of data quality, model interpretability, scalability, security against adversarial attacks, ongoing maintenance, and ethical implications when implementing machine learning systems for fraud detection and analysis in insurance claims.

## **IV. PROPOSED SYSTEM**

This system leverages machine learning to analyze insurance claim data and identify potential fraudulent activities. Here's a breakdown of the key components:

1. Data Collection and Preprocessing:

Gather data on insurance claims, including policyholder details, claim history, type of claim, amount claimed, timestamps, and any other relevant information.

Clean the data by handling missing values, inconsistencies, and outliers. This might involve data imputation or transformation techniques.

Feature engineering: Create new features from existing data that can be better

predictors of fraud. Examples include calculating ratios, identifying inconsistencies between claims, and deriving behavioral patterns.

2. Machine Learning Model Building:

Choose appropriate machine learning algorithms for classification. Common options include:

Logistic Regression: Identifies the probability of a claim being fraudulent based on various factors.

Decision Trees: Creates a tree-like structure to classify claims based on a series of decision rules.

Random Forests: Combines multiple decision trees for improved accuracy and robustness.

Gradient Boosting Machines (GBMs): Creates a series of models where each one improves upon the previous by focusing on the errors of the prior models.

Deep Learning: Particularly useful for complex data with intricate relationships between features.

Train the model on a labeled dataset where claims are already categorized as fraudulent or legitimate. Evaluate the model's performance using metrics like accuracy, precision, recall, and F1 score.

### 3. Fraud Detection and Analysis:

Apply the trained model to new, incoming claims to generate a fraud probability score for each claim. Set a threshold to classify claims as high-risk (potentially fraudulent) or low-risk. Analyze high-risk claims further. This might involve: Integrating with external data sources like social media or public records to verify information. Implementing network analysis to identify clusters of potentially fraudulent claims.

## 4. Benefits:

Improves efficiency and accuracy of fraud detection compared to manual methods. Reduces the cost associated with processing fraudulent claims. Enables faster claim processing for legitimate claims. Provides insights into fraudulent patterns to help refine future detection strategies.

# 5. Considerations:

Data quality is crucial for model performance. Ensure the data is accurate, complete, and representative. Model bias: Be mindful of potential biases in the training data that can skew the model's predictions. Regulatory compliance: Ensure the system adheres to data privacy regulations. By implementing a machine learning-based system for fraud detection, insurance companies can significantly improve their

60

ability to identify and address fraudulent claims, leading to cost savings and a more efficient claims process.

#### 4.1 Proposed System Advantages:

Here are the advantages of using machine learning for fraud detection in insurance claims:

## Advantage Description

Automated and Scalable Fraud Detection Machine learning models can analyze large volumes of claim data efficiently. identifying patterns and anomalies that might be missed by manual review. Improved Accuracy over Time As the model is exposed to more data and fraudulent claims, it continuously learns and refines its ability to detect fraud with higher accuracy. 24/7Real-time Monitoring Machine learning models can operate continuously, enabling real-time analysis of claims as they are submitted, reducing the window of opportunity for fraudsters. Identification of Complex Fraudulent Patterns Machine learning algorithms can uncover intricate relationships between data points that might be difficult to detect with traditional methods, leading to the identification of more sophisticated fraud schemes. Cost Reduction and Efficiency Gains By automating fraud detection and flagging

only high-risk claims for manual review, insurance companies can save time and resources associated with investigating nonfraudulent claims.

#### 4.2 Proposed System Limitations:

Here are some limitations to consider when using machine learning for fraud detection in insurance claims:

Data Dependence: Machine learning models are heavily reliant on the quality and quantity of data used for training. Inaccurate, incomplete, or biased data can lead to inaccurate fraud detection and even perpetuate existing biases.

Black Box Effect: Some machine learning models, particularly complex ones, can be difficult to interpret. This can make it challenging to understand why a particular claim is flagged as high-risk, hindering transparency and explainability in the decision-making process.

Evolving Fraud Tactics: Fraudsters constantly develop new methods to bypass detection systems. Machine learning models may struggle to identify these novel schemes until enough data on them becomes available to retrain the model.

False Positives and Negatives: Even with good training data, models can generate

false positives (flagging legitimate claims) and false negatives (missing fraudulent claims). This necessitates a balance between accuracy and efficiency when setting risk thresholds.

Regulatory Compliance: Data privacy regulations may restrict the collection and use of certain data points that could be beneficial for fraud detection. Insurance companies need to ensure their systems comply with relevant regulations.

Human Expertise Integration: While machine learning automates a significant portion of fraud detection, human expertise remains crucial. Investigators are needed to analyze high-risk claims, understand the context behind the data, and potentially uncover new fraud patterns that the model hasn't yet identified. In conclusion, machine learning offers significant advantages in fraud detection, but it's important to acknowledge and address its limitations. A successful system should combine machine learning's automation and analytical power with human expertise and ongoing monitoring to achieve optimal results.

## **V. SYSTEM ARCHITECTURE**



#### **VI. METHODOLOGY**

**Web Server:** This component likely provides a user interface for users to interact with the system.

Web Database: This stores the data used by the system, likely including historical insurance claim data.

**Rumols thar (unclear function):** The purpose of this component is unclear from the image.

#### Service Provider

**Login:** This functionality allows users to log in to the system.

**Train and Test Data Sets:** This function allows users to upload data sets for training and testing the ML model.

**View Train and Test Results:** This functionality allows users to view the results of training and testing the ML model on the data sets.

**View Trained and Tested Results By Bar Chart:** This functionality allows users to view the results in a bar chart format, likely to help visualize the model's performance.

**View All Bisurance Claim Prediction Type: (unclear function):** The purpose of this function is unclear from the text in the image.

**Find Insurance Claim Ratio Type:** (unclear function): The purpose of this function is unclear from the text in the image.

View Insurance Claim latio Rs (حالب) unclear function): The purpose of this function is unclear from the text in the image, and the last part appears to be in Arabic.

**Download Trained Data Sets:** This functionality allows users to download the trained data sets.

**View All Remote Users:** This functionality allows users to view all remote users of the system. Overall, this diagram appears to depict a system that uses machine learning to predict insurance claim types. Users can upload data sets to train and test the model, view the results, and download the trained data sets. The system also allows users to view all remote users. It's important to note that the specific functionality of some parts of the system is unclear from the image.

## ALGORITHM

1: Data Preprocessing

- Normalize data:  $X' = (X - \mu) / \sigma$  (mean and standard deviation)

- Handle missing values: imputation or interpolation

Step 2: Feature Engineering

Claim frequency:  $CF = \Sigma$  (claims per policyholder)

- Claim amount:  $CA = \Sigma$  (claim amounts)

- Claim type: CT = categorical variable (e.g., accident, theft, etc.)

- Policyholder age: PA

- Policyholder location: PL

- Claim submission time: CST

#### **Step 3: Anomaly Detection**

Statistical methods:

Z-score:  $Z = (X' - \mu) / \sigma$ 

- Modified Z-score: MZ = 0.6745 \* (X' μ) / MAD (median absolute deviation) - Machine learning methods:

Local Outlier Factor (LOF)

Isolation Forest

# **Step 4: Fraud Scoring**

Logistic Regression:

 $P(fraud) = 1 / (1 + e^{-z})$ 

 $z = \beta 0 + \beta 1 * CF + \beta 2 * CA + \dots + \beta n * CST$ 

- Decision Trees/Random Forest:

Fraud score =  $\Sigma$  (feature importance \* feature value)

# **Step 5: Thresholding and Ranking**

Set threshold (e.g., 0.5)

Rank claims by fraud score

# **Step 6: Investigation and Verification**

Investigate top-ranked claims

Verify fraud through additional data or manual review

# **Step 7: Model Updating**

Retrain model with new data and feedback

Update feature engineering and anomaly detection methods as needed

Mathematical techniques used: Probability theory (Bayes' theorem, conditional probability)Statistical inference (hypothesis testing, confidence intervals)Linear algebra (vector operations, matrix multiplication)

Optimization techniques (gradient descent, stochastic gradient descent)

Note: This is a simplified example, and actual implementation may vary based on specific project requireme

# RESULT







# -VIII CONCLUSION

The machine learning models that square measure mentioned which square measure applied on these datasets were able to determine most of the fallacious cases with low false positive rate which suggests with cheap exactness. Certain knowledge sets had severe challenges around data quality, resulting in comparatively poor levels of prediction.

Given inherent characteristics of varied datasets, it would not be sensible to outline optimum algorithmic techniques or use feature engineering process for a lot of higher performance. The models would then be used for specific business context and user priorities. This helps loss management units to specialize in a replacement fraud situations and then guaranteeing that models square measure adapting to spot them. However, it might be cheap to counsel that supported the model performance on backtesting and talent to spot new frauds, the set of models work the cheap suite to use within the space of the insurance claims fraud detection.

## **IX FUTURE ENHANCEMENT**

The future of AI-powered fraud detection in insurance claims is bright. Here are some key enhancements: Advanced Analytics: Leverage deep learning for complex pattern recognition and explore Explainable AI (XAI) for better decision transparency.

Data Enrichment: Integrate external data sources and explore real-time analysis for faster fraud identification.

Adaptive Systems: Develop self-learning models and human-in-the-loop systems for continuous improvement.

These enhancements, along with a focus on network analysis and robust cybersecurity, will strengthen fraud detection and benefit both insurers and honest policyholders.

#### REFERENCE

[1] K. Ulaga Priya and S. Pushpa, "A Survey on Fraud Analytics Using Predictive Model in Insurance Claims," Int. J. Pure Appl. Math., vol. 114, no. 7, pp. 755–767, 2017.

[2] E. B. Belhadji, G. Dionne, and F. Tarkhani, "A Model for the Detection of Insurance Fraud," Geneva Pap. Risk Insur. Issues Pract., vol. 25, no. 4, pp. 517–538, 2000, doi: 10.1111/1468-0440.00080.

[3] "Predictive Analysis for Fraud Detection."

https://www.wipro.com/analytics/comparati

veanalysis-of-machine-learning-techniquesfor-%0Adetectin/.

[4] F. C. Li, P. K. Wang, and G. E. Wang, "Comparison of the primitive classifiers with extreme learning machine in credit scoring," IEEM 2009 - IEEE Int. Conf. Ind. Eng. Eng. Manag., vol. 2, no. 4, pp. 685– 688, 2009, doi: 10.1109/IEEM.2009.5373241.

[5] V. Khadse, P. N. Mahalle, and S. V. Biraris, "An Empirical Comparison of Supervised Machine Learning Algorithms for Internet of Things Data," Proc. - 2018
4th Int. Conf. Comput. Commun. Control Autom. ICCUBEA 2018, pp. 1–6, 2018, doi: 10.1109/ICCUBEA.2018.8697476.

[6] S. Ray, "A Quick Review of Machine Learning Algorithms," Proc. Int. Conf. Mach. Learn. Big Data, Cloud Parallel Comput. Trends, Prespectives Prospect. Com. 2019, pp. 35–39, 2019, doi: 10.1109/COMITCon.2019.8862451.

## **AUTHOR PROFILE**

Mr. K. Jaya Krishna, currently working as an Associate Professor in the Department of Master of Computer Applications, QIS College of Engineering and Technology, Ongole, Andhra Pradesh. He did his MCA from Anna University, Chennai, M.Tech (CSE) from JNTUK, Kakinada. He published more than 10 research papers in reputed peer reviewed Scopus indexed journals. He also attended and presented research papers in different national and international journals and the proceedings were indexed IEEE. His area of interest is Machine Learning, Artificial intelligence, Cloud Computing and Programming Languages.

SHAIK.GOUSE LAZAM currently pursuing Master of Computer Applications at QIS College of engineering and Technology (Autonomous), Ongole, Andhra Pradesh. He Completed B.Sc. in Computer science Acharya Nagarjuna University Guntur, Andhra Pradesh. Her areas of interest are Java & Cloud comput