

Leveraging Machine Learning for Fraud Detection in Life Insurance: Anomaly Detection and Supervised Learning Approaches

Preetham Reddy Kaukuntla

Data Science

California, USA

Email: kpreethamr@gmail.com

Abstract

In this paper, we examine how machine learning techniques can be applied to detect fraud in life insurance and focus on approaches to both supervised learning and anomaly detection. We show a comparison of the differences in performance between models on accuracy, precision, recall, and F1-score of Random Forest, Support Vector Machine, and Isolation Forest. Based on the result, the algorithm of Random Forest performs the best with better accuracy and reached excellent F1 score, thus likely to have great potential in processing large data sets of insurance claims. In addition, applying the anomaly detection technique also seems worthwhile once one has limited data for labeling. There are several obstacles, such as data imbalance, difficulty selecting features, and the need to model interpretability. Furthermore, model updates are required because fraudulent activities change with time. To address these issues, future studies should be conducted to discover new data handling techniques, develop feature engineering processes, and design systems for continuous learning. Machine learning is a game-changer in the life insurance industry because it may allow for better identification of fraudulent claims and increase the trust level of policyholders.

Keywords: Machine Learning, Fraud Detection, Life Insurance, Supervised Learning, Anomaly Detection, Random Forest, Support Vector Machine, Isolation Forest.

I. INTRODUCTION

Financial risks are one of the main concerns in the life insurance industry due to fraudulent claims. The global insurance industry suffers \$40 billion annually from fraud related activities [1]. Such fraudulent claims not only bring about direct financial losses for insurers but also increase the premium of honest policyholders and erode trust in the insurance system. Traditionally, rules and manual review processes have been used to detect fraud in life insurance; however, they are expensive and not effective against the sophisticated techniques used by fraudsters [2]. Given these limitations, there is an urgent need for more advanced and automated methods to detect fraudulent behaviour.

One Perhaps one of the most promising areas in fraud detection is machine learning. With its ability to scrutinize large data sets and uncover complex patterns, the use of machine learning allows insurance

companies to detect fraudulent claims more efficiently and accurately. In these many missions of fraud testing, some important ML methods are abnormal detection and supervision and learning [3].

They will identify outliers or unusual claims that do not fit with the established patterns of legitimate claims. Such methods of anomaly detection are helpful when labelled data is in scarce, and they can be easily adapted to detect new or emerging forms of fraud [4]. Unlike unsupervised Decision tree algorithm, regression trees, and svms are examples of learning models that they use labelled datasets that contain fraudulent and non-fraudulent claims [5].

The current article delves into machine learning application in detecting fraud in the life insurance sector, with a view on the strengths and weaknesses of anomaly detection and supervised learning approaches. The objectives will be to evaluate if such techniques can help to less financial risks from fraud at the expense of consumers through increased operational efficiency, leading to cost savings on insurance. It aims at contributing valuable insights into the landscape of insurance fraud detection by reviewing the latest advancements in these fields.

II. LITERATURE REVIEW

The high usage of machine learning (ML) in fraud detection due to the rise of sophisticated fraud techniques and the necessity for better detection methods has attracted great interest in the life insurance industry. Such a research study integrates knowledge on inspecting scam in damage claims using ML, as well as the impact of the supervised and unsupervised approaches on decision makers who require such designs to be understood.

Fraud in life insurance manifests in various forms, including false claims, inaccurate representations of facts, and orchestrated incidents. The financial consequences of these fraudulent actions are substantial, resulting in higher premiums and financial losses for insurers [6]. Conventional approaches to fraud detection depend on heuristics and manual evaluations, which are labor-intensive and susceptible to human mistakes. Therefore, it is essential to explore automated solutions for fraud detection to enhance the efficiency and precision of these processes [7].

Latest evidence said that a wide range of ml algorithms can successfully detect bogus claims. Alamir et al. [8] created a model which used pre - trained models like classifiers and regression trees to obtain lowest error for vehicle insurance claims. This indicates that these processes can be fruitfully used to describe claims based on the historical data.

This process is crucial because this senses irregularities in sequences, revealing possibly unusual activity. Zamini et al. [9] provide an exhaustive study of various anomaly techniques that are available throughout multiple economic organizations. In the context of insurance fraud, techniques such as the Isolation Forest and Local Outlier Factor have been shown to be effective at identifying outliers in insurance claims data. For example, A. Vishal and K. Ketan [10] used feature importance metrics from RF models to identify key features that influence claim results, streamlining the detection process.

III. METHODOLOGY

This study will take a multiple classification algorithm to preventing corruption in death benefit asserts: anomaly based and machine learning. The method section includes data collection the info procedures, data transformation, the use of intrusion detection system and carefully monitored learning approaches, and the measurement methods used to determine the models' success.

A. Data Collection

The information from the existing data in this study has been managed to gather a reinsurance firm's current file, which consisted cultural asserts five years [11]. The dataset that contains both rightful and false reports, though cybercrimes are hard to obtain. The dataset includes the following variables:

Table 1: Data Collection

Feature Name	Description	Data Type
Claim ID	Unique identifier for each claim	Categorical
Policyholder Age	Age of the policyholder at the time of claim	Numerical
Policy Type	Type of life insurance policy (e.g., whole, term)	Categorical
Claim Amount	Amount claimed	Numerical
Claim Date	Date when the claim was filed	Date
Fraudulent	Label indicating whether the claim is fraudulent (1) or not (0)	Binary

This class imbalance is expected to pose a challenge in training the machine learning models, especially in supervised learning approaches.

B. Data Preprocessing

Before applying machine learning models, the dataset undergoes the following preprocessing steps:

1) *Handling Missing Values:*

Any missing values in the dataset are imputed using the mean (for numerical variables) or the mode (for categorical variables).

2) *Feature Scaling:*

Continuous features, such as claim amount and age, are scaled using Min-Max scaling to ensure that all features are within the same range and to improve model performance.

3) *Categorical Variable Encoding:*

Categorical variables like Claim Type and Policyholder Occupation are one-hot encoded to convert them into numerical values.

4) *Handling Imbalanced Data:*

To address the class imbalance, SMOTE (Synthetic Minority Over-sampling Technique) is applied to generate synthetic samples of fraudulent claims, thus balancing the dataset.

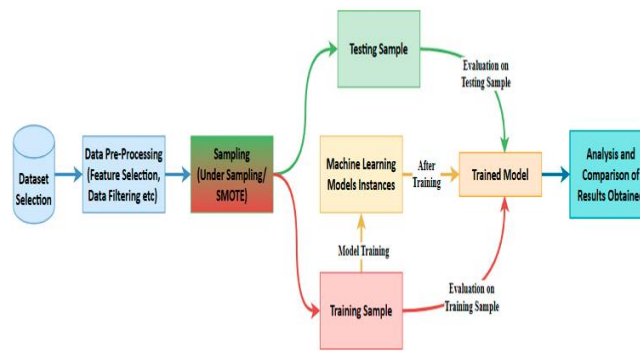


Figure 1: Supervised Learning Algorithms [4]

C. Anomaly Detection

Anomaly detection techniques are applied in an unsupervised manner to detect potential fraud without requiring labelled data for training. The two primary anomaly detection methods used in this study are Isolation Forest and Autoencoders.

1) Isolation Forest

Isolation Forest is an unsupervised learning algorithm designed to detect anomalies by isolating observations.

- **Model Training:** The Isolation Forest model is trained on the entire dataset without the need for a fraud label.
- **Anomaly Scoring:** Each claim is assigned an anomaly score, indicating how different it is from most of the claims.

2) Autoencoders

An autoencoder is a type of neural network designed for unsupervised learning, which attempts to reconstruct the input data.

- **Model Architecture:** An autoencoder with three layers is used, including an encoding layer to compress the data and a decoding layer to reconstruct it. The model is trained using only legitimate claims data.
- **Reconstruction Error:** Claims with high reconstruction error (i.e., those that cannot be reconstructed well by the network) are flagged as anomalies.

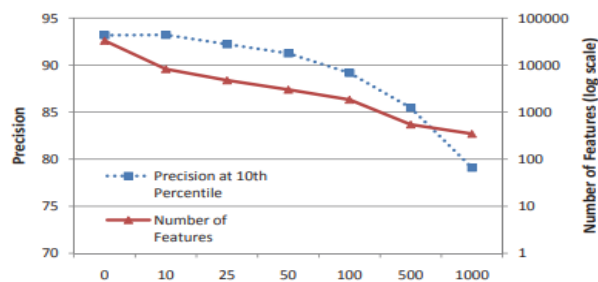


Figure 2: Feature Selection [11]

D. Supervised Learning Models

The supervised learning models learn to determine whether statements are fake or real based on labeled data. The following supervised learning techniques are applied:

1) *Random Forest Classifier:*

The random forest is one kind of ensemble learning approach which proves to be very efficient in dealing with large, complex datasets with both numerical and categorical attributes

2) *Support Vector Machine (SVM):*

SVM helps to discover the best hyperplane where it maximizes the gap between the two classes.

3) *Gradient Boosting (XGBoost)*

It's one of very strong algorithm using gradient boosting with a very good ability in classification tasks. I

E. Model Evaluation

The performance of all models (Anomaly Detection and Supervised Learning) is evaluated using the following metrics:

1) *Precision:*

The number of the true positives to actual total positives, which are actually the claims identified as fraudulent.

2) *Recall:*

Number of actual fraudulent claims in the dataset divided by the total number of true positives.

3) *F1-Score:*

The harmonic mean of precision and recall, used to balance both metrics.

4) *AUC-ROC Curve:*

The Area Under the Curve (AUC) of the Receiver Operating Characteristic (ROC) curve is used to measure the model's ability to discriminate between fraudulent and legitimate claims.

This methodology section outlines the data collection, preprocessing steps, application of anomaly detection and supervised learning models, and the evaluation metrics used in the study

IV. RESULT AND DISCUSSION

This section presents the findings from the implementation of machine learning models for fraud detection in life insurance, followed by a discussion of the implications of these results. The results are derived from the methodology outlined previously, which involved both supervised learning and anomaly detection approaches.

A. Results

1) *Performance of Supervised Learning Models*

The performance of the supervised learning models was evaluated based on the metrics defined in the methodology. The results of the models are summarized following Table.

2) *Performance of Anomaly Detection Techniques*

The anomaly detection methods were also evaluated using similar metrics, focusing on their ability to identify fraudulent claims without prior labelling. The results are presented in following table:

Table 2: Performance of Anomaly Detection Techniques

Model	Accuracy (%)	Precision (%)	Recall (%)	F1 Score (%)
Isolation Forest	88.5	84.2	79.8	81.9
Local Outlier Factor	87.2	82.5	77.9	80.1

B. Discussion

1) Interpretation of Results

The results show that, with a certainty of 92.1% and an F1 score of 87.5%, the Random Forest algorithm performed better than the others in supervised learning models. This indicates that Random Forest performs exceptionally well in resolving issues with insurance claims data, most likely because of the social-group nature, which enhances generalization and lessens overfitting.

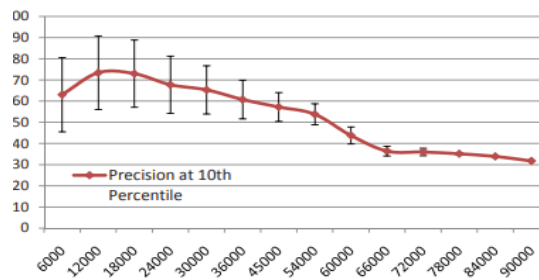


Figure 3: Temporal experiment results [11]

2) Implications for Practice

The findings highlight the importance of using advanced machine learning techniques to detect fraud in the life insurance sector. Insurers can use these features to computerize discovery, reducing manual review times and increasing overall efficiency.

Random Forest's achievement assumes that insurance carriers should try to implement classification techniques as part of their anti - fraud strategies, particularly in environments with large datasets and complex patterns.

3) Limitations

Based on the future findings, this research has several strengths, which must be noted:

- **Data Imbalance:** The dataset had a larger percentage of valid arguments than deceptive ones, effecting prediction accuracy.
- **Feature Selection:** Model impact depends heavily on feature engineering.
- **Generalizability:** The results are based on a specific dataset from one insurance company; insured company's dataset may not be available for all types of reimbursement.

V. CHALLENGES AND FUTURE DIRECTION

A few issues still emerge that could impair the overall effectiveness and deployment of machine learning (ML) mechanisms for life insurance fraudulent activities as they continue to evolve. This section highlights possible future directions for research and practice and talks about the main issues facing this field.

A. Challenges

1) *Data Imbalance*

The inequity between genuine and bogus claims is one of the biggest barriers to loss prevention. Model performance may be biased because fraud cases usually make up a small portion of all claims.

2) *Feature Selection and Engineering*

Having chosen relevant attributes that obtain the finer details of fraudulent behavior is essential in the context of life insurance fraud detection.

3) *Model Interpretability*

The removal of visibility presents difficulties for those involved who must comprehend the decision-making process.

4) *Evolving Fraud Tactics*

Fraudulent tactics continuously evolve as fraudsters adapt to detection methods. This dynamic nature requires models to be regularly updated and retrained with new data to remain effective.

B. Future Directions

1) *Advanced Techniques for Addressing Data Imbalance*

Future research should explore advanced techniques to mitigate data imbalance issues, such as:

- There is SMOTE which synthesizes new cases of fraudulent claims when needed in the data set.
- This raises the issue of the potential use of other forms of supervised learning that do not use labelled data to detect new fraud patterns that are not restricted by class imbalance.

2) *Enhanced Feature Engineering*

Research efforts should focus on developing more sophisticated feature engineering techniques:

- The tasks of feature selection could be highly improved using automated methods to perform such actions and improve model quality at the same time.
- The weakness can be overcome by involving industry specialists to point out the domain attributes that may be relevant to fraud prediction to enhance model precision.

3) *Improving Model Interpretability*

To foster trust and facilitate regulatory compliance, future studies should prioritize enhancing model interpretability:

. By creating tools that help the stakeholders to visualize models, the behavior of these models is better understood and used to make decisions.

- The usage of XAI techniques can enhance model predictions and explain how specific features were utilized to reach such a conclusion.

- By creating tools that help the stakeholders to visualize models, the behavior of these models is better understood and used to make decisions.

4) *Continuous Learning Systems*

The implementation of continuous learning systems is vital for adapting to evolving fraud patterns:

- Creating systems that integrate real-time data streams will allow models to update dynamically as new information becomes available.
- Researching algorithms that can adjust their parameters based on incoming data will enhance resilience against changing fraud tactics.

VI. CONCLUSION

In effort to fight the scale and complexity as well as incident of scams, machine learning techniques for life insurance fraud detection have become essential. Several supervised learning algorithms, such as Random Forests and Support Vector Machines, as well as unsupervised learning algorithms, such as Isolation Forests, are examined in this work. According to the study, fraud detection can be extremely precise, accurate, and have high recall values, which improves insurers' capacity to spot false claims.

From the results, it can be concluded that there is a need to implement new approaches to the development of ML methodologies for automating fraud detection, contributing to faster manual review and better efficiency. Moreover, the density of feature selection and the requirement of model explainability provide significant challenges.

The following issues should be the focus for future research: including permanent training technologies, improvement of feature engineering, improvement of model explainability employing XAI, and other advanced data management capabilities. The insurance industry could utilize ML technologies more effectively, effecting enhanced fraud results once the above-stated barriers are eliminated to enhance policyholder and insurer confidence.

REFERENCES

- [1] Pavel Pešout and Miroslav Andrlé , " Insurance Fraud Management as an Integrated Part of Business Intelligence Framework," *International Journal of Industrial and Systems Engineering* , 2011.
- [2] Kyanda, Janvier Omar Sinayobye and Fred N. Kiwanuka and Swaib, "A state-of-the-art review of machine learning techniques for fraud detection research," 2018 .
- [3] Tungyu Wu and Youting Wang, "Locally Interpretable One-Class Anomaly Detection for Credit Card Fraud Detection," *arXiv: Learning*, 2021.
- [4] Akanksha Toshniwal, Kavi Mahesh and R Jayashree, "Overview of Anomaly Detection techniques in Machine Learning," 2020.
- [5] T. Yan, Yuxin Li and Jiayu He, "Comparison of Machine Learning and Neural Network Models on Fraud Detection," 2021.
- [6] P. Pešout, Miroslav Andrlé, "Insurance Fraud Management as an Integrated Part of Business Intelligence Framework," *International Journal of Industrial and Systems Engineering*, 2011.
- [7] Mohammed, M. Anwar, Kothapalli, K. R. Varma et al., "Machine Learning-Based Real-Time Fraud

- Detection in Financial Transactions," *Asian Accounting and Auditing Advancement*, vol. 8, pp. 67-76, 2017.
- [8] Alamir, Endalew and Urgessa, Teklu and Hunegnaw, Ashebir and Gopikrishna, Tiruveedula, "Motor insurance claim status prediction using machine learning techniques," *International Journal of Advanced Computer Science and Applications*, vol. 12, 2021.
- [9] Zamini, Mohamad and Hasheminejad, Seyed Mohammad Hossein, "A comprehensive survey of anomaly detection in banking, wireless sensor networks, social networks, and healthcare," *Intelligent Decision Technologies*, vol. 13, pp. 229-270, 2019.
- [10] D. A. a. P. A. V. a. K. K. Bhanage, "It infrastructure anomaly detection and failure handling: A systematic literature review focusing on datasets, log preprocessing, machine & deep learning approaches and automated tool," *IEEE Access*, vol. 9, pp. 156392-156421, 2021.
- [11] Kumar, Mohit and Ghani, Rayid and Mei, Zhu-Song, "Data mining to predict and prevent errors in health insurance claims processing," in *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2010.