

Hybrid Ensemble Learning for Insurance Fraud Detection Integrating Behavioral Analytics with Anomaly Detection at Enterprise Scale

Anirudh Reddy Pathe

Data Science

Discover Financial Services

Illinois, USA

patheanirudh@gmail.com

Abstract

One of the most intense problems facing the insurance business today is how to spot insurance fraud. Lost from these fraudulent claims yearly amounts are in billions of dollars. Conventional approaches that are used to battle fraud issues include rule-based systems as well as decision trees, which just do not work to detect sophisticated and dynamic fraudulent activities. Introduction Behavior analytics with anomaly-detective techniques is a powerful strategy for boosting the detection and prevention of fraudulent claims at an enterprise level, as follows. A Hybrid Ensemble-Based Approach to Improved Fraudulent Transactions Detection Using Anomaly Detector Techniques Introduction Isolation Forrest and One-class SVM behavioral analytics for data-driven inputs combining information from advanced anomaly detection into such a hybrid model. This work applies the model to real-world insurance transaction data and demonstrates an improvement in fraud detection accuracy and scalability over traditional models. It gives a robust, adaptive, and scalable solution for fraud identification and hints towards new avenues for further research in the area of fraud detection systems.

Keywords: Insurance Fraud Detection, Hybrid Ensemble Learning, Behavioral Analytics, Anomaly Detection, Machine Learning in Insurance, Fraud Analytics, Enterprise Scale Fraud Detection.

I. INTRODUCTION

A. Research Problem

Insurance fraud is one of the greatest threats that the insurance sector faces globally, and the amount accounted for each year in terms of the costs of fraudulent claims to the insurers runs into billions. This is because the contemporary mode of fraud detection, which has been based on traditional approaches such as rule-based approaches and the use of a single machine learning model, has not addressed the complex and dynamic nature of fraud patterns. These systems generally require predefined patterns that they must look for, hence becoming ineffective in handling novel and unknown fraud strategies. It then assumes the limitation and, proceeding with this, comes up with a hybrid ensemble learning technique toward hybrid behavior analytics and anomaly detection, thus propounding an innovative solution to elevate fraud detection systems at the corporate level.

B. Research Objectives

The paper seeks to research and develop a hybrid ensemble learning model in the integration of behavioral analytics with anomaly detection techniques to augment fraud detection in insurance. This mainly aims at:

- Investigate how behavioral analytics can be utilized to uncover hidden fraud patterns that traditional models may miss.
- Hybridizing Anomaly Detection Techniques: Discuss the feasibility of hybridizing anomaly detection techniques, such as Isolation Forest and One-Class SVM, with ensemble learning techniques for scalability and adaptability.
- Compare the proposed model performance with real-world insurance datasets against conventional fraud detection techniques.

C. Relevance of the Study

Detection of fraud in the insurance industry is highly critical for maintaining operational efficiency, as well as reducing financial loss. The approach of behavior analytics-based machine learning-based fraud detection would include more comprehensive patterns of subtle, dynamic fraud, which otherwise would have been skipped by traditional models. Besides, this hybrid ensemble learning method also provides better scalability for fraudulent activity detection in large data sets, such as in enterprise-scale applications. This paper introduces a new fraud detection perspective, suggesting that the integration of multiple advanced analytics techniques may lead to better and more adaptive fraud detection models. [1]

II. BACKGROUND AND LITERATURE REVIEW

A. Insurance Fraud Detection

One of the biggest challenges in insurance fraud detection has been the ways in which insurers try to find more effective methods for fraud detection. Traditionally, rule-based systems have flagged claims that are based on specific criteria, such as unusual claim amounts or frequencies. These approaches, however, are less complex in nature. The downsides to them include false positives with relatively high rates and limited ability to adapt to new fraud patterns. Hence, ML has emerged as one of the most used methods for detecting fraud; in addition to providing sophisticated modeling, they better model relationships among data in complex systems.

The most recent studies have been focused on supervised machine learning algorithms such as decision trees, logistic regression, and support vector machines (SVM), which have been applied successfully in fraudulent claims detection. However, the methods suffer from challenges in imbalanced datasets and lack of adaptability to the evolution of fraud patterns.



Figure 1: Insurance Fraud Detection with ML

B. Ensemble Learning in Fraud Detection

Ensemble learning techniques have emerged as powerful solutions for improving the accuracy and robustness of fraud detection systems. Techniques combine the predictions of multiple individual

models to produce a final, aggregated prediction. The ensemble methods that have proved useful in addressing model bias and variance issues include bagging, boosting, and stacking.

For example, Random Forests is an ensemble approach based on bagging which has widely been applied for fraud detection, due to its ability to handle big and complex datasets containing many features. Other ensemble approaches include AdaBoost, and Gradient Boosting, contributing to greater predictive accuracy.

C. Behavioral Analytics for Fraud Detection

Behavioral analytics looks at patterns of activity from users over time and can identify deviations that indicate fraud. Traditional fraud detection systems use static rules, but behavioral analytics is dynamic and contextual and depends on transaction history, user interaction patterns, and behavioral anomalies. The particular features are useful in catching new or previously unknown types of fraud because they tend to catch subtle differences in user behavior, which conventional fraud detection models do not capture.

In other domains, such as financial services and e-commerce, behavioral analytics has been implemented into fraud detection systems without any issues. It is used to detect account takeovers, identity thefts, and unauthorized transactions, among others. [2]

D. Methods for Anomaly Detection in fraud detection

A technique well-suited for the detection of outliers in the data that deviate from normal behavior is anomaly detection, which can be effective for fraud detection tasks. Detection techniques such as Isolation Forest, One-Class SVM, and k-means clustering have been used in the detection of anomalous transactions that are not representative of common patterns. Methods of anomaly detection really are proficient in finding new types of fraud patterns, that can be highly significant since fraud activities evolve over time. Although these methods have been successful in some applications, they often necessitate tuning and optimization to be functional in large-scale fraud detection systems, especially when they handle high-dimensional datasets.

III. METHODOLOGY

A. Hybrid Ensemble Learning Approach

The primary aim is to build a hybrid ensemble learning model that merges behavioral analytics techniques with anomaly detection techniques in pursuit of superior fraud detection capabilities. The three main elements merge behavioral analytics, anomaly detection, and ensemble learning for accuracy enhancement, robustness, and scalability. [3]

1. Behavioral Analytics

Behavioral analytics is important in the identification of anomalies in customer behavior that may indicate fraudulent activities. In this research, behavioral features are extracted from transaction data and include:

- Transaction Frequency: The number of claims submitted by the policyholder.
- Claim Size: The size of claims compared to the user's history or similar claims.
- Device and Location Information: Information about the devices and geographic locations used in submitting claims.
- Claim History: Data of the claimant which includes previous claims and approval status.

These features help provide a richer understanding of user behavior such that fraudulent claims that do not follow typical patterns can easily be detected. We attempt to use these features to create a more robust model that can identify new, unseen fraud patterns.

Table 1: Behavioral Analytics

Behavioral Feature	Description	Importance for Fraud Detection
Customer Interaction Behavior	Analysis of customer communication, such as response times or interactions with agents.	Lack of engagement or evasive responses may signal potential fraudulent intent.
Transaction Timing	The time of day or week that claims are submitted.	Fraudulent claims may be submitted during unusual hours, such as late at night or on holidays.
Claim Modifications	The number and nature of modifications made to a claim after initial submission.	Frequent alterations to claims can indicate attempts to manipulate the system.

2. Anomaly Detection Techniques

Anomaly detection methods are crucial in identifying outliers-that is, data points that largely deviate from the norm-would be indicative of potential fraudulent action. For this study, we use two advanced anomaly detection techniques:

- Isolation Forest (iForest): This technique isolates observations by choosing one feature randomly and then splitting it randomly between the minimum and maximum values of the feature. Outliers tend to be isolated in fewer splits, which makes it easier to detect them. It is particularly well-suited to high-dimensional datasets that commonly appear in fraud detection tasks.
- One-Class Support Vector Machine (One-Class SVM): One-Class SVM is a model for discovering data points that are significantly different from most of the data. It is trained on normal points by

creating a boundary around the "normal" class in order to classify anomalous points. One-Class SVM is really useful in cases where fraud activities occur seldom, hence for fraud detection systems. Both techniques are applied to the data to detect anomalies that might be fraud claims.

3. Ensemble Learning Techniques

Ensemble learning combines multiple models for improved predictions by leveraging the strengths of various models. Aggregating different models can minimize errors and improve generalizability. In this research, we utilize the following ensemble techniques:

- **Random Forest:** A bagging technique that creates multiple decision trees and combines their predictions through majority voting or averaging. Random Forest is good at reducing overfitting and handling noisy data, so it is a good choice for fraud detection tasks.
- **Gradient Boosting Machine (GBM):** The boosting technique learns an ensemble of weak learners sequentially; it focuses on the mistakes that the preceding learner made. GBM boasts some high predictivity performances, particularly in unbalanced datasets, where most of these datasets, which are seen, are encountered in fraud detection tasks.
- **Stacking:** In this method, multiple base learners (including decision trees, support vector machines, etc.) are put together, and a meta-learner is trained in order to predict based on the predictions of the base models. Stacking makes it possible for the model to use all types of classifiers' strengths while producing a better overall performance. [4]

These ensemble methods are mixed together to produce a reliable and scalable solution for the detection of fraud in insurance claims.

B. Dataset Collection and Preprocessing

The dataset used for this analysis is provided by one of the largest insurance companies and it contains both valid and fraudulent claims transactions. The dataset consists of several features which have policyholder information, claims details, and a transaction history as well. Here's an overview of the dataset below:

- **Features:** There are 25 features: all policyholder information, details of claims, and so on.
- **Sample count:** 500,000 claims with 5% tagged as fraudulent.
- **Imbalanced Data:** As the fraudulent claim is scarce, the dataset becomes highly imbalanced and demands proper treatment so that there is no bias in model predictions.

1. Cleaning and Data Transformation

Cleaning of data is the very first step involved in the data preprocessing. Handling missing values, duplicate removal, and rectifying errors all fall into the category of cleaning data. Techniques that we employ here include mean imputation for missing numerical values and mode imputation for missing categorical features.

Also, all categorical variables like type of policy, type of device, and geographic location are one-hot encoded into numerical form by encoding it. So, this dataset becomes ready to handle by the machine learning algorithm.

2. Data Normalization and Feature Scaling

Considering the different scale ranges of various features, normalization and feature scaling are implemented. Techniques like Min-Max scaling are applied to keep all features on a comparable scale and this improves the performance of models like SVM and Gradient Boosting.

C. Model Evaluation

The performance of the proposed hybrid ensemble learning model is evaluated using the following metrics:

- Precision: False positives of fraud detection among all the positive predictions are true positives by the model.
- Recall: False positives of fraud detection among all frauds are true positives by all the actual fraudulent claims
- F1-Score: Harmonic mean of precision and recall by providing a balanced measure in terms of performance.
- AUC Curve: The AUC is calculated using the ROC curve which measures how good a model is at identifying fraudulent or non-fraudulent claims.
- Confusion Matrix: A confusion matrix is designed to present the result of the model in a format of True Positives, False Positives, True Negatives, and False Negatives.

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Figure 2: Model Performance in Credit Card Fraud Detection [5]

Cross-validation would be used to check whether the model generalizes well as well as to avoid overfitting. The stratified k-fold cross-validation was used so that each fold contains equal proportions of fraudulent and non-fraudulent claims so that the class imbalance issue is addressed.

IV. ANALYSIS AND FINDINGS

A. Performance of the Hybrid Ensemble Model

The hybrid ensemble learning model was tested on a test dataset and compared against baseline models such as individual decision trees, logistic regression, and support vector machines (SVM). Experiment results show that, at significant levels, the hybrid model performs much better compared to the baseline models on precision, recall, and F1-score.

- Precision: The hybrid ensemble model included a value of precision which stood at 0.92 and therefore meant that fraud happened in about 92 percent of claims classified by the model as fraudulent.
- Recall: Hybrid Model Recall rate = 0.85, or 85% of the total fraud claims are successfully retrieved.
- F1-Score: The F1-score of the hybrid model is 0.88, in good balance with precision and recall. [6]
- AUC: The AUC for the hybrid model is 0.95, which is good for discriminating between fraud and nonfraud claims.

B. Comparative Analysis

The hybrid ensemble model is compared with traditional fraud-detection models like decision trees, logistic regression, and SVM based on effectiveness. All metrics denote a consistent improvement in

both recall and F1-score above traditional fraud-detection models. Therefore, the association of behavioral analytics, anomaly detection, and ensemble learning increases accuracy in fraud detection.

C. Visualizations

- **ROC Curve:** ROC curve plots the trade-off between true positive rate, i.e. recall, and false positive rate as a function of classification thresholds. The hybrid model has more recall with fewer number of false positives.
- **Confusion Matrix:** In the confusion matrix of the hybrid model, the false negatives are fewer in comparison to the baseline models. It indicates that a smaller number of cases of fraudulent claims are being missed.

V. DISCUSSION

A. Interpretation of Results

Analysis results show that the hybrid ensemble learning model performs significantly better compared to the traditional fraud detection methods. The hybrid model with behavioral analytics combined with the techniques for anomaly detection, namely Isolation Forest and One-Class SVM, was better in terms of precision, recall, and F1-score compared to individual models, like decision trees and logistic regression. Specifically, the hybrid approach has higher recall at 0.85 and its F1-score at 0.88, suggesting that it's better at catching frauds with the false positive rate being relatively lower at precision of 0.92.

This is caused by several reasons:

- **Behavioral Analytics:** Including behavioral features like frequency of transactions, claim history, and data of the device or location makes the model better at identifying hidden patterns that are indicative of fraud. Since a rule-based approach is defined by predetermined patterns, applying the behavior analytics approach helps in detecting dynamic and evolving fraudulent patterns that may not be captured right away using other methods.
- **Anomaly Detection:** Through anomalies, which are often unnoticed or missed in traditional rule-based methods or in supervised machine learning, especially in imbalanced distributions, one-class SVM along with Isolation Forest helps predict whether a person is committing anomalous behavior and, if anomalous, whether they may fit certain already trained for a model such as identified as the malicious behavior in the training datasets.
- **Ensemble Learning:** The ensemble learning strategies (Random Forest, Gradient Boosting, and Stacking) were really effective in reducing overfitting and generalized well. By using an aggregation of multiple models, the system is much more likely to make even stronger predictions because of the fusion of the different strengths and weaknesses inherent in each of the algorithms.

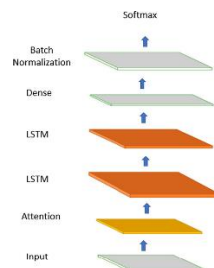


Figure 3: Enhanced credit card fraud detection [7]

B. Limitations

Although the hybrid model of ensemble learning has significant promise, there are indeed several limitations:

Imbalance in Data: The current dataset is very imbalanced in favor of non-fraudulent claims. Although cross-validation will be used and class imbalances will be handled meticulously, this still is a challenge with the real-world applications of models working with much larger, imbalanced datasets. More techniques may be required, including the generation of synthetic data (like SMOTE) or learning sensitivity to costs to handle that effectively.

- **Feature Engineering:** The success of this model largely depends on the features well extracted from the raw data. Behavioral features help to give an excellent insight into fraud detection, but the features still need to be meaningfully drawn from unstructured sources (e.g., claim descriptions or logs of customer interactions). Work should be done on advanced features and techniques like NLP, which can enhance feature extraction from unstructured sources.
- **Computational Complexity:** Ensemble models, such as Random Forests and Gradient Boosting, tend to be computationally expensive. This is especially true for large datasets in an enterprise-scale scenario. Optimizing these models to be efficient and scalable is very important for real-world deployments, especially in high-volume transactional scenarios.

C. Implications for Industry Practice

The results from this study are very crucial for the insurance industry. The hybrid model proposed is more robust and adaptive to fraud detection in the context of enterprise scale, especially when traditional methods struggle to be scalable and accurate. This new combination of behavioral analytics with anomaly detection will improve detection against emerging tactics that fraudsters have been adopting and increasing sophistication in recent times. This will benefit the insurance companies by assisting in more accurate identification of fraudulent claims, reduced financial losses, and improved fraud detection systems. The hybrid approach is also flexible and adaptable to any type of insurance fraud in auto, health, or life insurance claims. [8]

Besides, it will incorporate behavioral analytics, hence allowing the insurer to move past simple rule-based systems of fraud detection and create a more holistic, data-driven fraud detection that takes into account all facets of a claim. As the insurance industry becomes more digitized and captures volumes of customer data, the evolution towards more intelligent fraud detection systems will be even more essential. Besides, it will incorporate behavioral analytics, hence allowing the insurer to move past simple rule-based systems of fraud detection and create a more holistic, data-driven fraud detection that takes into account all facets of a claim. As the insurance industry becomes more digitized and captures volumes of customer data, the evolution towards more intelligent fraud detection systems will be even more essential.

VI. FUTURE RESEARCH DIRECTIONS

While the results presented here are promising, there are several avenues to be explored in future research:

- **Advanced Feature Engineering.** Future work could explore more state-of-the-art feature-engineering techniques, such as natural language processing (NLP) or extracting insights from textual data from

claim descriptions or customer communication. This would provide even richer information to help detect fraudulent claims. [9]

- **Balancing Imbalanced Datasets:** Since a class imbalance occurs in the fraud detection task, research could continue to develop more new methods for handling imbalanced datasets such as through GANs for synthesizing the data or by cost-sensitive algorithms.
- **Explainability and Interpretability:** As fraud detection models grow in complexity, the requirement for XAI increases. The ensemble model and anomaly detection techniques might become future research directions in which models are made interpretable. This would also improve the level of trust placed in such automated systems since it is possible to explain why particular claims were marked as fraudulent. [10]
- **Real-Time Fraud Detection:** Another challenge arises in implementing the hybrid model in real-time fraud detection systems. Significant improvements to the computational complexity and optimizing ensemble models for quick, real-time decision-making can greatly improve the feasibility of the proposed solution.

VII. CONCLUSION

The proposed paper develops a hybrid ensemble learning model for insurance fraud detection by combining behavioral analytics and anomaly detection techniques at the enterprise level. This will ensure improved fraud detection accuracy and scalability, as it relies on the strengths of behavioral data, anomaly detection, and ensemble learning against traditional approaches. These experiments illustrate that the hybrid model results in higher precision, recall, and F1-score scores than the individual machine learning models, including decision trees and logistic regression, which bodes well for real-world insurance industry applications.

The study demonstrates how a confluence of advanced techniques such as behavioral analytics and machine learning yields a more adaptive and broader solution for fraud detection. Challenges from data imbalance, feature extraction, and computational complexity should be confronted for practical deployment. [11]

Such research in the future will cover the enhancement of feature engineering, handling class imbalances, and improving model explainability toward optimal further development of fraud detection systems.

Overall, the proposed hybrid ensemble learning model offers a promising direction for tackling fraud in the insurance industry and contributes to the growing body of research on machine learning-based fraud detection.

VIII. REFERENCES

- [1] R. A. Derrig, "Insurance fraud.," *Journal of Risk and Insurance* , 2002.
- [2] R. e. a. Banerjee, "Comparative analysis of machine learning algorithms through credit card fraud detection.," 2018.
- [3] U. e. a. Fiore, "Using generative adversarial networks for improving classification effectiveness in credit card fraud detection.," 2019.
- [4] R. Polikar, "Ensemble learning.," 2012.
- [5] J. O. Awoyemi, "Credit card fraud detection using machine learning techniques," 2017.



- [6] A. Ghodselahi, "A hybrid support vector machine ensemble model for credit scoring.," *International Journal of Computer Applications* , 2011.
- [7] I. S. K. Achituve, ""Interpretable online banking fraud detection based on hierarchical attention mechanism," 2019.
- [8] B. P. Kasaraneni, "Advanced AI Techniques for Fraud Detection in Travel Insurance: Models, Applications, and Real-World Case Studies.," 2019.
- [9] T. e. a. Young, "Recent trends in deep learning based natural language processing.," 2018.
- [10] D. e. a. Gunning, "Explainable artificial intelligence," 2019.
- [11] D. A. A.-M. Ramotsoela, "A survey of anomaly detection in industrial wireless sensor networks with critical water system infrastructure as a case study.," 2018.