Vehicle Insurance Fraud Detection System

Pushpavalli P¹, Shamitha M², Aishwarya S³, Shyleshwari M Shetty⁴

^{1, 2, 3} Dept of CSE ⁴Assistant Professor, Dept of CSE ^{1, 2, 3, 4} GSSSIETW, Mysuru, India

Abstract- Concern over insurance fraud in the machine sedulity is growing, as it can affect in significant financial losses and advanced decorations for law- abiding policyholders. This study offers a machine knowledgepredicated system for further directly relating false claims. To address class imbalance, we use a Kaggle dataset and the SMOTE fashion. With a Python Flask- predicated web interface that lets stoners enter claim details and get immediate fraud discovery results, our system predicts fraud using Random Forest, Decision Tree, and Logistic Regression models.

Keywords- Vehicle Insurance, Fraud Detection, SMOTE, Random Forest, Flask, Machine Learning.

I. INTRODUCTION

The purposeful form of false or inflated claims for financial gain is known as insurance fraud, and conventional discovery ways constantly prove shy, particularly as fraud ways advance in complication. By using SMOTE to balance the dataset, this design addresses the oddity of fraudulent exertion by using machine knowledge (ML) to find patterns in data that indicate it. To find the swish strategy, we compare the performance of several machine knowledge models, analogous as Random Forest, Decision Tree, and Logistic Regression. The finished product is incorporated into an intuitive Teacup- erected web operation that lets stoners enter claim information with ease and get realtime prognostications about the validity or fraud of a claim. pivotal issues like data imbalance, real- time discovery, scalability. This system successfully addresses important issues like data imbalance, real- time discovery, scalability, and user vacuity. SMOTE integration improves overall discovery delicacy by adding the model's perceptivity to uncommon fraud cases. The system makes sure the swish algorithm is applied for accurate prognostications by comparing several models. The system is useful for everyday use thanks to the Flask web operation's responsive and intuitive interface.

II. LITERATURE SURVEY

Insurance fraud remains a significant challenge in the automobile insurance sector, causing substantial financial

losses and increasing operational inefficiencies. Traditional manual methods and rule-based systems have proven insufficient in detecting sophisticated and evolving fraudulent patterns. In response, recent studies have leveraged machine learning (ML) and data-driven approaches to build more accurate and scalable fraud detection systems.

[1]Detection of Insurance Fraud in Automobile Claims Using Machine Learning by Kumar, R., Sharma, S., and Jain, A. (2022) proposes a machine learning-based approach for detecting fraudulent automobile insurance claims. The authors highlight the importance of data preprocessing, model selection, and evaluation metrics to enhance prediction accuracy. Their study shows how the integration of intelligent systems can significantly minimize fraudulent activity and improve operational efficiency in the insurance sector.

[2] Fraud Detection in Insurance Claims Using Machine Learning Algorithmsby Singh, N., Patel, R., and Mehta, M. (2021) explores the use of machine learning models such as Decision Trees and Logistic Regression to detect fraudulent insurance claims. The research focuses on evaluating model performance using precision, recall, and F1-score, emphasizing the importance of choosing the right algorithm for timely and accurate fraud prediction.

[3] Insurance Fraud Detection: A Comparative Study of Machine Learning Modelsby Gupta, A., Verma, P., and Deshmukh, S. (2022) presents a comparison of multiple machine learning models for fraud detection. The study concludes that ensemble learning techniques outperform single classifiers in terms of accuracy and stability, making them more suitable for practical fraud detection systems.

[4]Anomaly Detection in Insurance Claims Using Machine Learning Techniquesby Sharma, S., Kapoor, R., and Agarwal, P. (2021) investigates the effectiveness of anomaly detection techniques using unsupervised and semi-supervised machine learning. The study emphasizes how these methods can uncover new fraud patterns that rule-based systems might miss, thereby improving the adaptability of fraud detection frameworks. [5] Predicting Fraudulent Automobile Insurance Claims Using Ensemble Learningby Kumar, J., Mehta, S., and Verma, A. (2022) introduces an ensemble learning method to improve fraud prediction accuracy. By combining multiple base models, their approach increases robustness and reduces the risk of overfitting, making it effective in identifying fraudulent claims in varied scenarios.

[6] Fraudulent Claim Prediction in Insurance Using Data Mining Techniquesby Patel, S., Gupta, R., and Mehra, M. (2020) applies data mining approaches such as classification and clustering to detect insurance fraud. The study highlights the role of pattern discovery in automating the claims evaluation process and reducing manual oversight.

[7]A Hybrid Approach to Insurance Fraud Detection Using Deep Learning and Ensemble Methodsby Verma, M., Jain, P., and Kumar, A. (2023) proposes a hybrid model that integrates deep learning with ensemble techniques for improved fraud detection. Their method effectively handles complex, high-dimensional datasets and adapts to evolving fraud behaviors, offering better scalability and accuracy.

[8]Auto Insurance Claims Data (Kaggle Dataset) retrieved from Kaggle (https://www.kaggle.com) provides the dataset used for training and evaluation in the fraud detection system. The dataset, initially imbalanced, was processed using the Synthetic Minority Oversampling Technique (SMOTE), which helped improve model sensitivity to rare fraudulent cases and overall prediction performance.

III. METHODOLOGY

The methodology for the Vehicle Insurance Fraud Detection System involves several systematic steps, starting with data acquisition, preprocessing, model training, evaluation, and deployment. The first step in the process is the collection of data from a publicly available Kaggle dataset that contains real-world auto insurance claim records. This dataset includes various features such as claim ID, customer demographics, vehicle information, incident details, and claim status. Upon initial examination, it was observed that the dataset was highly imbalanced, with a small proportion of fraudulent claims compared to legitimate ones, which poses a challenge for traditional machine learning models.

To overcome this imbalance, the Synthetic Minority Oversampling Technique (SMOTE) is applied. SMOTE is a data augmentation method that synthetically generates new instances of the minority class (fraud cases) by interpolating between existing examples. This results in a more balanced dataset, improving the ability of the model to learn and identify fraudulent patterns effectively. The data then undergoes extensive preprocessing, including handling missing values, encoding categorical variables using techniques like one-hot encoding or label encoding, normalization of numerical features, and removal of outliers to ensure data quality and consistency.

Once the data is prepared, it is split into training and testing sets, typically in an 80:20 ratio, to train and validate the machine learning models. In this project, three supervised learning algorithms—Logistic Regression, Decision Tree, and Random Forest—are selected due to their suitability for binary classification tasks and interpretability. Each model is trained using the balanced dataset, and hyperparameter tuning is performed using grid search or cross-validation to improve performance.

The trained models are then evaluated based on key performance metrics, including accuracy, precision, recall, and F1-score. These metrics help in determining the model's effectiveness, especially in detecting minority class instances (fraud cases) without misclassifying genuine claims. Among the models, the one with the highest F1-score and balanced precision-recall tradeoff is selected as the optimal model for deployment.

To provide a practical and accessible interface, the selected model is integrated into a Python Flask-based web application. The web interface allows insurance claim officers or system administrators to register, log in, and input claim details through structured forms. Upon submission, the data is preprocessed in real time, passed to the trained model, and the prediction output—whether the claim is likely fraudulent or legitimate—is displayed to the user instantly. This real-time fraud prediction capability helps in accelerating the claim verification process and minimizes manual intervention.

Furthermore, the system architecture is designed to support two phases: a training phase, where models are periodically retrained on new data to keep them updated with emerging fraud trends, and a prediction phase, which processes real-time inputs and displays outputs to the user. The system also includes provisions for system administrators to monitor performance, manage users, and trigger model updates when needed.

The complete system is built using Python, Flask, HTML/CSS/JavaScript for the frontend, and Scikit-learn, Pandas, and MySQL for machine learning and backend data handling. The deployment environment supports scalability and portability and runs efficiently on systems with moderate hardware requirements. This comprehensive methodology ensures that the system is not only accurate and robust but also user-friendly, scalable, and adaptable to real-world insurance fraud scenarios.

The methodology adopted for the Vehicle Insurance Fraud Detection System effectively combines data preprocessing, machine learning, and web development to create a reliable and efficient solution for identifying fraudulent claims. By addressing key challenges such as data imbalance and real-time prediction, the system enhances the decision-making process for insurance companies. Its scalable architecture and user-friendly interface ensure that it can be practically implemented in real-world environments to reduce financial losses and improve operational transparency.

IV. SNAPSHOTS



Snapshot 1: Home Page



Snapshot 2: Register Page



Snapshot 3: Database Table

Insurance Fraud Detection		Home Register Login
Login Enai: puripa31389(ggmal.com	Password:	
		Activate Windows Gelo Teringo to attuate Windows,
D Seach	()) 🖬 🖞 🖑 🖬 刘 😁 🔉	∧ 6 ¹⁰⁵ ⊕ 0 10 ¹⁰⁵⁴ 0 ∧
Snaps	hot 4: Login Page	



Snapshot 5: Fraud Prediction Page



Snapshot 6: Fraud Prediction Result

→ C 0	127881 (COlipedictical	-	distant.		0 9 9	•
	INCIDENT TYPE:		COLLISION TYPE.			
	Multi-vehicle Collision	~	Rear Collision			
	INCIDENT SEVERITY:		AUTHORITIES CONTACTED			
	Major Damage	*	Fite	*		
	INCIDENT HOUR OF DAY.		NUMBER OF VEHICLES			
	13		2			
	PROPERTY DAMAGE		BODILY INJURIES			
	NO	*	0			
	WITNESSES		POLICE REPORT AVAILABLE			
	3		NO	*		
	TOTAL CLAIM AMOUNT		INJURY CLAIM			
	500		100			
	PROPERTY CLAIM		VEHICLE CLAIM	Actually		
	200		200	Col - Second		





Snapshot 8: Fraud Prediction Result

V. CONCLUSION

This project, the Vehicle Insurance Fraud Detection System, offers a practical and intelligent solution to help insurance companies tackle the growing issue of fraudulent claims. Using machine learning models like Logistic Regression, Decision Tree, and Random Forest, the system learns from past data to accurately predict whether a new claim is likely to be genuine or fake. One of the standout features of this project is the use of SMOTE (Synthetic Minority Oversampling Technique), which balances the dataset and ensures the model can detect even the rare cases of fraud effectively.

From collecting and cleaning the data to building and evaluating models, and finally integrating everything into a user-friendly web application using Flask, every step in this project is aimed at making fraud detection faster and more reliable. The web interface makes it easy for claim officers to input details and get real-time predictions, saving time and reducing the need for manual checking.

The system is designed to be flexible and can be improved over time by retraining it with new data, making it better at catching new types of fraud as they emerge. This not only helps reduce financial losses but also improves trust in the claim process. By automating the detection process, the system takes some pressure off employees and helps companies make quicker, more confident decisions.

In essence, this project shows how machine learning and technology can be used together to solve real-world problems. While the current version already performs well, future upgrades—like using deep learning, analysing claim descriptions with text analysis, or connecting the system to live databases—could make it even more powerful. This project is a solid step toward smarter, more secure insurance systems.

With the increasing number of digital transactions and online insurance processes, the risk of fraud is only expected to grow. Implementing such a fraud detection system helps insurance companies stay ahead by proactively identifying suspicious claims before they are processed. Additionally, the use of interpretable models ensures that users can understand the reasoning behind predictions, which adds transparency to the decision-making process. The system can also be customized for different insurance providers based on their specific claim patterns and requirements.

REFERENCES

- [1] Kumar, R., Sharma, S., & Jain, A. (2022). Detection of Insurance Fraud in Automobile Claims Using Machine Learning.
- [2] Singh, N., Patel, R., & Mehta, M. (2021). Fraud Detection in Insurance Claims Using Machine Learning Algorithms.
- [3] Gupta, A., Verma, P., & Deshmukh, S. (2022). Insurance Fraud Detection: A Comparative Study of Machine Learning Models.
- [4] Sharma, S., Kapoor, R., & Agarwal, P. (2021). Anomaly Detection in Insurance Claims Using Machine Learning Techniques.
- [5] Kumar, J., Mehta, S., & Verma, A. (2022). Predicting Fraudulent Automobile Insurance Claims Using Ensemble Learning.
- [6] Patel, S., Gupta, R., & Mehra, M. (2020). Fraudulent Claim Prediction in Insurance Using Data Mining Techniques.
- [7] Verma, M., Jain, P., & Kumar, A. (2023). A Hybrid Approach to Insurance Fraud Detection Using Deep Learning and Ensemble Methods.
- [8] Kaggle Dataset: Auto Insurance Claims Data. Retrieved from: https://www.kaggle.com