

Insurance Claim Fraud Detection

Objective

To develop a machine learning model that accurately detects fraudulent insurance claims by analyzing customer claim data. The aim is to reduce fraud-related losses by identifying suspicious claims based on multiple relevant features.

Dataset Description

A synthetic dataset containing **1000 insurance claim records** was generated. It includes both fraudulent and non-fraudulent cases, with an increased fraud ratio (30%) to better train classification models.

Key Features:

- **Claim Amount** – Total value of the insurance claim.
- **Injury Type** – Type of injury reported: None, Minor, or Major.
- **Repair Shop** – Whether the repair was done at a partnered or non-partnered shop.
- **Claim Delay** – Number of days between the incident and the claim filing.
- **Vehicle Age** – Age of the vehicle in years.
- **Previous Claims** – Number of claims previously filed by the claimant.
- **Region** – Geographical region where the claim originated.
- **Label** – Output class indicating if the claim is **Fraud** or **Not Fraud**.

Data Preprocessing

- Categorical features such as injury type, repair shop, and region were encoded using one-hot encoding.
- The target variable was converted to binary: Fraud = 1, Not Fraud = 0.
- Numerical features were standardized to improve model performance.

Data Splitting

The dataset was split into:

- **80% Training Data** – Used to train machine learning models.
- **20% Testing Data** – Used to evaluate model performance.

Models Used

1. Isolation Forest (Unsupervised)

- Used to identify outliers and anomalies without prior fraud labels.
- Effective in flagging unusual patterns in the claim data.

2. Random Forest Classifier

- Supervised ensemble model that uses multiple decision trees.
- Achieved high accuracy and recall on fraudulent claims.

3. XGBoost Classifier

- Gradient boosting model designed for performance and interpretability.
- Provided precise results and was further interpreted using SHAP values.

Evaluation Metrics

- **Confusion Matrix** – Showed true/false positives and negatives.
- **Precision & Recall** – Focused on minimizing false positives and false negatives.
- **F1-Score** – Balanced metric for evaluating fraud detection accuracy.

Model Explainability with SHAP

- SHAP (SHapley Additive exPlanations) was used to interpret predictions made by XGBoost.
- Provided insights into which features contributed most to identifying fraudulent claims.
- The most impactful features included: **claim_amount**, **claim_delay**, and **repair_shop**.

Fraud Hotspot Analysis

- Regional fraud frequency was analyzed to detect **fraud-prone regions**.
- Visualizations highlighted certain areas (e.g., South or East) with higher fraud activity, useful for risk mitigation strategies.

Project Outcome

- Successfully built models to detect insurance fraud with high precision.
- Enhanced interpretability through SHAP.
- Identified feature patterns and regions where fraud is more likely to occur.

Future Scope

1. Integrate real-time data from insurance portals.
2. Apply deep learning models like Autoencoders for anomaly detection.
3. Build a dashboard for live fraud alerts and claim tracking.

Conclusion

This project applied machine learning to detect fraudulent insurance claims using a synthetic dataset. Models like Isolation Forest, Random Forest, and XGBoost delivered reliable classification performance. SHAP analysis helped explain key fraud-driving features like claim amount and delay. The approach demonstrates the potential of ML in reducing insurance fraud effectively.