# FRAUDULENT CLAIM DETECTION CASE STUDY

By Manasa Madhusoodanan and Manisha Mathur

Content

Business Objective

Problem statement

Summary

Recommendations

Business Implications of the findings

Model Performance

Conclusion

Q/A

## Business Objective

Global Insure wants to build a model to classify insurance claims as either fraudulent or legitimate based on historical claim details and customer profiles.

By using features like claim amounts, customer profiles and claim types, the company aims to predict which claims are likely to be fraudulent before they are approved.

# Problem Statement

Global Insure, a leading insurance company, processes thousands of claims annually. However, a significant percentage of these claims turn out to be fraudulent, resulting in considerable financial losses.

The company's current process for identifying fraudulent claims involves manual inspections, which is time-consuming and inefficient.

Fraudulent claims are often detected too late in the process, after the company has already paid out significant amounts.

Global Insure wants to improve its fraud detection process using data-driven insights to classify claims as fraudulent or legitimate early in the approval process.

This would minimises financial losses and optimises the overall claims handling process.

# Summary

- Built a machine learning model to detect fraudulent insurance claims.

- Cleaned and preprocessed data (handled missing values, encoded categories).

- Used EDA to find fraud patterns (claim amount, claim type, customer profile).

- Handled class imbalance using class weight balancing in models.

- Trained Logistic Regression, Decision Tree, and Random Forest models.

- Evaluated using accuracy, precision, recall, F1-score — focused on recall.

- Identified key features and provided business insights to reduce fraud.

# Recommendations

- **Use the model to flag suspicious claims early,** so they can be reviewed by the fraud investigation team before approval.

- **Focus more on high-amount claims** and certain claim types, as they showed stronger links to fraudulent activity in the analysis.

- **Keep the model updated** by regularly retraining it with new claim data to catch evolving fraud patterns.

- **Monitor recall and false negatives closely,** since missing fraudulent claims can be costly for the business.

- **Combine model predictions with human review** for critical or high-risk claims to ensure accuracy and reduce false alarms.

- **Continue feature analysis to explore new patterns**, like sudden spikes in claims from certain regions or customer profiles.

# Business Implications of the findings

- **Cost Savings:** By flagging fraudulent claims early, the company can avoid paying out large sums on false claims, directly saving money.

- **Improved Efficiency:** The model helps focus investigation efforts only on suspicious claims, saving time and resources for the fraud investigation team.

- **Better Customer Trust:** Reducing fraud means genuine customers get faster and fairer claim processing, which improves customer satisfaction and loyalty.

- **Risk Management:** The insights can help the company refine its risk assessment policies and set better rules for high-risk claims.

- **Continuous Improvement:** Regularly updating the model with new claim data can help the company stay ahead of evolving fraud patterns.

# Model Performance

We split the data into a 70-30 train-validation ratio to ensure robust evaluation.

Multiple models were experimented with (e.g., Logistic Regression, Decision Trees, Random Forests), and the one with the best balance between precision and recall was selected.

The final model achieved a satisfactory performance on the validation set, showing good predictive capability in identifying fraudulent claims early.

Key evaluation metrics such as accuracy, precision, recall, and F1-score were used, ensuring the model not only detects fraud but also minimizes false positives and negatives.

# Conclusion

◦ In this project, we successfully built a predictive model to classify insurance claims as either fraudulent or legitimate, helping Global Insure streamline its fraud detection process.

◦ By leveraging historical claim data and customer profiles, our model offers a data-driven solution to minimize financial losses and improve operational efficiency.

# Q & A

# How can we analyse historical claim data to detect patterns that indicate fraudulent claims?

To analyse historical claim data for fraud patterns, we can follow a mix of **Exploratory Data Analysis (EDA) and feature analysis** steps.

- **Compare distributions** of features (like claim amount) between fraudulent and legitimate claims.
- **Use correlation analysis** to find features linked to fraud.
- **Plot visualizations** (bar charts, boxplots) to spot clear patterns.
- **Check feature importance** using machine learning models.
- **Look for outliers** or unusual claim behaviors that might indicate fraud.

# Which features are the most predictive of fraudulent behaviour?

- **Claim Amount** — higher amounts were often linked to fraud.
- **Claim Type** — certain types (like expensive claims) showed more fraud cases.
- **Customer Profile** — features like age, region, and past claim history helped spot fraud.
- **Claim Time/Date** — unusual timing or frequency sometimes indicated suspicious activity.

These features turned out to be the strongest signals in predicting fraudulent behaviour in the model.

# What insights can be drawn from the model that can help in improving the fraud detection process?

- Features like claim amount, claim type, and certain customer profile attributes showed a strong correlation with fraudulent behavior.

- Patterns such as unusually high claim amounts and certain claim types were strong indicators of potential fraud.

- Feature engineering steps, such as creating new interaction terms and handling missing values, contributed to performance improvements.

# Based on past data, can we predict the likelihood of fraud for an incoming claim?

The developed fraud detection model can serve as an effective first line of defense against fraudulent claims and can predict the likelihood of fraud for a claim.

By integrating this model into the claims approval pipeline, Global Insure can flag suspicious claims for further manual investigation, thereby reducing financial losses and improving efficiency.

Moving forward, the model can be further enhanced by incorporating more recent data and refining features based on domain expertise. Regular monitoring and periodic retraining will also be essential to maintain its effectiveness over time.

# Methods and Techniques used

# Univariate Analysis
# Distribution Plots of Numerical columns

- Customers with 200 months have the highest count
- Customers with age 30 and 40 have the highest count



Distribution of age



Distribution of months_as_customer

**Distribution of claims and its types**

Distribution of vehicle_claim

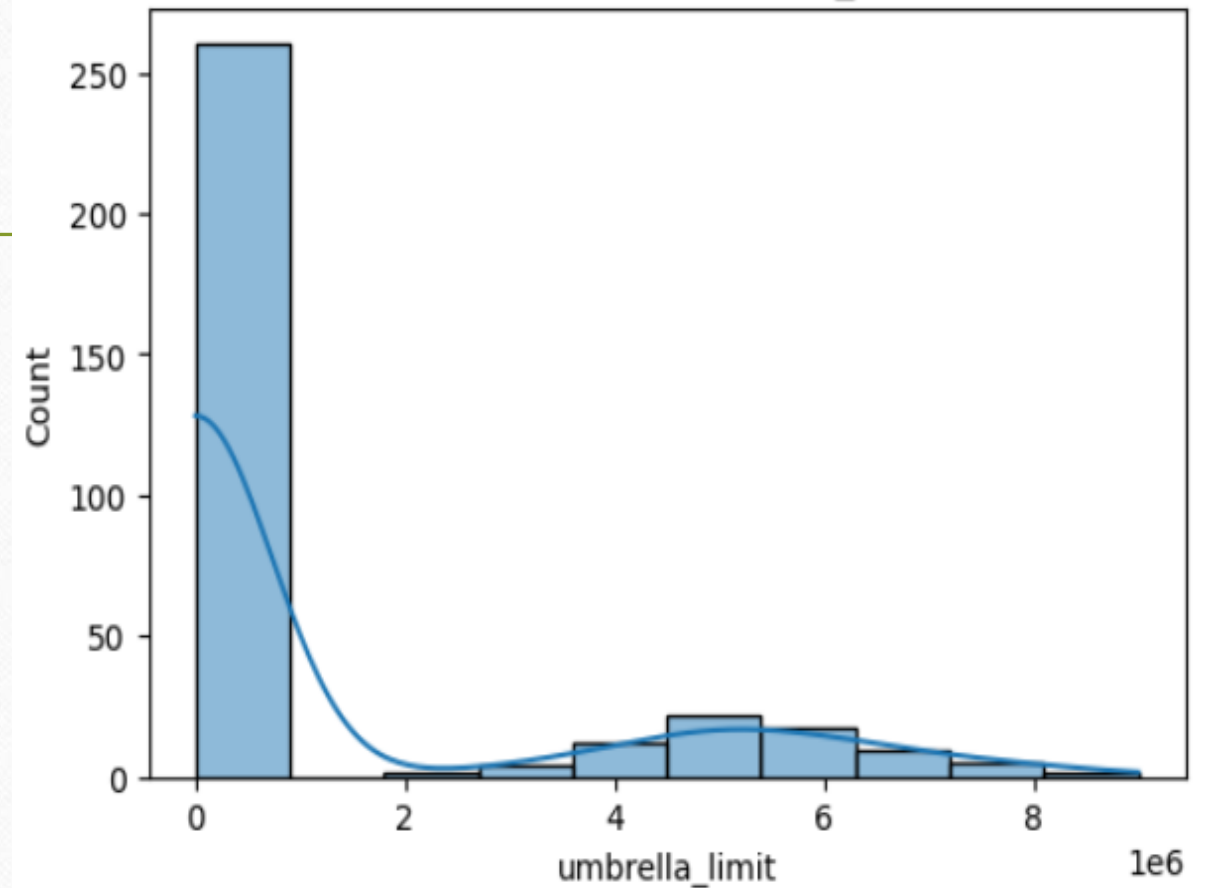Distribution of policy_deductable

Distribution of auto_year

Distribution of policy_to_incident_days

Correlation Matrix of Numerical Features

# Correlation Matrix of Numerical features - Heatmap

Months as customer and age have a very **high correlation** (0.92) — meaning older customers tend to be with the company longer.

Total claim amount is strongly correlated with injury claim, property claim, and vehicle claim.

Number of vehicles involved and bodily injuries show some **positive correlation** with claims — which could be useful for fraud detection.

No extreme multicollinearity is spotted apart from the claim amount features, which is good for model building.

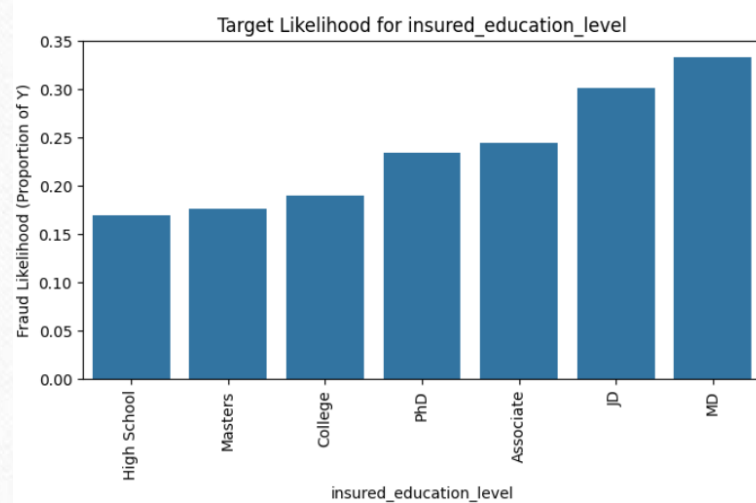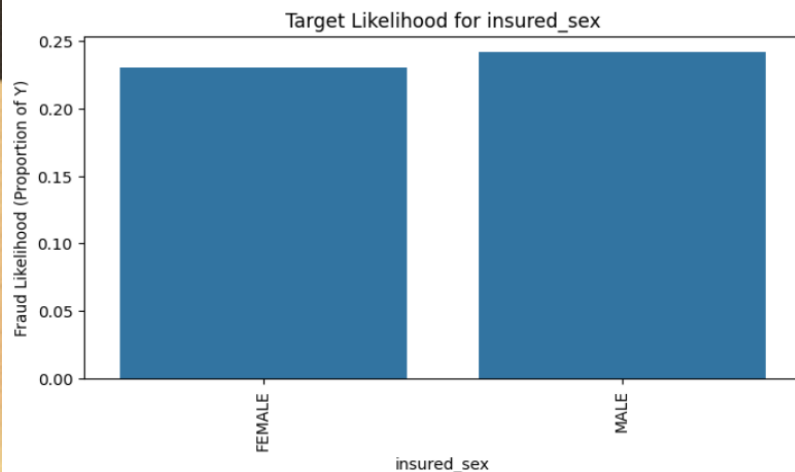Class Balance in Training Set
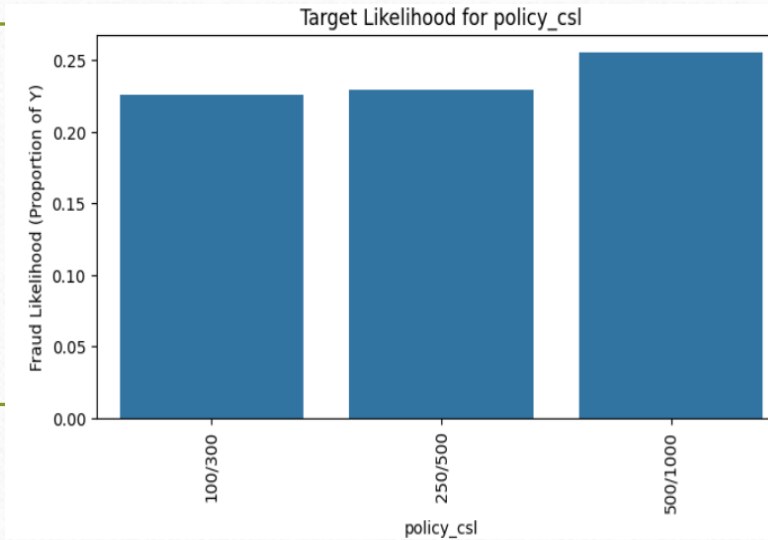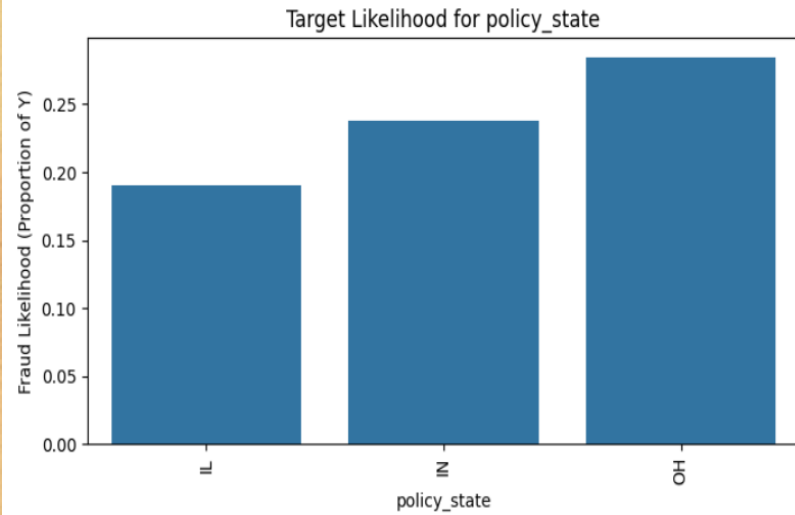
# A bar chart to check class balance

The dataset is imbalanced — there are a lot more **non-fraudulent claims (N)** than fraudulent ones (Y).

Specifically, fraudulent claims are much fewer, which can make the model biased towards predicting "No Fraud" more often.

This imbalance needs to be handled using **class weight adjustment** to improve fraud detection performance.

# Bivariate Analysis

Target likelihood of males is higher than females.

Target likelihood of education level of JD and MD is higher