# Title: Fraudulent Insurance Claim Detection

**Submitted by:** Saniya Baig and Chirag Panchal

**Course/Module:** Machine learning

**Date:** 11th May 2025

# Table of Contents

- Problem Statement

- Objective

- Dataset Overview

- Methodology

- Exploratory Data Analysis (EDA)

- Feature Engineering

- Model Building and Evaluation

- Key Insights and Visualizations

- Assumptions

- Recommendations and Conclusion

# Problem Statement

➢ Insurance fraud is a major challenge impacting the financial performance of insurance companies. Detecting fraudulent claims is crucial to prevent unnecessary payouts and improve operational efficiency.

➢ This project focuses on building a classification model to identify whether an insurance claim is fraudulent or legitimate based on historical data.

# Objective

The primary objective of this assignment is to:

- Analyze historical insurance claims data.

- Identify patterns indicative of fraudulent behavior.

- Build and evaluate classification models for fraud detection.

- Compare model performances and recommend the best approach.

# Dataset Overview

The dataset consists of 1000 insurance claims with 40 features, including:

- Customer demographics

- Policy information

- Claim details

- Incident descriptions

- Target variable: fraud_reported (Yes/No)

# Methodology

•Data Cleaning & Preprocessing: Handling missing values, encoding categorical variables.

•Exploratory Data Analysis: Visual exploration of fraud patterns.

•Feature Engineering: Creation of relevant dummy variables and selection of significant features.

•Modeling: Logistic Regression and Random Forest models were trained.

•Evaluation: Based on accuracy, precision, recall, specificity, sensitivity, F1-score.

# Exploratory Data Analysis (EDA)

Key findings:

- Fraudulent claims were around 20% of total data.

- Certain states and occupations showed higher fraud rates.

- Visual trends observed in claim amounts and incident types.

- ➡ Insert bar plots or pie charts showing fraud distribution, state-wise incidents, etc.
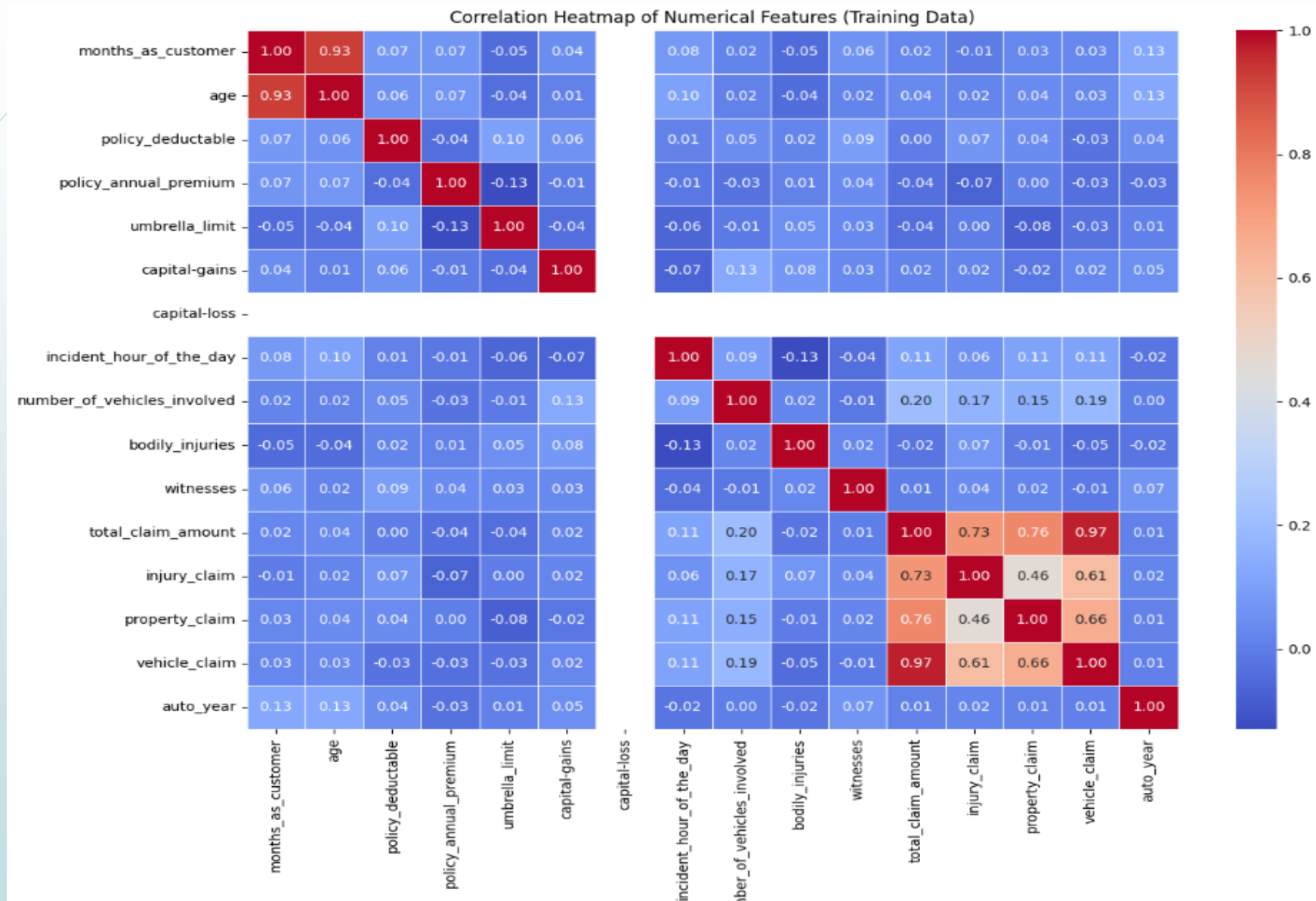
Class Distribution in Training Set

Class distribution:
 fraud_reported
N    225
Y     76
Name: count, dtype: int64

➧ The target variable is imbalanced, with significantly more non-fraudulent claims (around 73%) than fraudulent ones (about 27%). This may affect model performance, particularly recall for the minority class (fraud). Handling techniques such as resampling or class weighting may be considered in later steps.

Given the imbalance in the fraud_reported variable (approx. 73% No vs. 27% Yes), we will initially retain the original distribution for EDA and baseline model development. To mitigate bias during training, class weights will be applied. If performance on fraudulent claims remains poor, oversampling techniques such as SMOTE may be applied to the training data.

# Correlation heatmap of numerical features



Correlation Heatmap of Numerical Features (Training Data)

**Correlations Identified:**

1. **Very strong correlations:**

   ○ csl_per_person ↔ csl_per_accident → **0.99**

   ○ vehicle_claim ↔ total_claim_amount → **0.98**

   ○ property_claim ↔ total_claim_amount → **0.77**

   ○ injury_claim ↔ total_claim_amount → **0.76**

2. **Strong pairwise correlations:**

   ○ vehicle_claim ↔ injury_claim → **0.67**

   ○ vehicle_claim ↔ property_claim → **0.67**

   ○ property_claim ↔ injury_claim → **0.47**

3. **Other notable one:** months_as_customer ↔ age → **0.92** → This is unexpectedly high and worth checking if there's data leakage or a data engineering artifact.
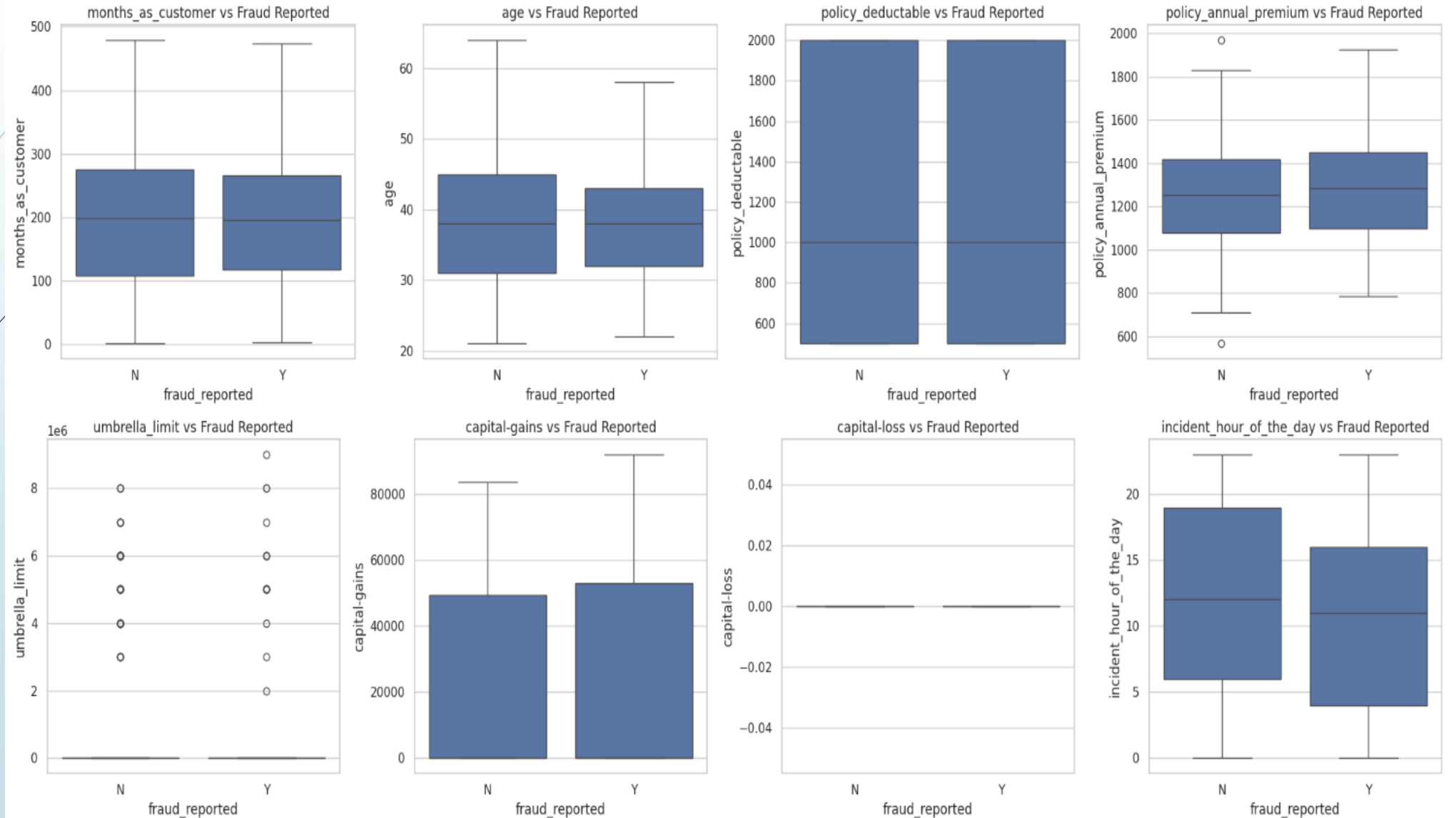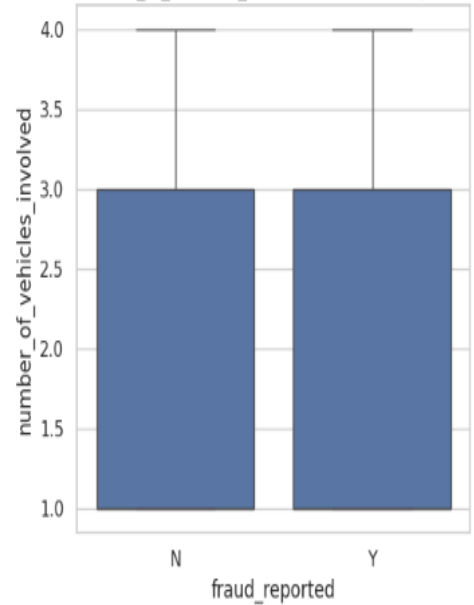
# Feature Engineering

• Created dummies for multi-category variables like policy_csl, incident_severity.

• Selected top features using p-values and correlation with the target.

• Final selected features included:

  • policy_csl_incident_severity_250/500_Major Damage

  • policy_csl_incident_severity_500/1000_Major Damage

  • insured_occupation_handlers-cleaners
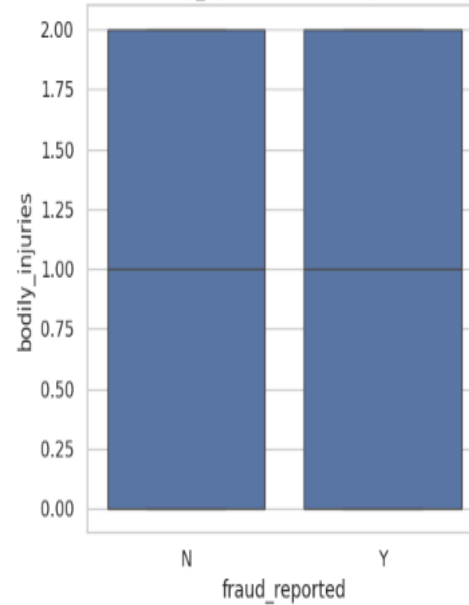
  • incident_state_OH

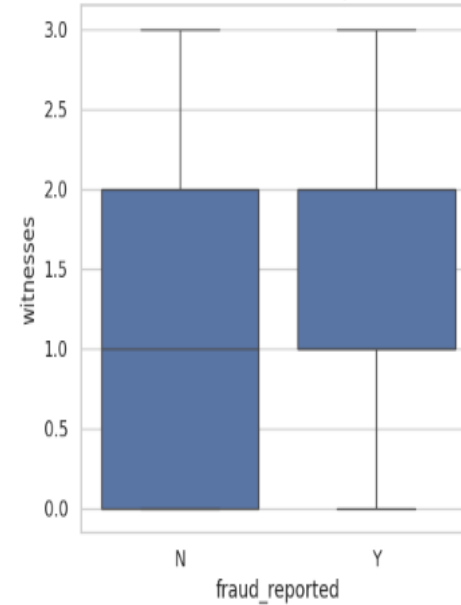# Box plot of numerical features vs Target variable

1. **Strong Predictors**:

    o **total_claim_amount**, **witnesses**: Large differences between fraudulent and non-fraudulent claims, particularly with witnesses showing a **clear separation**. These features could be strong indicators of fraud.

2. **Moderate Predictors**:

    o **capital_loss**, **vehicle_claim**, **property_claim**: These features show **moderate differences** in medians and include outliers, which may suggest they have some predictive power but require further investigation.

3. **Weak Predictors**:

    o **auto_year**, **policy_annual_premium**: These features show only **slight differences** in medians between fraud and non-fraud, making them weaker predictors for the target variable.

4. **Outliers**:

    o Several features, including **vehicle_claim**, **property_claim**, **total_claim_amount**, and **umbrella_limit**, have **outliers**

# Model Building and Evaluation

- **Model 1: Logistic Regression**

  - Accuracy: 79.2%

  - F1-Score: 0.526

  - Good balance of recall and precision.

- **Model 2: Random Forest**

  - Accuracy: 74.6%

  - F1-Score: Poor (due to overfitting, only predicted one class)

  - Despite high training performance, failed on validation.

➡ **Add confusion matrices and a comparison table of metrics here.**

## Model Performance Comparison on Test Data

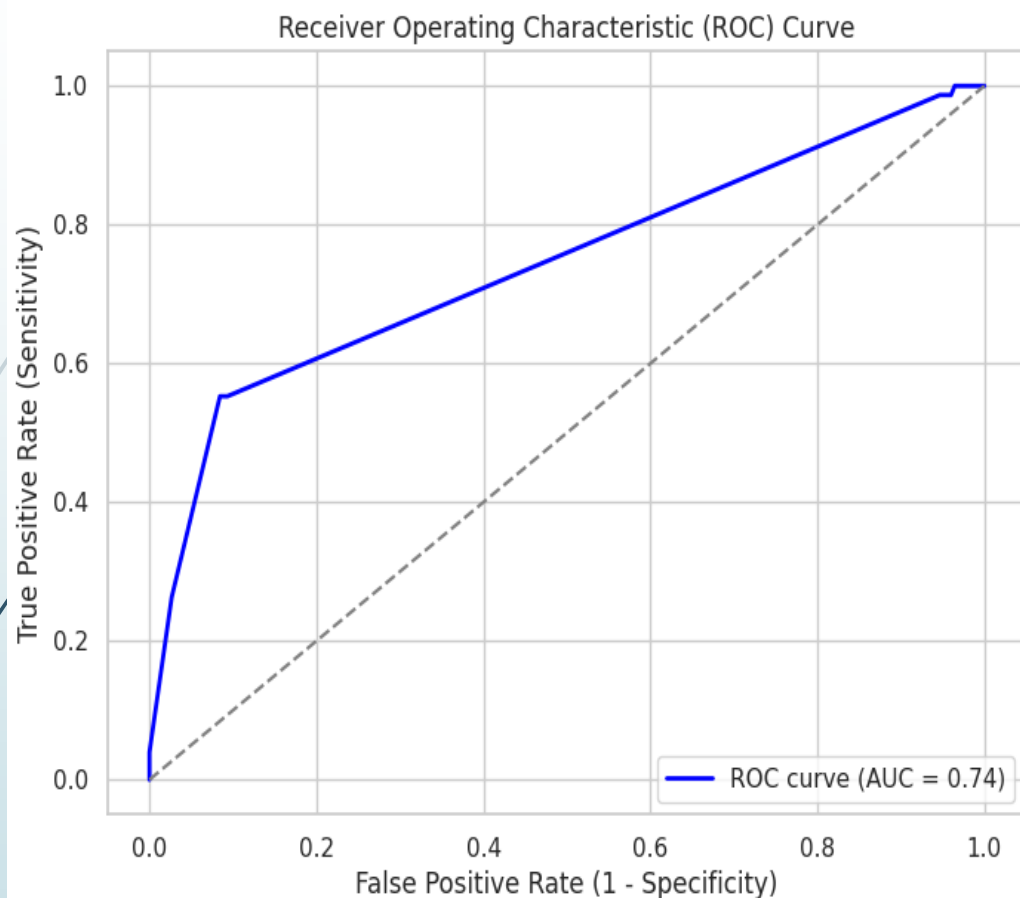| Model | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| Logistic Regression | 0.79 | 0.625 | 0.45 | 0.53 |
| Random Forest | 0.75 | 0.00 | 0.00 | 0.00 |

**Random Forest** in terms of **precision, recall, and F1 score** on the test data. Random Forest failed to identify any positive (fraudulent) cases,
leading to zero values for all key classification metrics. This indicates that it is either **overfitting** the majority class or **completely biased** toward predicting
only legitimate claims. On the other hand, Logistic Regression demonstrates a much better balance, correctly identifying a reasonable portion of fraud cases
while maintaining decent precision. Therefore, **Logistic Regression is the preferred model** for this task, especially given the class imbalance and the business importance of detecting fraud.
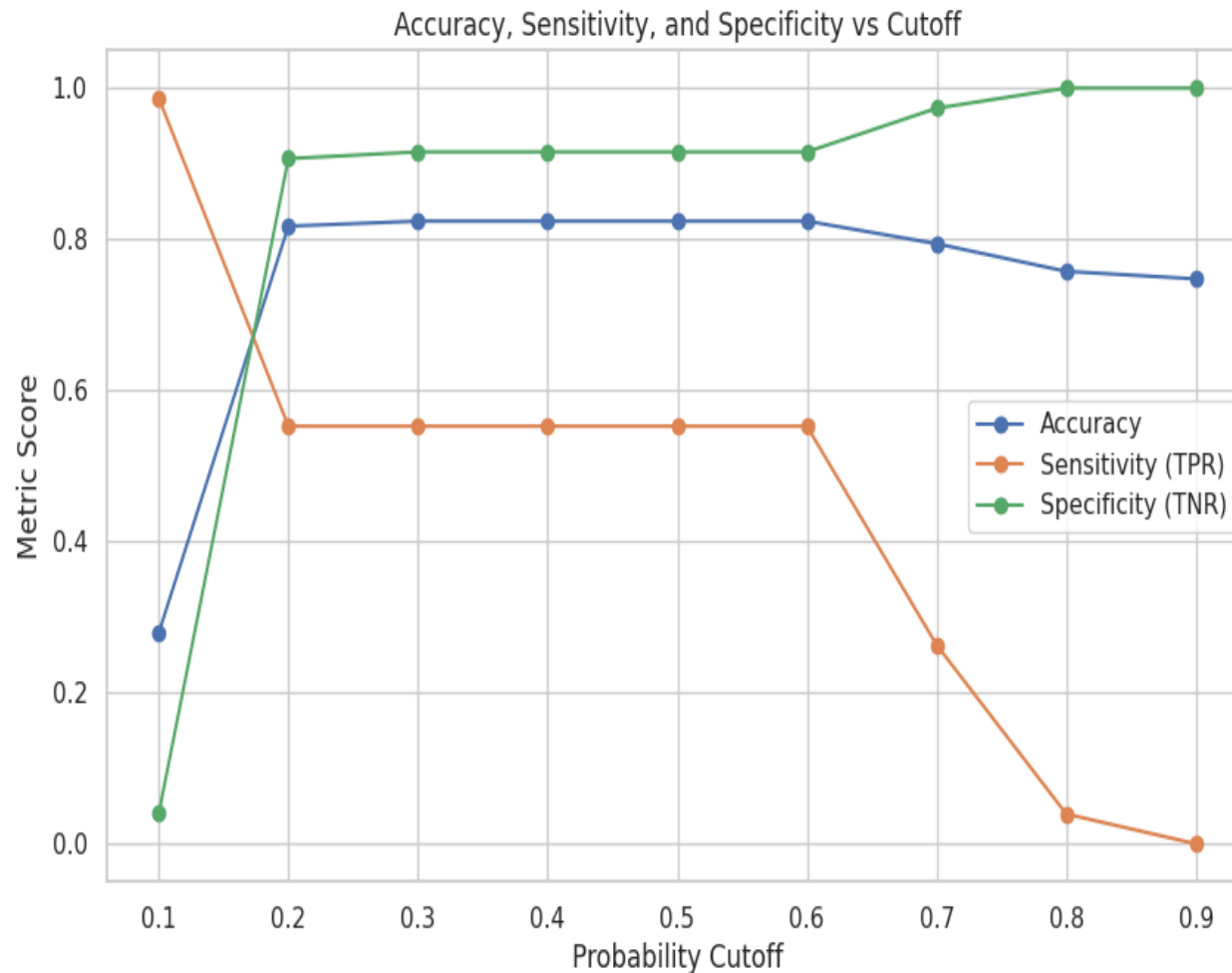
# Key Insights and Visualizations

- Logistic regression outperformed Random Forest in balanced performance.

- Random Forest overfitted on training data.

- Features involving claim severity and occupation are strong fraud predictors.

# ROC Curve



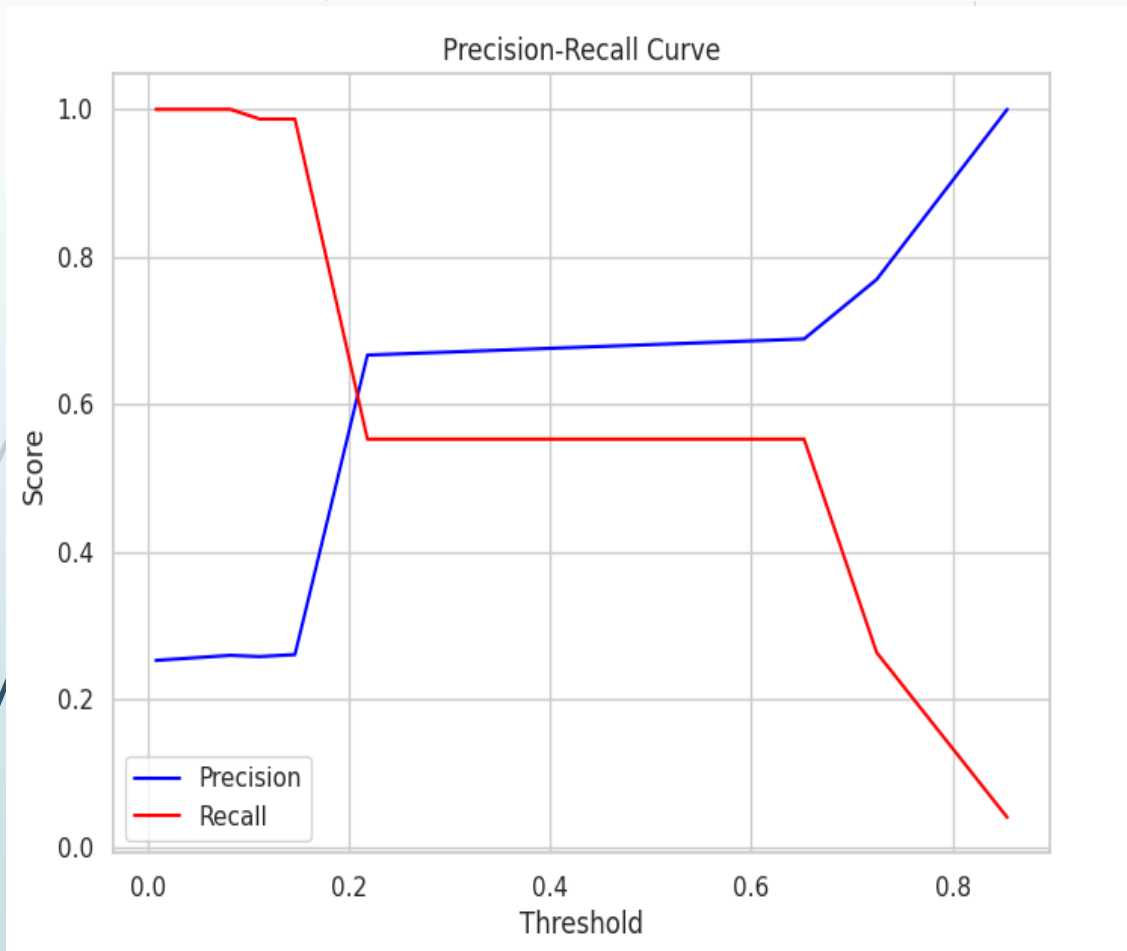Receiver Operating Characteristic (ROC) Curve

The ROC curve demonstrates that the logistic regression model has a fair ability to distinguish between fraudulent and legitimate claims, with an AUC of 0.74. This indicates a 74% probability that the model ranks fraudulent cases higher than legitimate ones. While this is a respectable performance for an initial model, especially in the presence of class imbalance, there is potential to improve the model's discriminatory power, possibly by exploring more complex algorithms, rebalancing techniques, or threshold tuning

# Accuracy, sensitivity, specificity at different values of probability cutoffs



Accuracy, Sensitivity, and Specificity vs Cutoff

The cutoff analysis reveals that a threshold between 0.2 and 0.3 offers the best balance between sensitivity and specificity. At this range, the model maintains high accuracy (~82%), captures a reasonable portion of fraud cases (sensitivity ~55%), and correctly identifies most legitimate claims (specificity ~91%). Beyond this range, increasing the threshold improves specificity but sharply reduces the model's ability to detect fraud, making it unsuitable for a fraud detection scenario where recall is crucial

# Precision-recall curve



Precision-Recall Curve

- There's a clear **trade-off** between precision and recall.

- The **threshold range around 0.2 to 0.6** appears to give a **reasonable balance**:
  - Precision ~0.66
  - Recall ~0.55

- **Beyond 0.6**, precision improves, but **recall collapses**, making it less useful for detecting fraud.

As the threshold increases, precision improves but recall drops significantly. An optimal balance occurs around a threshold of 0.2 to 0.3, where both precision and recall are reasonably high (~0.66 and ~0.55, respectively), maximizing the F1 score. Beyond this range, recall drops too low, compromising the model's ability to detect fraud effectively."

# Assumptions

- Validation data has similar distribution and features as training data.

- Cost of false negatives (fraud missed) is higher than false positives.

- The cutoff threshold for logistic regression was optimized based on F1-score.

# Recommendations and Conclusion

- Use logistic regression for deployment due to better generalization.

- Focus fraud investigations on high-risk occupations and states.

- Implement a periodic retraining pipeline as fraud patterns evolve.

➡ Final Thought: A simple, interpretable model like logistic regression can offer reliable fraud detection when combined with domain insights.