

# Claim fraud detection report

## Introduction

The report will include the overall approach of the assignment, covering the problem statement, methodology, techniques used and key insights.

## Overall Approach

The overall approach was to do analysis, clean-up, modelling and predictions. Follow various steps learnt so far in ML journey.

The objective is to build a model to classify insurance claims as either fraudulent or legitimate based on historical claim details and customer profiles. By using features such as claim amounts, customer profiles, claim types and approval times, the company aims to predict the claims that are likely to be fraudulent before they are approved.

## Methodology

Steps to solve the problem:

### Logistics regression

- Raw Data
- Data cleaning : handle missing values, outliers
- Train-Test Split
- Encoding categorical features : get\_dummies
- Scale numeric features
- Features selection using RFECV
- Resampling
- Train final model
- Predict
- Evaluate performance : accuracy, F1, AUC, confusion matrix

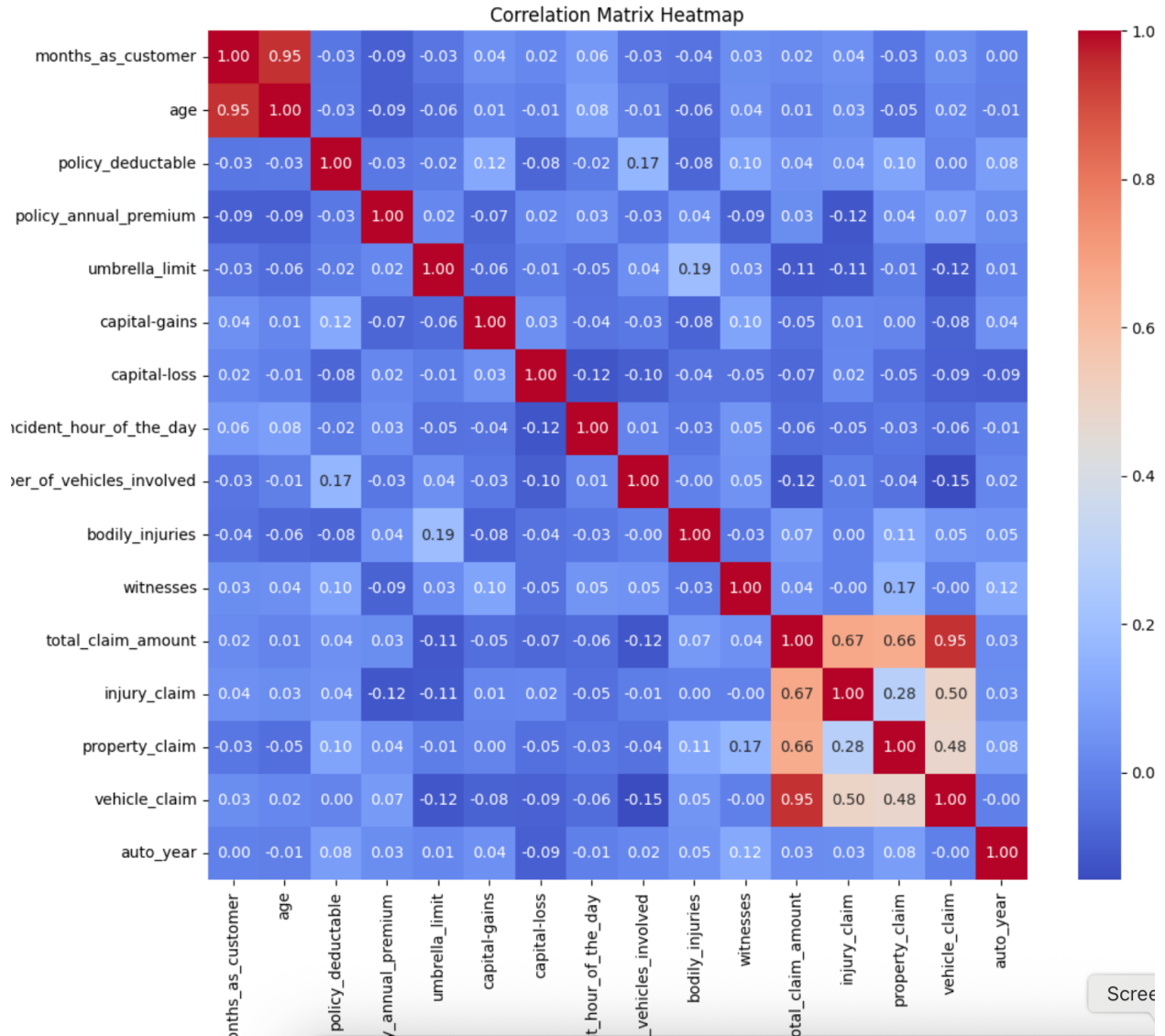
## Random Forrest

- Raw Data
- Data cleaning : handle missing values, outliers
- Train-Test Split
- Encoding categorical features : get\_dummies
- Features selection using feature importance
- Resampling
- Train final model
- Predict
- Evaluate performance : accuracy, F1, AUC, confusion matrix

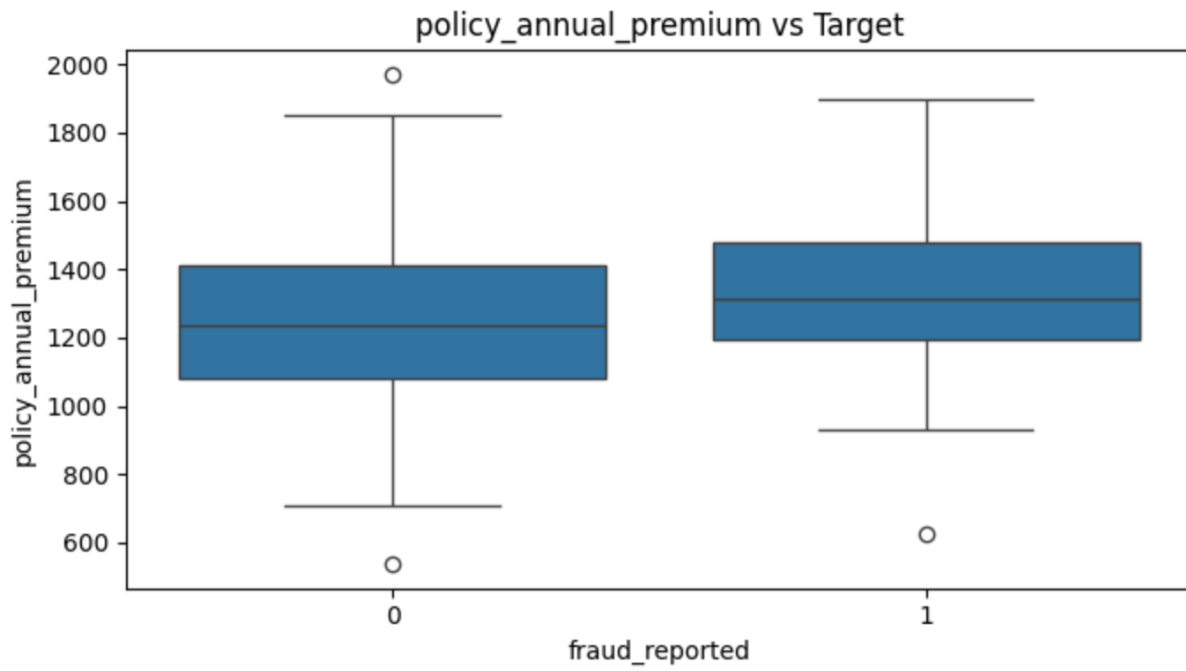
## Key Insights

It was lengthy & tough assignment. It has to be done carefully considering so many steps. Mixing the steps can create much different results. Metrics used to find different aspects of data, model & prediction:

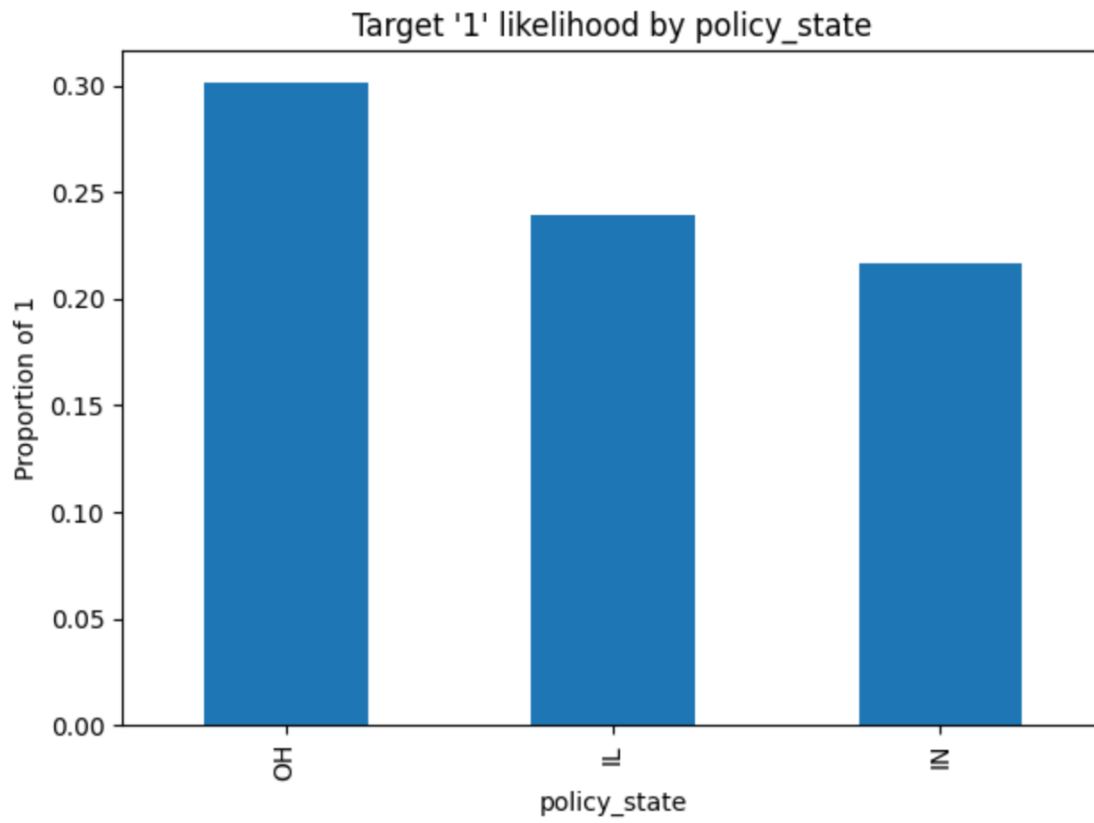
## Heatmap



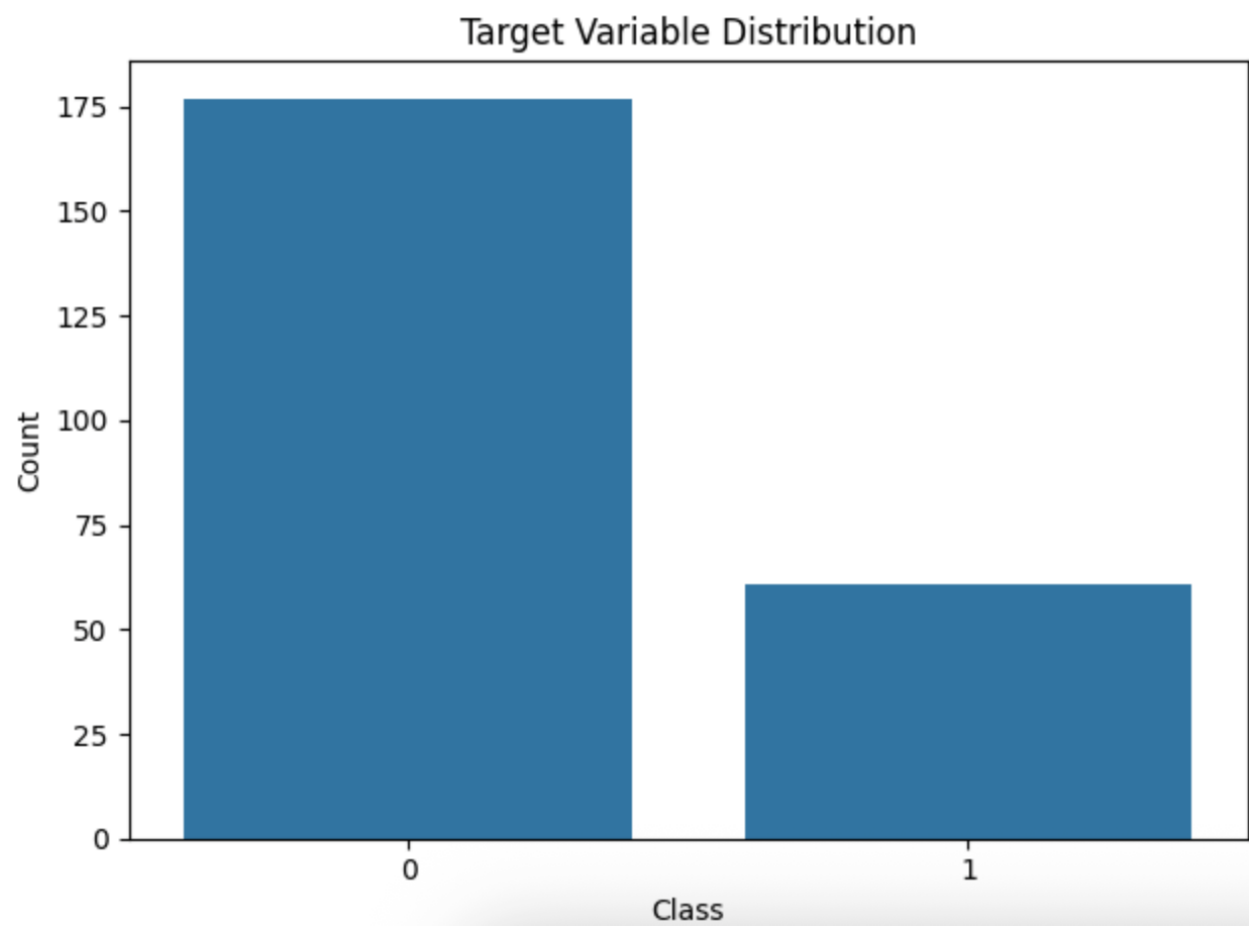
Relationship between numerical & target variable



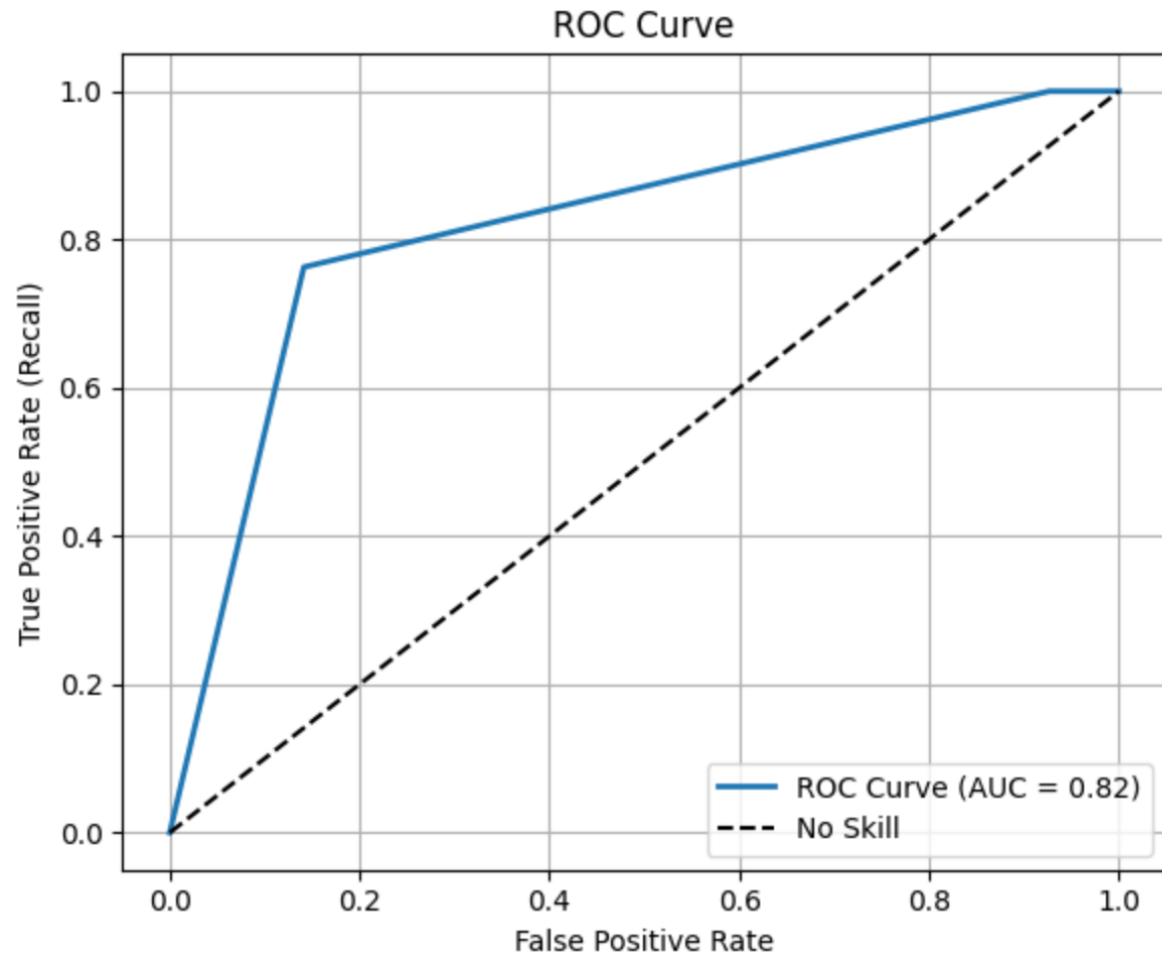
## Categorical feature likelihood



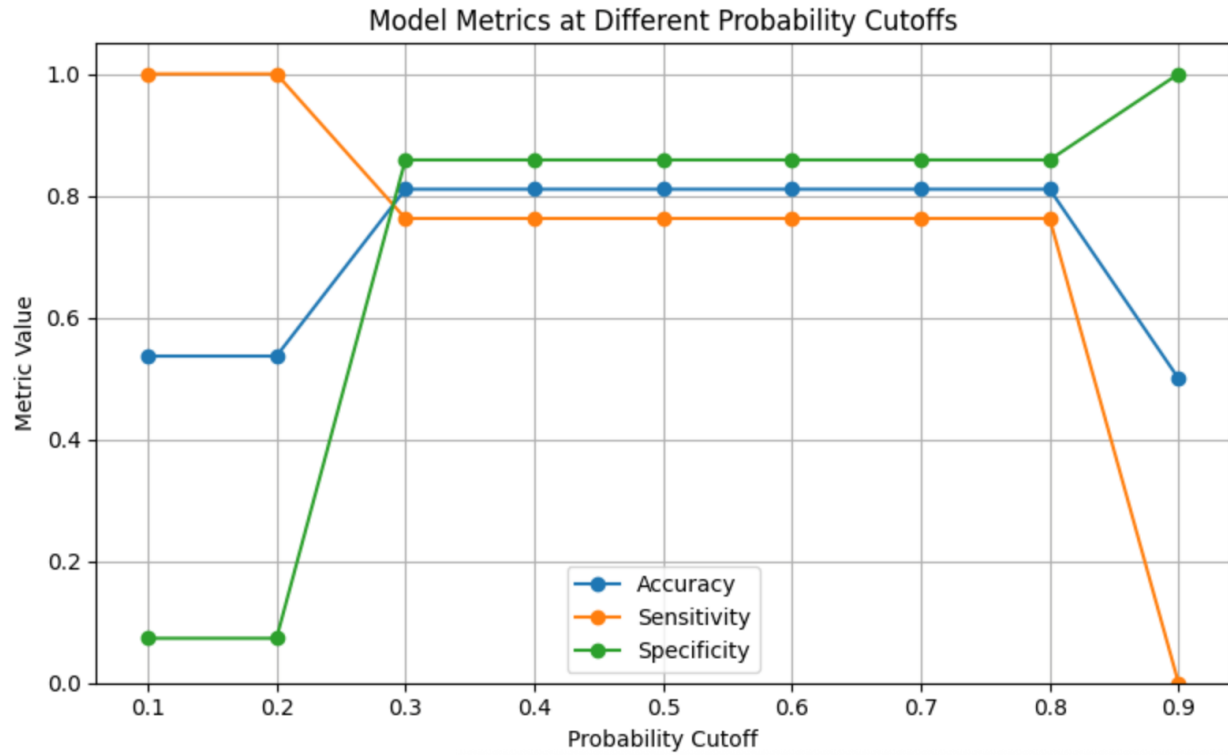
## Class balance



## ROC Curve



## Probability Cutoff



## Precision-Recall Tradeoff

