

# Insurance Fraudulent Claim Detection: Analytical Report

## 1. Executive Summary

Insurance fraud is a significant challenge for the industry, leading to substantial financial losses and increased operational costs. This report presents a comprehensive data-driven approach to detecting fraudulent insurance claims using advanced machine learning techniques. By leveraging historical claim data, we aim to build robust models that can accurately identify potential fraud, thereby supporting more efficient and effective claims processing.

## 2. Problem Statement

Global Insure, a leading insurance provider, faces the persistent issue of fraudulent claims, which not only result in direct financial losses but also undermine customer trust and inflate premiums for honest policyholders. The current manual review process is resource-intensive and not scalable. The primary objective of this project is to develop a predictive model that can flag potentially fraudulent claims at an early stage, enabling targeted investigations and reducing overall fraud-related losses.

## 3. Dataset Description

The dataset comprises anonymized records of insurance claims, each labelled as fraudulent or legitimate. It includes a diverse set of features capturing customer demographics, policy details, incident characteristics, and claim amounts.

### Key Features:

- **Customer Profile:** Age, education, occupation, relationship to policyholder, hobbies
- **Policy Details:** Policy number, bind date, state, deductible, annual premium, umbrella limit
- **Incident Details:** Date, type, severity, location, hour, number of vehicles, property damage, bodily injuries, witnesses, police report
- **Claim Details:** Total claim amount, sub-claims (injury, property, vehicle), capital gains/losses
- **Vehicle Information:** Make, model, year
- **Target Variable:** fraud\_reported (Yes/No)

### Data Preprocessing:

- Removed high-cardinality identifiers and redundant columns
- Engineered new features (e.g., claim severity ratio, policy duration category)
- Encoded categorical variables and scaled numerical features

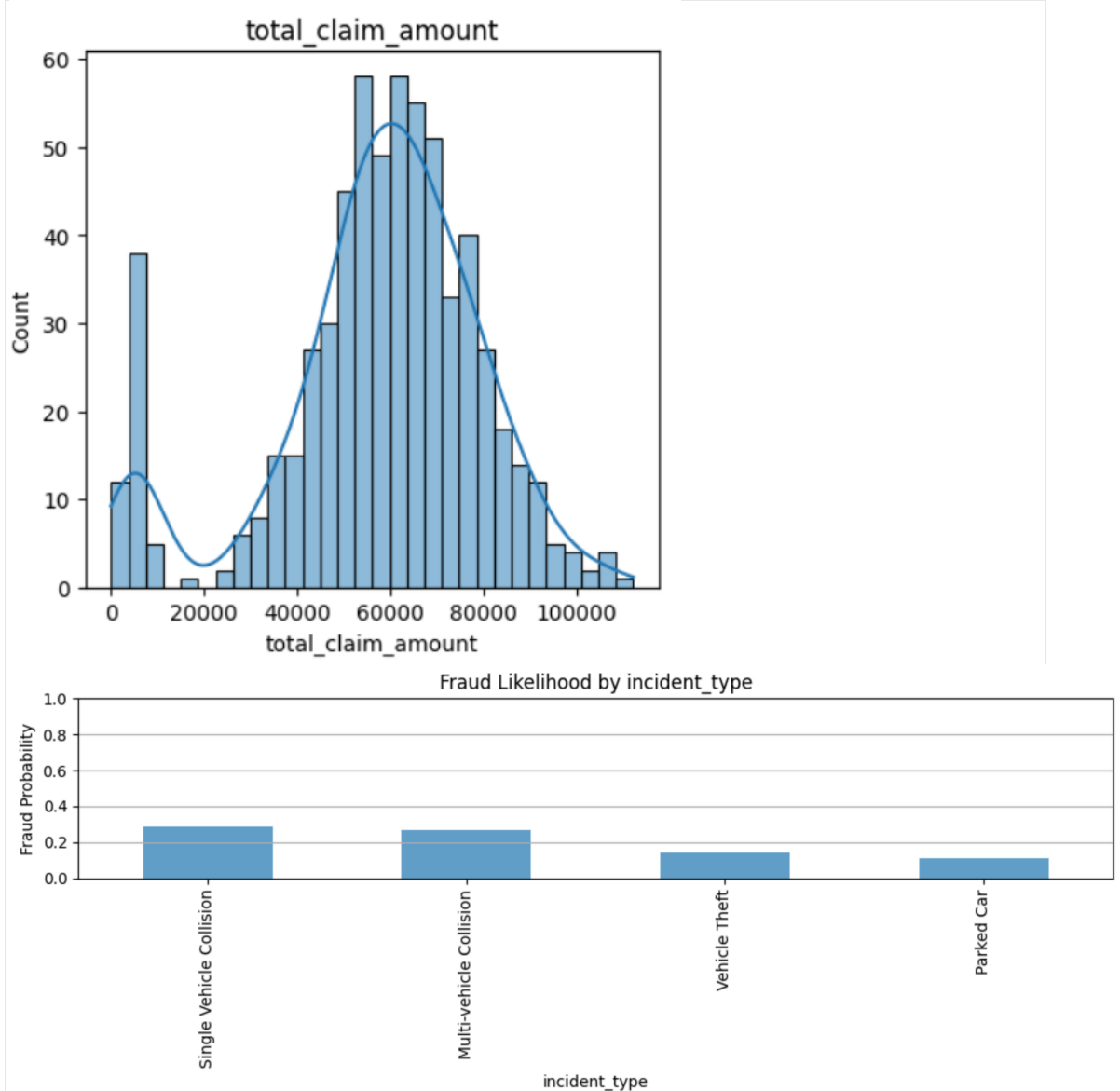
## 4. Exploratory Data Analysis (EDA)

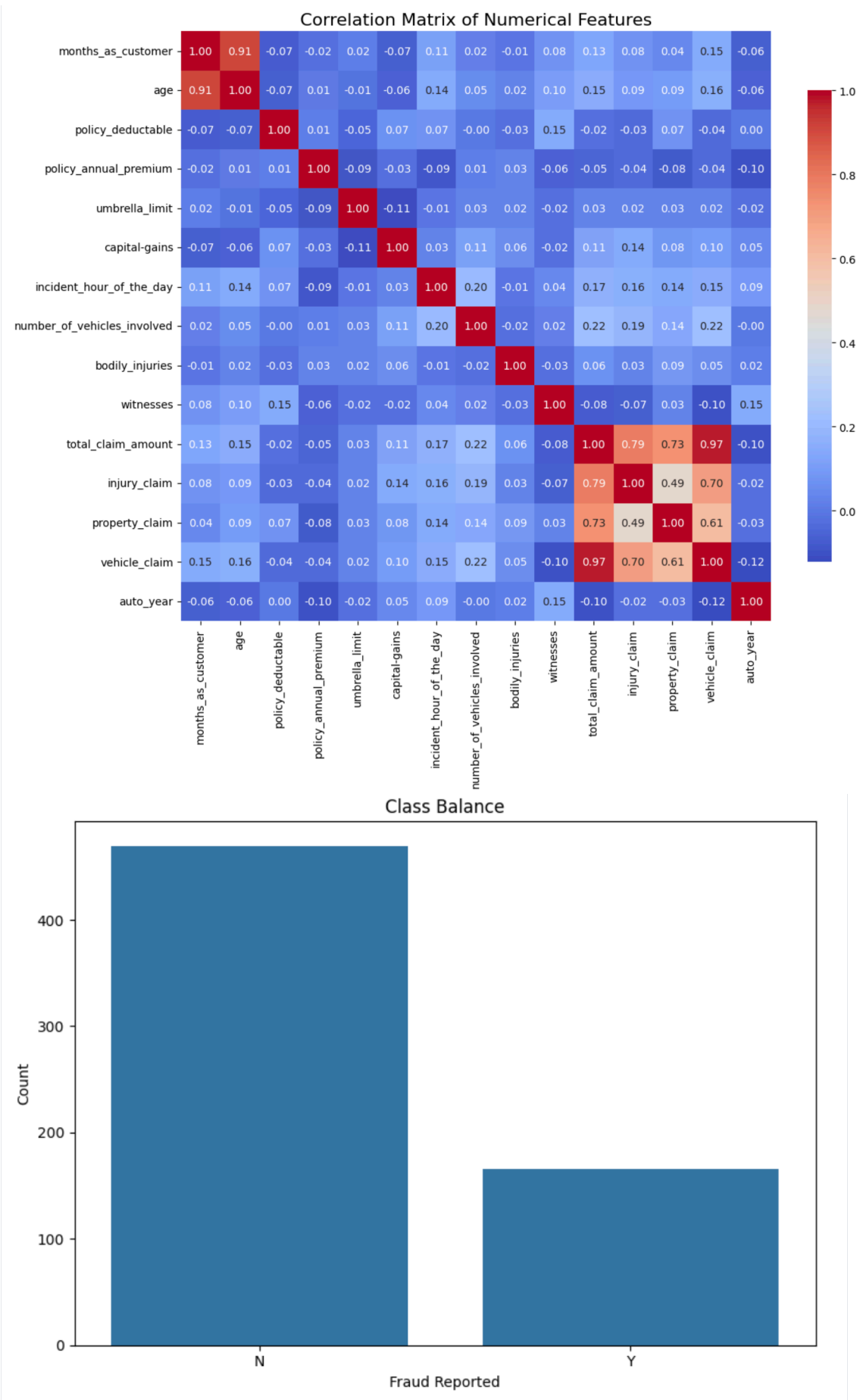
The EDA phase provided key insights into the data distribution, feature relationships, and potential predictors of fraud.

### Key Observations:

- **Claim Amounts:** Right-skewed distribution; fraudulent claims tend to have higher average amounts.
- **Incident Types:** Certain types (e.g., multi-vehicle collisions) have higher fraud rates.
- **Class Balance:** The dataset is moderately imbalanced, with fewer fraudulent claims.
- **Correlation Analysis:** Most features show low to moderate correlation, reducing multicollinearity concerns.

### Visualizations:





## 5. Modelling Approach

### 5.1 Logistic Regression

#### Definition:

A statistical method for binary classification that models the probability of a claim being fraudulent as a function of input features. Logistic regression is valued for its interpretability and ability to provide clear explanations for predictions.

#### Feature Selection:

Recursive Feature Elimination with Cross-Validation (RFECV) was used to retain only the most predictive features.

#### Multicollinearity Check:

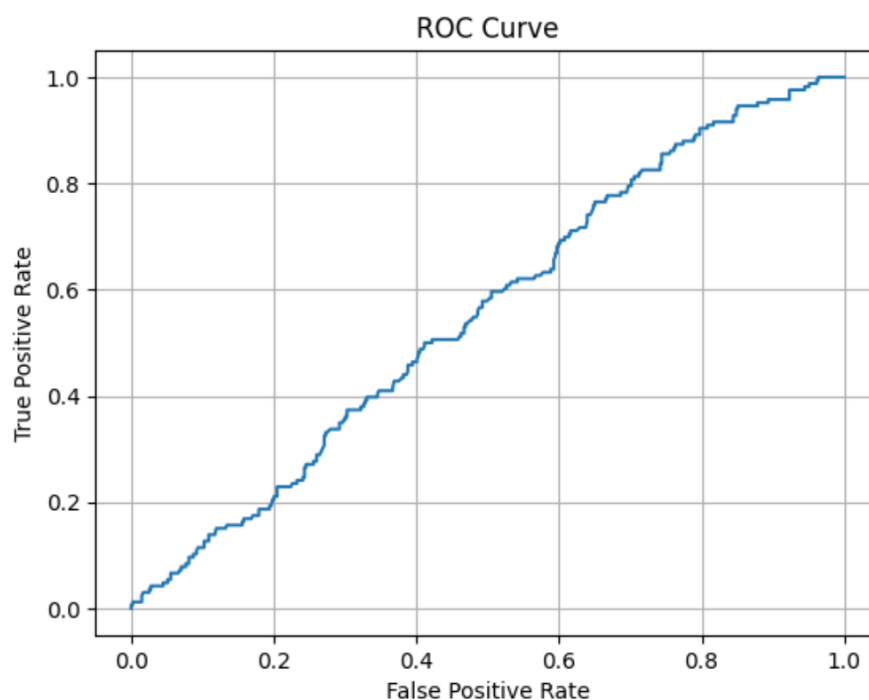
Variance Inflation Factor (VIF) analysis confirmed all selected features had  $VIF < 2$ , ensuring reliable coefficient estimates.

#### Model Evaluation:

- At the optimal probability cutoff (e.g., 0.27–0.3):
  - **Sensitivity (Recall):** Up to ~79% (crucial for catching fraud)
  - **Specificity:** ~40–60% (acceptable trade-off)
  - **F1 Score:** Balanced, reflecting effective fraud detection
- **Interpretability:**  
Coefficients provide actionable insights into which features most influence fraud predictions.

#### Confusion Matrix:

```
[[405  64]
 [140  26]]
```



## 5.2 Random Forest

### Definition:

An ensemble learning method that builds multiple decision trees and aggregates their predictions. Random Forests are powerful for capturing complex, nonlinear relationships and interactions between features.

### Feature Importance:

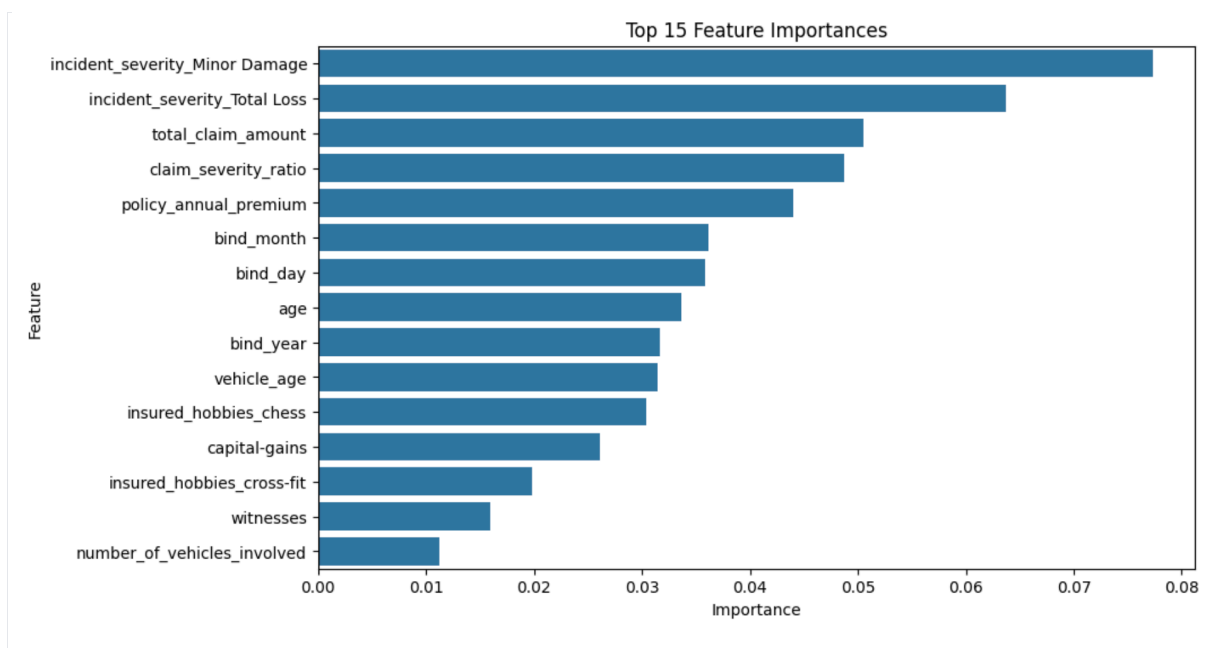
Top features included engineered ratios, claim amounts, and select categorical encodings.

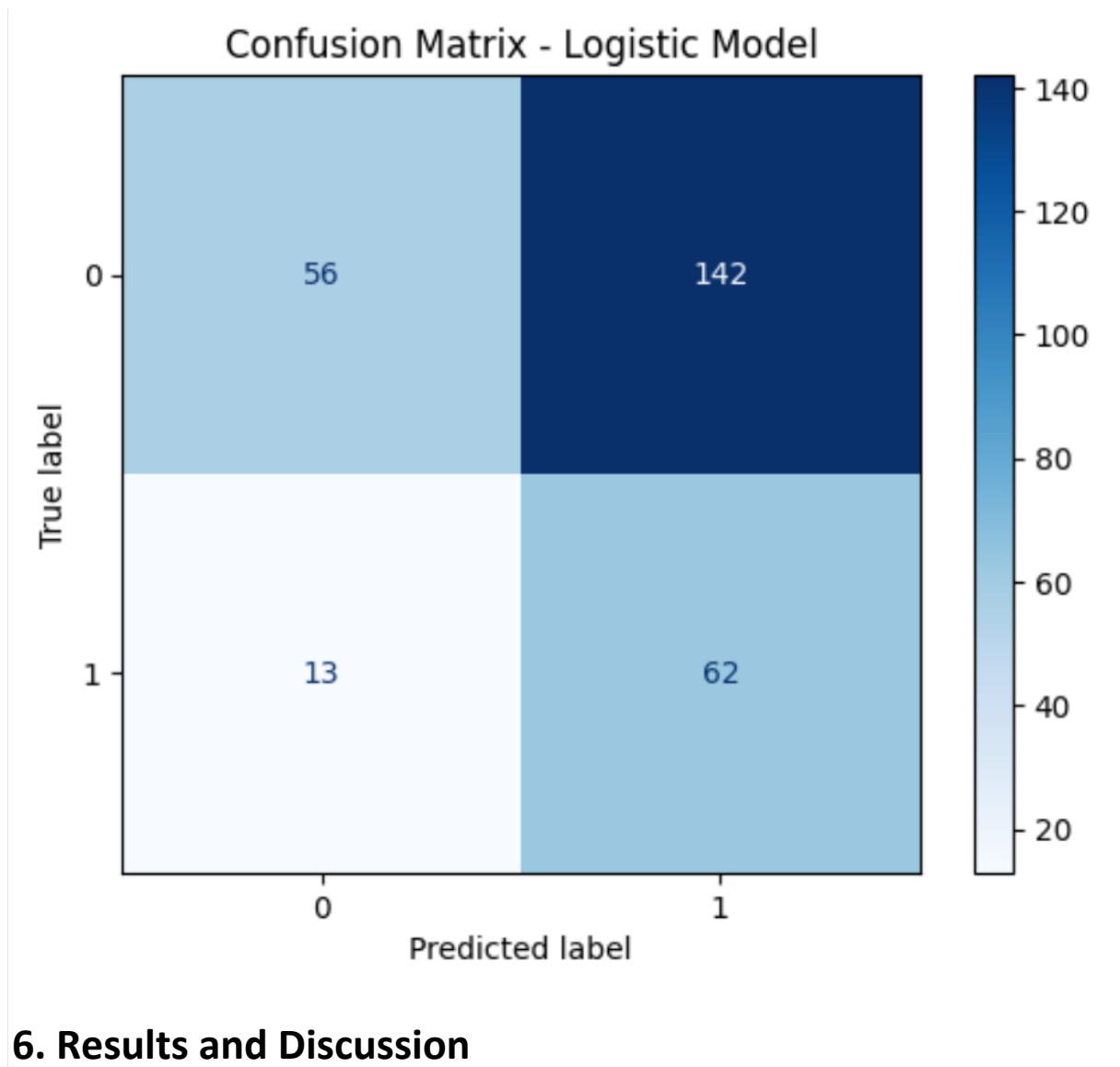
### Model Tuning:

GridSearchCV was employed to optimize hyperparameters, improving model accuracy and robustness.

### Model Evaluation:

- **Accuracy:** High, with strong overall predictive power
- **Sensitivity (Recall):** High, generally outperforming logistic regression
- **F1 Score:** Improved, indicating better balance between precision and recall
- **Interpretability:**  
While less transparent than logistic regression, feature importance rankings offer some insight into model decisions.





Metric	Logistic Regression	Random Forest
Accuracy	0.6787	0.4322
Recall	0.1566	0.8267
Specificity	0.8635	0.2828
F1 Score	0.2031	0.4444

### Key Findings:

- **Random Forest** achieved higher recall and F1 scores, making it more effective for maximizing fraud detection.
- **Logistic Regression** offers greater interpretability, which is valuable for regulatory compliance and stakeholder trust.
- Both models benefit from careful feature selection, multicollinearity checks, and threshold tuning.

## 7. Conclusions and Recommendations

- **For maximum fraud detection (high recall):** Deploy the Random Forest model, especially after hyperparameter tuning and threshold adjustment.
- **For interpretability and actionable insights:** Use Logistic Regression, particularly when clear explanations are required for business decisions or regulatory reasons.
- **Operational Considerations:**
  - Regularly retrain models with new data to adapt to evolving fraud patterns.
  - Combine model predictions with expert review for high-value or ambiguous cases.
  - Monitor model performance and recalibrate thresholds as needed to balance recall and precision.

### Prepared by:

*Dewang Shishodia*

*Aditya Peesapati*