

# Cybersecurity Incident Classification Using Machine Learning

## 1. Introduction

Security Operation Centers (SOCs) process a large volume of cybersecurity incidents daily. This project utilizes machine learning to automate incident classification, helping analysts prioritize threats effectively. The dataset used is the GUIDE dataset, which categorizes cybersecurity incidents into True Positive (TP), Benign Positive (BP), and False Positive (FP) classes.

## 2. Data Preprocessing

### 2.1 Data Exploration & Cleaning

Before feeding the data into machine learning models, thorough preprocessing steps were conducted:

- **Handling Missing Values:** Missing values were imputed using median values for numerical features and the most frequent category for categorical features.
- **Feature Engineering:** New features were created based on domain knowledge, such as aggregated log counts and timestamps.
- **Encoding Categorical Variables:** One-hot encoding was used for nominal features, while label encoding was applied to ordinal features.

### 2.2 Handling Imbalanced Data

Cybersecurity datasets are often imbalanced, with certain attack types being rarer than others. This project addressed the imbalance using:

- **Synthetic Minority Over-sampling Technique (SMOTE):** SMOTE generated synthetic examples for underrepresented classes.
- **Class Weighting:** Models were trained with adjusted class weights to reduce bias towards majority classes.

## 2.3 Data Splitting

The dataset was divided to ensure generalization and robust model evaluation:

- **Training Set:** 70% of the dataset used to train the models.
- **Validation Set:** 20% used for hyperparameter tuning.
- **Test Set:** 10% held out for final performance evaluation.

## 3. Model Selection & Training

Several machine learning models were trained and evaluated:

### 3.1 Logistic Regression

- A linear model used as a baseline.
- Struggled with class imbalance, leading to lower recall for minority classes.

### 3.2 Decision Tree

- Non-linear model capable of capturing complex decision boundaries.
- Prone to overfitting, especially on high-dimensional datasets.

### 3.3 Random Forest

- An ensemble of multiple decision trees.
- Reduced overfitting but computationally expensive.

### 3.4 Gradient Boosting & XGBoost

- Boosting algorithms that iteratively improve weak learners.
- XGBoost outperformed other models due to its ability to handle complex patterns.

### 3.5 LightGBM

- Faster than XGBoost while maintaining similar accuracy.
- Well-suited for large-scale datasets with high dimensionality.

## 4. Model Evaluation

Each model was evaluated using precision, recall, and F1-score.

## 4.1 Classification Reports

### Logistic Regression:

	precision	recall	f1-score	support
Class 0	0.60	0.48	0.54	427481
Class 1	0.35	0.31	0.33	214576
Class 2	0.52	0.69	0.59	347935
accuracy			0.52	989992
macro avg	0.49	0.49	0.49	989992
weighted avg	0.52	0.52	0.51	989992

### Decision Tree:

	precision	recall	f1-score	support
Class 0	0.99	0.99	0.99	427481
Class 1	0.98	0.98	0.98	214576
Class 2	0.99	0.99	0.99	347935
accuracy			0.99	989992
macro avg	0.98	0.98	0.98	989992
weighted avg	0.99	0.99	0.99	989992

### XGBoost:

	precision	recall	f1-score	support
Class 0	0.84	0.96	0.90	427481
Class 1	0.93	0.78	0.85	214576
Class 2	0.94	0.87	0.91	347935
accuracy			0.89	989992
macro avg	0.91	0.87	0.89	989992
weighted avg	0.90	0.89	0.89	989992

## 5. Best Performing Model

- **XGBoost** was chosen as the best-performing model due to its high precision-recall balance and feature importance analysis.
- Key contributing factors to classification were network-based and behavioral attributes.

## 6. Deployment Strategy

- **API Deployment:** Model integrated into a REST API for real-time predictions.
- **Cloud Deployment:** Hosted on AWS Lambda for scalable inference.
- **Dashboard Integration:** Results visualized for SOC analysts.

## 7. Conclusion

The project successfully demonstrated the use of machine learning for cybersecurity incident classification. The XGBoost model provided a significant improvement in classification accuracy, aiding in better incident prioritization. Future work will explore deep learning approaches and real-time threat intelligence integration.

Thank You...

Documentation Report Submitted By Sakthi Krishna Kumar (Aspiring Data Scientist)