

Project Report

Microsoft Cybersecurity Incident Classification with Machine Learning

Introduction:

The rise in cyber threats has significantly increased the workload for Security Operations Centers (SOCs), which are responsible for monitoring, detecting, and responding to cybersecurity incidents. SOCs deal with a vast number of alerts daily, many of which are false positives or benign incidents that do not require immediate action. This creates a challenge in efficiently triaging and prioritizing real threats, leading to alert fatigue and slower response times.

This project aims to address these challenges by developing a machine learning model capable of classifying cybersecurity incidents into three categories: True Positive (TP), False Positive (FP), and Benign Positive (BP). By automating the classification process, the model will assist SOC analysts in identifying true threats faster, improving the overall security posture of organizations while reducing the manual effort involved in managing false alarms.

Problem Statement:

SOCs face the challenge of managing an overwhelming number of alerts daily. Most of these alerts are not critical, but manually identifying which ones require attention is time consuming and prone to errors. This project seeks to solve this problem by creating a machine learning model that automatically classifies incidents as True Positive, False Positive, or Benign Positive. The model will reduce manual efforts, improve response times, and help SOC teams focus on actual threats.

Data Exploration:

- **Data Loading:** The dataset was loaded in chunks to handle its large size.
- **Summary Statistics:** The data was analysed to check its structure, data types, and missing values.
- **Visualizations:** The distribution of key features, including the target variable (Incident Grade), was visualized to understand class imbalances.
- **Class Imbalance:** The Benign Positive class is significantly overrepresented compared to the other classes.

Data Preprocessing:

- **Handling Missing Data:** Missing values were imputed using forward fill and mean imputation. Columns with more than 50% missing values were dropped.
- **Feature Engineering:** Derived timestamp-based features and removed redundant columns.

- **Encoding Categorical Variables:** Categorical features were converted into numerical formats using encoding technique.
- **Scaling:** Standardized numerical features to ensure equal contribution during model training.

Data Splitting:

Split the data into training and validation sets to evaluate model performance.

Train-Validation Split: Data was split into 80% for training and 20% for validation, while maintaining the balance between classes.

Model Selection and Training:

- **Logistic Regression:** A simple model used as a baseline for comparison.
- **Decision Tree:** A non-linear model that works well for small datasets and easy interpretability.
- **Random Forest:** An ensemble of decision trees that provides more accuracy and stability.
- **XGBoost:** A powerful algorithm that handles large datasets efficiently.

✓ **Random Forest** was the best-performing model, achieving high accuracy and macroF1 scores.

✓ **XGBoost** also performed well, though slightly lower than Random Forest.

Model Evaluation and Tuning:

Evaluate the model's performance using cross-validation and optimize it using hyperparameter tuning.

Metrics Used

- **Accuracy:** Measures overall correctness.
- **Precision:** Measures how many positive predictions were correct.
- **Recall:** Measures how well the model identifies actual positives.
- **Macro-F1 Score:** A balanced metric that treats all classes equally.

Hyperparameter Tuning: RandomizedSearchCV was used to find the best settings for Random Forest and XGBoost

Key Outcomes

- **Random Forest** performed best, achieving an accuracy of 99% and a macro-F1 score of 97%.

Classification Report:

	precision	recall	f1-score	support
0	0.94	0.97	0.96	11720
1	0.96	0.97	0.97	13464
2	1.00	0.99	1.00	126942
accuracy			0.99	152126
macro avg	0.97	0.98	0.97	152126
weighted avg	0.99	0.99	0.99	152126

Evaluation on Test Set:

- Test the final model on unseen data to ensure it generalizes well.
- The Random Forest model was evaluated on the test set, achieving high precision, recall, and macro-F1 scores.

Classification Report (Test Data)

The classification report below provides detailed performance metrics for each category on the test dataset

Classification Report on Test Data:

	precision	recall	f1-score	support
0	0.68	0.90	0.77	24124
1	0.56	0.78	0.65	21252
2	1.00	0.94	0.97	303765
accuracy			0.93	349141
macro avg	0.74	0.88	0.80	349141
weighted avg	0.95	0.93	0.93	349141

Macro-F1 Score: 0.80

Macro Precision: 0.74

Macro Recall: 0.88