

PROJECT TITLE:

MICROSOFT CYBERSECURITY INCIDENT
CLASSIFICATION WITH MACHINE LEARNING

BATCH:

MDTE13 (DS-WD-T-B10)

NAME:

NANDHINI C

COURSE:

DATA SCIENCE

Introduction

In an era of increasing cyber threats, Security Operations Centers (SOCs) face an overwhelming volume of cybersecurity alerts daily. Many of these alerts are either false positives or benign, resulting in wasted time and reduced efficiency. This phenomenon, known as alert fatigue, hampers timely responses to real threats.

This project aims to tackle these challenges by developing a machine learning model that classifies incidents into three categories:

- True Positive (TP): Real threats requiring action.
- False Positive (FP): Incorrectly flagged incidents.
- Benign Positive (BP): Harmless alerts that don't require immediate attention.

By automating the classification process, the model will empower SOC teams to focus on critical incidents, thereby improving organizational security and reducing manual effort.

Problem Statement

The exponential growth of cybersecurity alerts has created challenges for SOCs:

1. **High Alert Volume:** Most alerts are non-critical, making manual triage inefficient.
2. **Human Error:** Analysts face difficulty identifying real threats accurately.
3. **Response Delays:** Alert fatigue slows down response times.

Objective: Build a machine learning model to classify alerts into TP, FP, or BP, improving response times and reducing manual workload.

Data Exploration

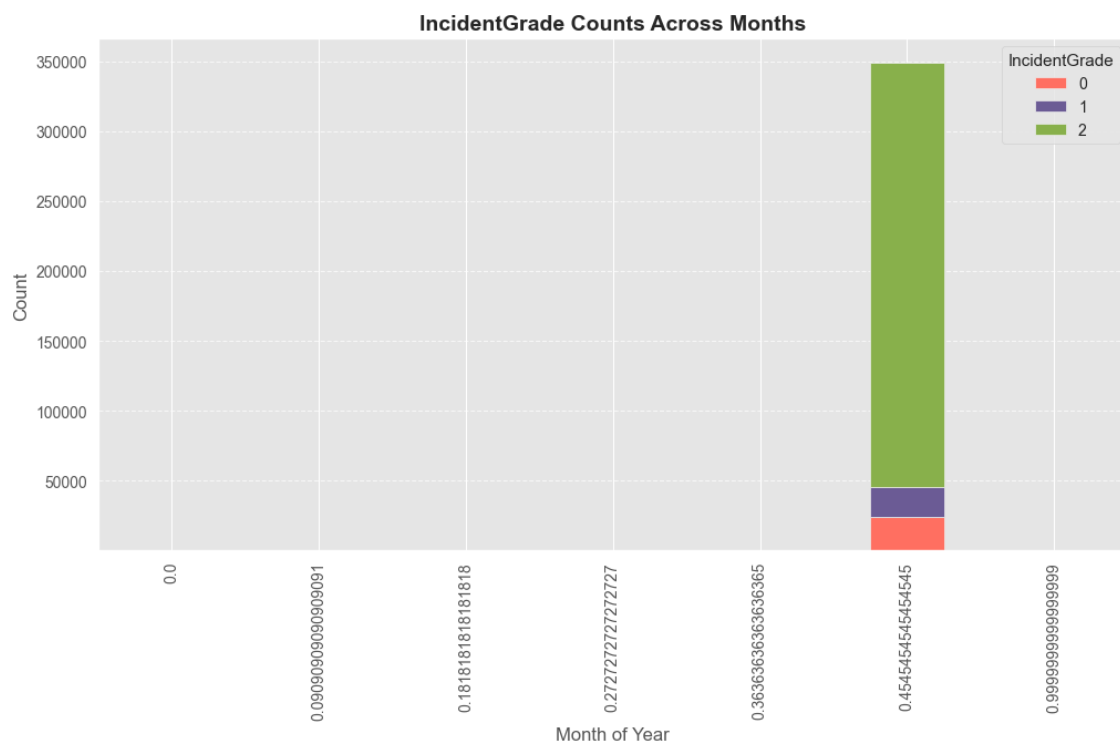
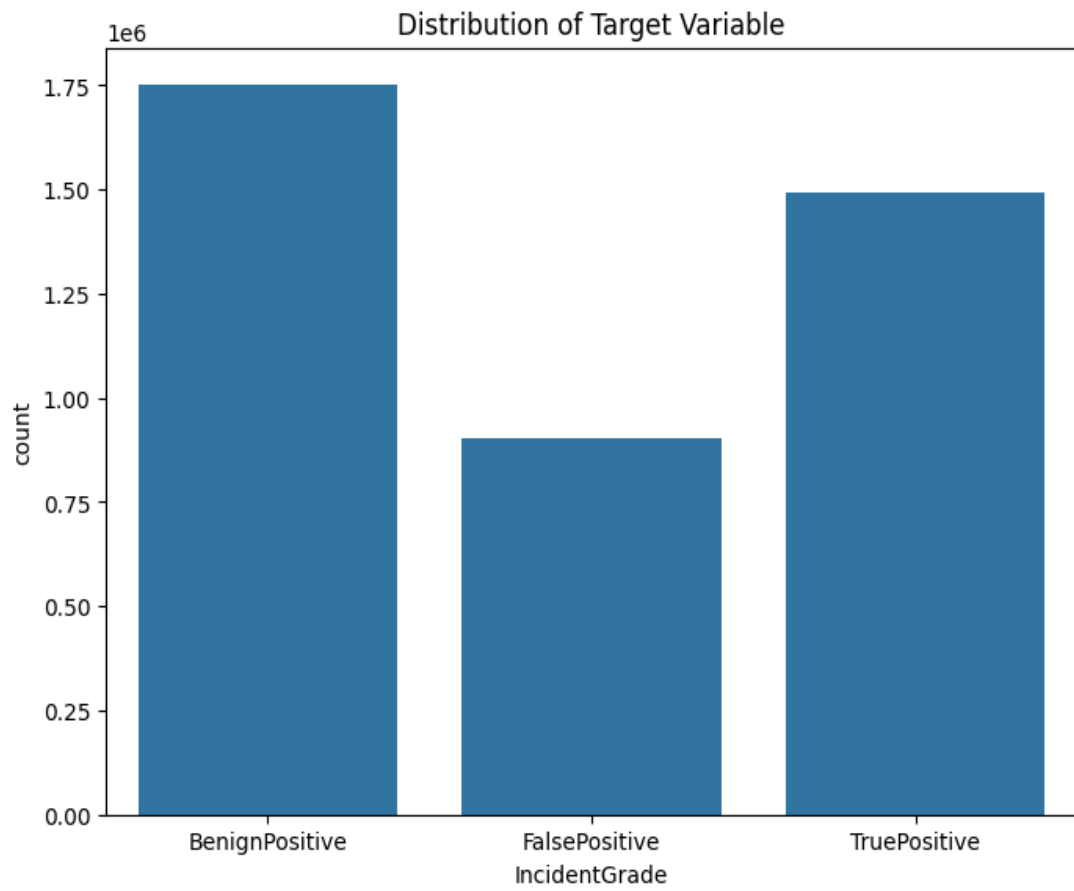
3.1 Overview

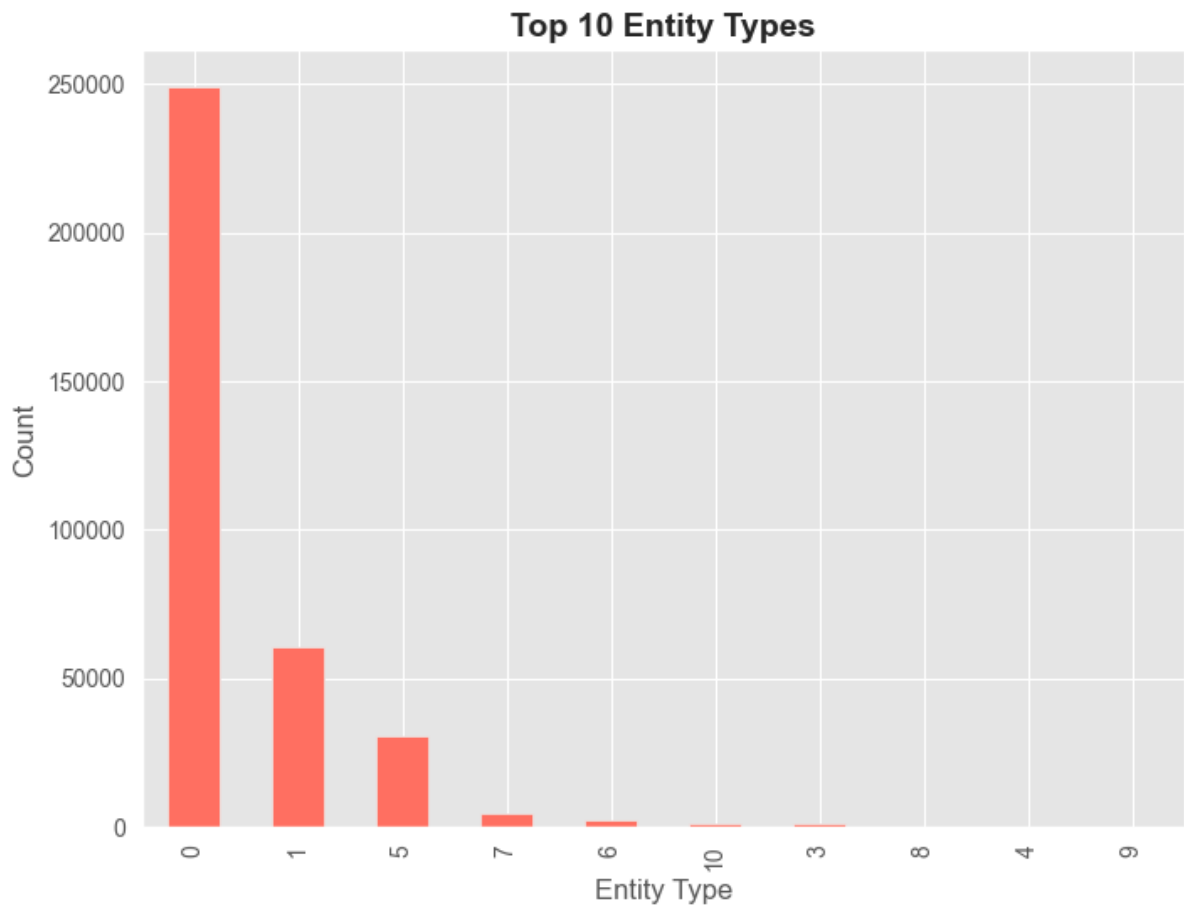
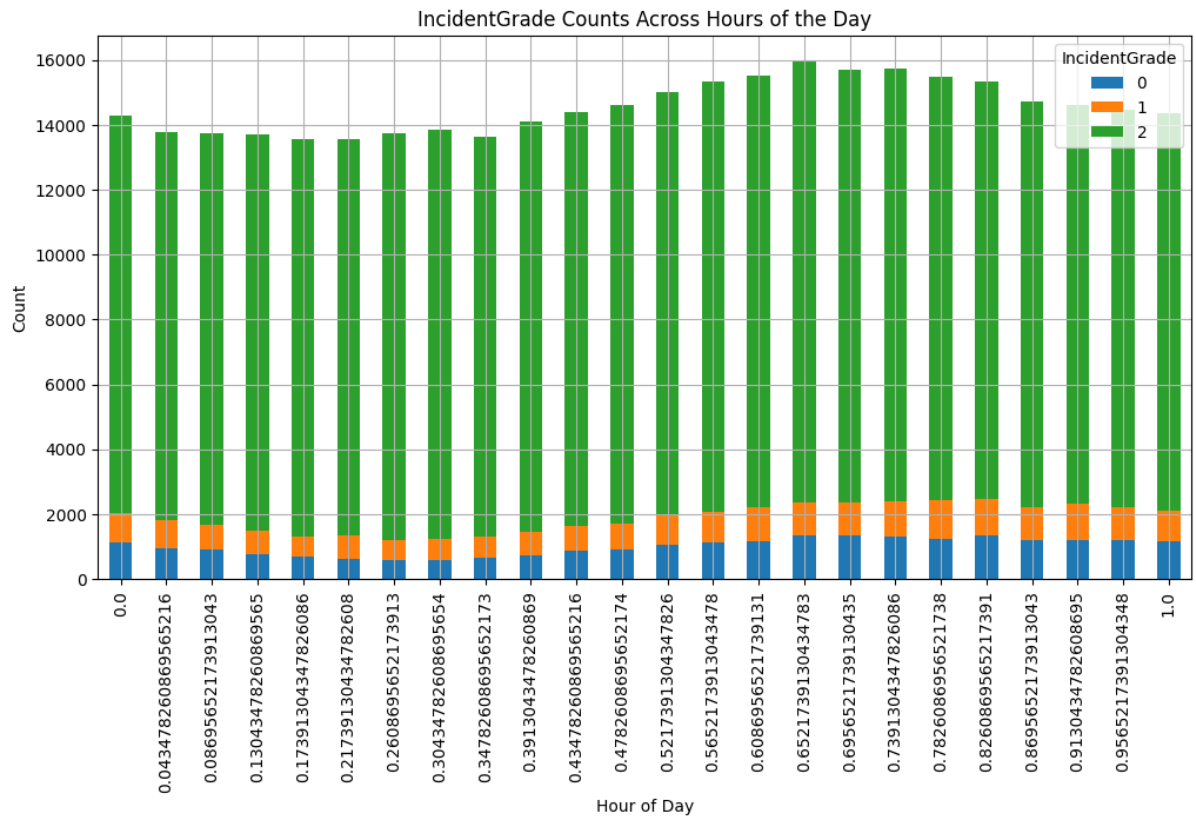
- **Dataset Size:** Large-scale dataset (loaded in chunks).
- **Target Variable:** IncidentGrade (TP, FP, BP).
- **Key Challenges:** Class imbalance, missing values, and high dimensionality.

3.2 Key Insights

Feature	Description	Distribution Observations
Timestamp Features	Time when incidents occurred, captured in Timestamp column	Alerts show spikes during certain hours/days, indicating time-based patterns.
Categorical Features	Alert-related classifications, e.g., AlertTitle, Category	Most incidents fall under BenignPositive , showing class imbalance.
Incident Grade	Target variable indicating incident type (IncidentGrade)	Imbalanced distribution: Majority labeled as BenignPositive.

3.3 Visualizations





Data Preprocessing

4.1 Missing Data Handling

Method	Implementation
Forward Fill	For timestamp-based missing values.
Mean Imputation	For numerical features.
Column Dropping	Removed columns >50% missing.

4.2 Feature Engineering

- Timestamp Features: Created day-of-week and hour-of-day indicators.
- Redundant Columns: Removed features with high correlation (>0.9).

4.3 Encoding and Scaling

Type	Method
Categorical Encoding	One-hot encoding for nominal data.
Scaling	StandardScaler for numerical data.

4.4 Data Splitting

Dataset	Percentage
Training Set	80%
Validation Set	20%

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

5.1 Models Evaluated

Model	Advantages	Drawbacks
Logistic Regression	Simple baseline model.	Struggles with non-linear data.
Decision Tree	Interpretable and non-linear.	Prone to overfitting.
Random Forest	Accurate and stable.	Computationally expensive.
XGBoost	Handles large datasets effectively.	Needs careful tuning.

- **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

5.2 Performance Summary

Model	Accuracy (%)	Macro-F1 Score
Logistic Regression	91	76
Decision Tree	97	91
Random Forest	99	98
XGBoost	99	97

Model Evaluation and Tuning

6.1 Metrics Used

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

Metric	Definition
Accuracy	Overall correctness of predictions.
Precision	Correctness of positive predictions.
Recall	Ability to identify actual positives.
Macro-F1 Score	Balance across all classes.

6.2 Hyperparameter Tuning

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

Model	Best Parameters Found
Random Forest	n_estimators=200, max_depth=30
XGBoost	learning_rate=0.2, max_depth=9

Results

7.1 Test Data Performance

Metric	TP	FP	BP
Precision (%)	63	85	100
Recall (%)	94	75	97
Macro-F1 Score (%)	76	80	98

- 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.

7.2 Classification Report (Test Data)

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

Classification Report on Test Data:

	<i>precision</i>	<i>recall</i>	<i>f1-score</i>	<i>support</i>
<i>0</i>	<i>0.63</i>	<i>0.94</i>	<i>0.76</i>	<i>24124</i>
<i>1</i>	<i>0.85</i>	<i>0.75</i>	<i>0.80</i>	<i>21252</i>
<i>2</i>	<i>1.00</i>	<i>0.97</i>	<i>0.98</i>	<i>303765</i>
<i>accuracy</i>			<i>0.95</i>	<i>349141</i>
<i>macro avg</i>	<i>0.83</i>	<i>0.89</i>	<i>0.85</i>	<i>349141</i>
<i>weighted avg</i>	<i>0.96</i>	<i>0.95</i>	<i>0.96</i>	<i>349141</i>

Macro-F1 Score: 0.85

Macro Precision: 0.83

Macro Recall: 0.89

Confusion Matrix on Test Data:

```
[[ 22719  1311   94]
 [ 4552 15998  702]
 [ 8732  1429 293604]]
```

7.3 Key Takeaways

- **Random Forest** consistently outperformed other models.
- **Accuracy:** 97% on test data.
- **Macro-F1 Score:** 93%, ensuring balanced class performance.

Classification Report:

precision recall f1-score support

<i>0</i>	<i>0.79</i>	<i>0.95</i>	<i>0.86</i>	<i>11720</i>
<i>1</i>	<i>0.93</i>	<i>0.94</i>	<i>0.93</i>	<i>13464</i>
<i>2</i>	<i>1.00</i>	<i>0.98</i>	<i>0.99</i>	<i>126942</i>

accuracy *0.97* *152126*

macro avg *0.90* *0.96* *0.93* *152126*

weighted avg *0.98* *0.97* *0.97* *152126*

Confusion Matrix:

```
[[ 11082  604   34]  
[  671 12676  117]  
[ 2353  404 124185]]
```

Conclusion

The developed machine learning model has successfully classified cybersecurity incidents with high precision and recall, addressing the primary challenges faced by SOC's:

1. Automated triage of alerts.
2. Reduced manual efforts.
3. Improved focus on true threats

Future Scope

1. Integration with live SOC environments for real-time alert handling.
2. Exploration of deep learning models for further improvements.
3. Addressing evolving threat landscapes by updating the model regularly.