

Intrusion Detection System (IDS) with NSL-KDD Dataset & XGBoost

1. Introduction

This project implements an Intrusion Detection System (IDS) using the NSL-KDD dataset. The IDS distinguishes between normal and malicious network traffic using machine learning, focusing on the XGBoost algorithm. Objectives include building a robust pipeline, evaluating model performance, and interpreting results with SHAP.

2. Dataset

The NSL-KDD dataset is an improved version of the KDD Cup 1999 dataset. It contains features representing network connections (categorical and numerical) and labels indicating normal or attack traffic. We preprocess the data by encoding categorical features, scaling numerical values, and handling class imbalance via SMOTE.

3. Methodology

- Data preprocessing: Encoding, scaling, handling imbalance. - Baseline model: Random Forest for initial benchmarking. - Main model: XGBoost, tuned for hyperparameters. - Evaluation: Precision, Recall, F1-score, Accuracy, False Positive Rate, Confusion Matrix. - Interpretability: SHAP analysis to identify influential features. - Ethical considerations: False positives, data privacy, bias in detection.

4. Results

The XGBoost model outperformed the baseline Random Forest in precision and recall, achieving a strong F1-score and robust accuracy. SHAP analysis highlighted critical features in distinguishing attacks from normal traffic.

5. Conclusion & Future Work

The IDS demonstrates strong performance with XGBoost, showing its effectiveness in cybersecurity applications. Future improvements may include real-time traffic analysis, deep learning approaches, or deployment in streaming systems.