# Network Intrusion Detection using UNSW-NB15 Dataset

By-

Avantika Deshmukh
Daniel Cersosimo
Harshini Akuleti
Himachand Pasupuleti
Purna Jadhav

# Introduction

Importance of network security in today's digital landscape:

❖ In today's interconnected digital landscape, evolving cyber threats pose significant risks to organizations and individuals, making the security of computer networks paramount.

❖ Network security breaches can result in severe consequences, including data breaches, financial losses, and reputational damage for businesses.
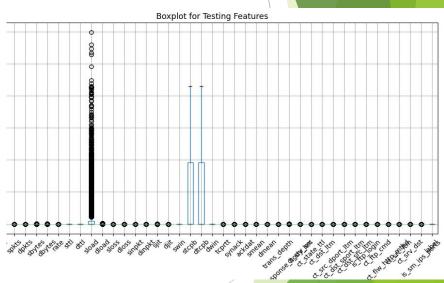
Objectives:

❖ Conduct comprehensive data exploration, preprocessing on the UNSW-NB15 dataset to utilize and assess multiple models for network intrusion detection.

# Data Acquisition and Preprocessing

❖ Data Acquisition:  Secure the UNSW-NB15 dataset

❖ Preprocessing:
  ➢ Identify nulls (none in this dataset)
  ➢ Assess statistical distributions and properties for outliers and noise
  ➢ Handle discrepancies in the categorical features
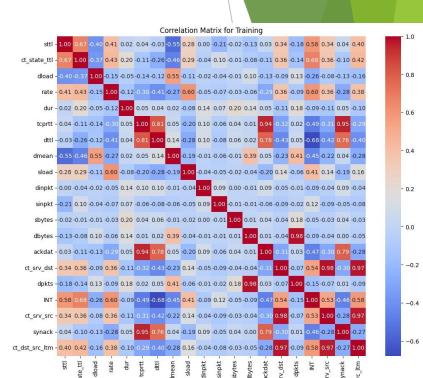  ➢ Transforming categorical data to numerical data using one-hot encoding

# Data Acquisition and Preprocessing

# Exploratory Data Analysis

❖ Performed Correlation Analysis
  ➢ variables with very high positive correlation
    ■ ct_state_ttl and sttl
    ■ ackdat and smean

  ➢ variables with very high negative correlation
    ■ dmean and sttl



Correlation Matrix for Training

# Overview of Machine Learning Models

❖ **Machine Learning for Network Intrusion Detection**
- ➤ Comprehensive exploration of various algorithms to identify effective classifiers.
- ➤ Systematic approach considering each algorithm's strengths and limitations.
- ➤ Aim for optimal performance using appropriate tools and techniques.
- ➤ Conduct a run of each model for both high dimensional and PCA reduced feature inputs. Top accuracies from overall runs details below
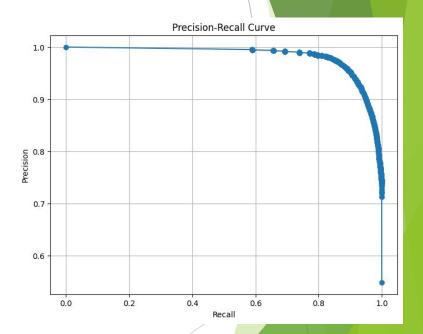
❖ **Machine Learning Algorithms:**
- ➤ k-Nearest Neighbors (k-NN) | top accuracy: .79 | Poor
- ➤ Logistic Regression | top accuracy: .71 | Poor
- ➤ SVM | top accuracy: .77 | Poor
- ➤ Random Forest | top accuracy : .92 | Strong
- ➤ XGBoost | top accuracy: .93 | Strong
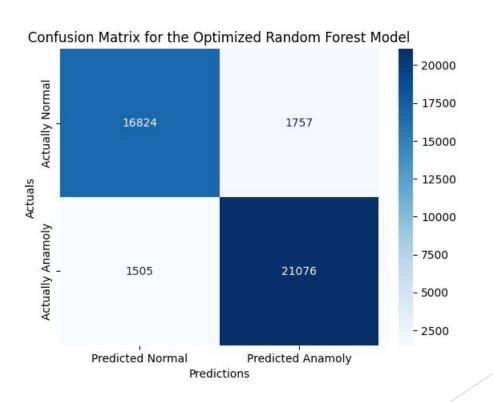- ➤ LightGBM | top accuracy: .94 | Strongest

# Random Forest Optimization

► Three primary methods used to optimize the

Random Forest

  ► Feature Selection
  ► n trees hyperparameter tuning
  ► Anomaly classifying threshold optimization
► Results

|   | precision | recall | f1-score | support |
|---|-----------|--------|----------|---------|
| 0 | 0.92 | 0.91 | 0.91 | 18581 |
| 1 | 0.92 | 0.93 | 0.93 | 22581 |
| ► accuracy | | | 0.92 | 41162 |



Precision-Recall Curve

# Random Forest Visualization



Confusion Matrix for the Optimized Random Forest Model
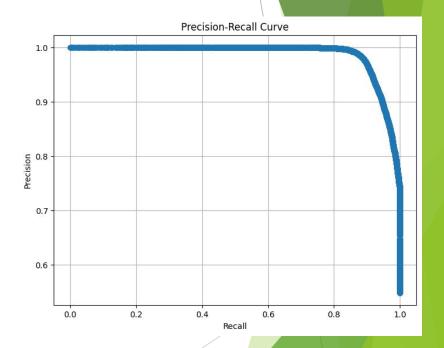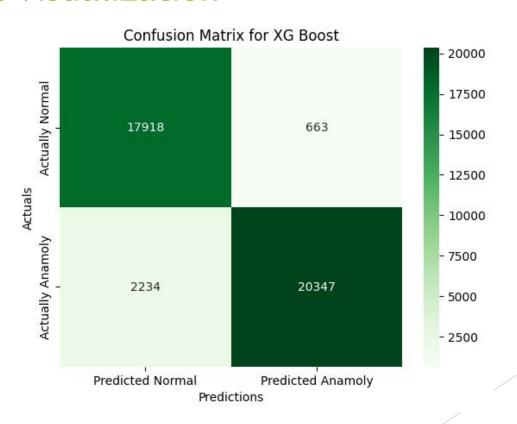
# XG Boost Optimization

- ► Optimized the model via similar anomaly classifying threshold optimization
- ► Results

|   | precision | recall | f1-score | support |
|---|-----------|--------|----------|---------|
| 0 | 0.89 | 0.96 | 0.93 | 18581 |
| 1 | 0.97 | 0.90 | 0.93 | 22581 |
| accuracy | | | 0.93 | 41162 |



Precision-Recall Curve

# XGBoost Visualization



Confusion Matrix for XG Boost
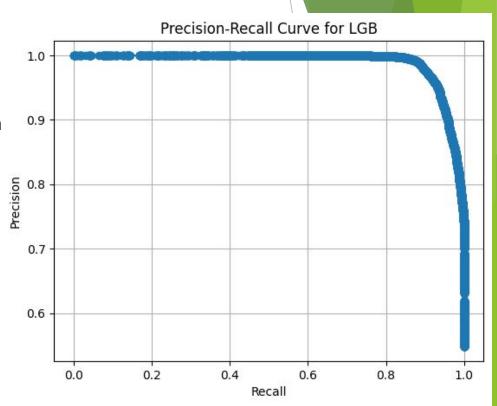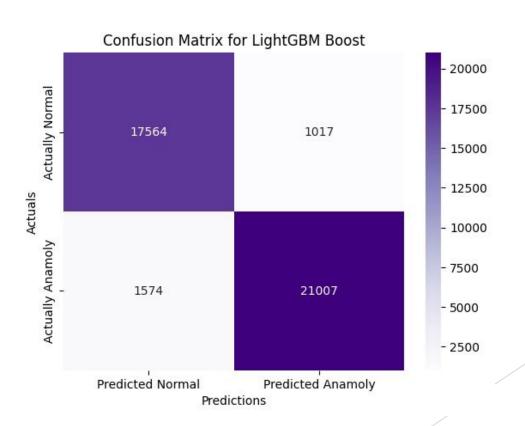
# LightGBM Optimization

► Optimized the model via similar anomaly classifying threshold optimization

► Results

|   | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.92 | 0.95 | 0.93 | 18581 |
| 1 | 0.95 | 0.93 | 0.94 | 22581 |
| accuracy | | | 0.94 | 41162 |

► This is our top performing model



Precision-Recall Curve for LGB
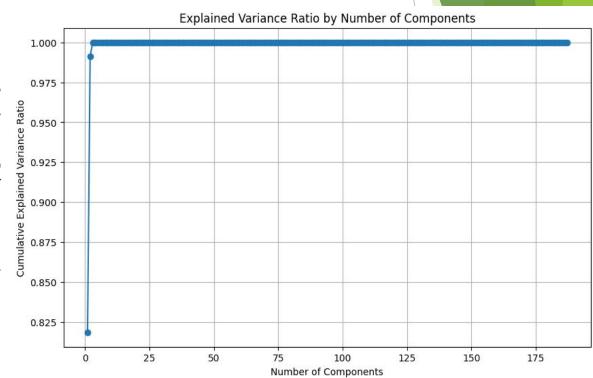
# LightGBM Visualization



Confusion Matrix for LightGBM Boost

# Principal Component Analysis

❖ Performance Comparison: Conducted a side-by-side analysis of models trained with original vs. PCA-transformed features.

❖ Sought to consider PCA to potentially enhance viability of non decision tree based models such as k-NN and even be able to attempt SVM

❖ Utilized the elbow method to determine the optimal amount of principal components as 3

❖ Gridsearch employed for k-NN however improvement was minimal and no PCA input models saw strong performance

Explained Variance Ratio by Number of Components

# Conclusion

- ❖ Identified three robust models with efficient accuracy, recall, precision, and F1-score metrics, showcasing their effectiveness in anomaly detection.

- ❖ LightGBM model outperformed others marginally followed by XGBoost and Random Forest. These models could be recommended for Intrusion Detection System to support Security Information and Event Management (SIEM) systems especially for systems which operate in real time due to the combination of lightweight overhead and performance.