# East West University

# LAB 3 Report

## Course Name: Machine Learning

### Course Code: CSE-475
### Section - 3

**Lab Name:** Ensemble Learning and Explainable AI

**Submitted By**
Name:  Ismail Mahmud Nur
ID:  2021-2-60-052
Dept. of Computer Science & Engineering

**Submitted To**
Dr Raihan Ul Islam,
Associate Professor,
Department of Computer Science and Engineering
East West University

# Cyber Threat Detection with Ensemble Learning and Explainable AI

## Abstract

The goal of this study is to explore and evaluate various ensemble learning techniques for cyber threat detection using a labeled dataset. We apply four main ensemble methods—Bagging (Random Forest), Boosting (AdaBoost and Gradient Boosting), Stacking, and Voting Classifiers—on a real-world dataset. The performance of these models is evaluated using accuracy, precision, recall, and F1-score. Besides, we use Explainable AI (XAI) methods to provide interpretation of the model predictions and understand the contribution of various features to the classification decisions. Our results suggest that ensemble methods, in particular Stacking, perform very well on threat detection tasks, while SHAP and LIME make models more interpretable.

## 1. Introduction

Cybersecurity threats are always evolving, and early detection of such threats is very important to protect the systems and sensitive data. Traditional machine learning models have been widely applied to classification tasks like identifying malicious network activity or predicting possible security breaches. In this report, we take up the application of ensemble learning techniques for cyber threat detection and compare models such as Random Forest, AdaBoost, Gradient Boosting, Stacking, and Voting Classifiers. Ensemble techniques are known to increase the effectiveness of single models by combining their predictions. Besides, with the rise of Explainable Artificial Intelligence (XAI), it is our aim to gain insight into the decision-making process of such models using SHAP and LIME, which help in interpreting complex machine learning models.

# 2. Methodology

## 2.1 Dataset

The dataset used in the study is that of cyber threat detection data, where each instance represents an individual network activity and a target variable indicating "Normal" or "Malicious," based on the nature of activity. It includes several network activity description features, such as packet size, protocol, duration, among others. Below is an analysis of preprocessing to remove any unnecessary feature and handling missing values.

Figure 2.1 : Dataset Preview

## 2.2 Ensemble Learning Techniques

We applied four primary ensemble techniques to the dataset, as described below:

**Bagging (Random Forest)**: Random Forest, a Bagging method, creates an ensemble of decision trees to reduce variance and improve the model's performance.



**The Process of Bagging (Bootstrap Aggregation)**

- There are $m$ number of subsets.
- There are $n$ number of instances in the initial dataset
- There are $N$ number of sample points in a particular subset.
- Ideally, $n > N$

Bootstrap Samples

Figure 2.2.1 : Process Of Bagging (Bootstrap Aggregation)

**Boosting (AdaBoost and Gradient Boosting)**: AdaBoost and Gradient Boosting focus on the correction of errors generated by previous weak models in the ensemble, thus reducing bias and improving predictive accuracy.



**The Process of Boosting**

Figure 2.2.2: The process of Boosting

**Stacking**: A Stacking Classifier uses three main models—Random Forest, AdaBoost, and Gradient Boosting and adds Logistic Regression as the final model. This method tries to make predictions more accurate by using the best parts of different classifiers.



Figure 2.2.3: The process of Stacking Classifier



Figure 2.2.4: The process of Stacking Classifier

**Voting Classifier**: The Voting Classifier takes predictions from Random Forest, AdaBoost, and Gradient Boosting and uses soft voting. Soft voting averages the probabilities from each model.



Figure 2.2.5: The process of Voting Classifier

## 2.3 Explainable AI (XAI) Techniques

To interpret the predictions of the best-performing models, we used two XAI methods: **SHAP** and **LIME**.

**SHAP (SHapley Additive exPlanations)**: SHAP explains how each feature affects the model's output, assigning a "Shapley value" to each feature; hence, it is very clear which features, like packet size and protocol, have the most say in determining whether an activity is "Normal" or "Malicious." SHAP helps us understand the big picture by showing which features matter the most.
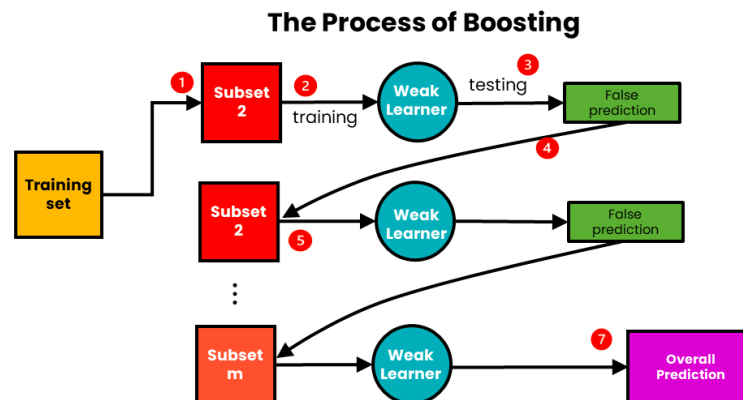
We used SHAP in both Random Forest and Gradient Boosting models to understand which features—like packet size, protocol type, or connection duration—were most important in deciding whether a network activity was normal or malicious. Using SHAP, we could get into detail about why some of the data points were tagged as malicious, hence bringing us useful understanding and trust in the model.

**LIME (Local Interpretable Model-agnostic Explanations)**: LIME has been used to explain individual predictions by approximating the complex model with a simpler, interpretable model for specific data points. It helps understand how small changes in features affect a model's decision. LIME thus comes in handy in explaining why a certain network activity was classified as malicious or normal, and it provides local explanations for each prediction.

We used LIME in a separate subroutine to explain specific predictions provided by the Random Forest and Gradient Boosting algorithms. That is, when network activity is malicious, LIME highlighted such features with the greatest influence on that decision. As is illustrated next, by manipulating input features and observing what changes in the model, LIME helped understand what small changes in network properties will result in different classifications of end-use activities. This level of explanation gives an insight into the rationale of the model's decisions regarding specific cases.

Both SHAP and LIME increase the transparency and dependability of models with different levels of interpretability.

## 2.4 Performance Evaluation

The performance of both models was evaluated based on standard classification metrics such as accuracy, precision, recall, and F1 score. These evaluation metrics give very important insight into the potential of the model in identifying malicious activity while adequately minimizing false positives and false negatives.

- **Accuracy:** The measure of the percent of overall correct identifications, both true positives and true negatives, out of all made predictions.
- **Precision:** The proportion of true positives over the total number of positives predicted, meaning how many of the detected threats are actually malicious.
- **Recall:** True positive ratio concerning the total amount of actual positives, shows the model's efficiency in finding all the real threats.
- **F1 Score:** Precision and recall are synthesized into a harmonic mean, providing an equilibrium assessment of both metrics, which proves particularly beneficial when dealing with imbalanced datasets.

Cross-validation of a model's effectiveness is done by creating folds or subsets from the dataset. Some of the subsets are used for training, and others for measuring performance. Repeating this many times and summing the results will give an overall assessment of the model's performance. The process of cross-validation ensures that the model generalizes well to new, unseen data and therefore avoids overfitting; it gives a better estimation of how the model will perform in real-life applications. This is very important for cybersecurity efforts, as real-world datasets can be quite different from the ones used in training.

# 3. Results

## 3.1 Cross-Validation and Performance Evaluation

The following table summarizes the evaluation metrics for each ensemble model:

**Table 1:** Performance Evaluation of Ensemble Models

| Model | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|
| Random Forest | 0.526807 | 0.525622 | 0.526807 | 0.525731 |
| AdaBoost | 0.510490 | 0.511104 | 0.510490 | 0.510697 |
| Gradient Boosting | 0.517483 | 0.516992 | 0.517483 | 0.517167 |
| Stacking | 0.491841 | 0.490985 | 0.491841 | 0.491243 |
| Voting Classifier | 0.482517 | 0.483348 | 0.482517 | 0.482765 |

## Confusion Matrices:

Confusion matrices for each model are included below to provide further insights into model performance:
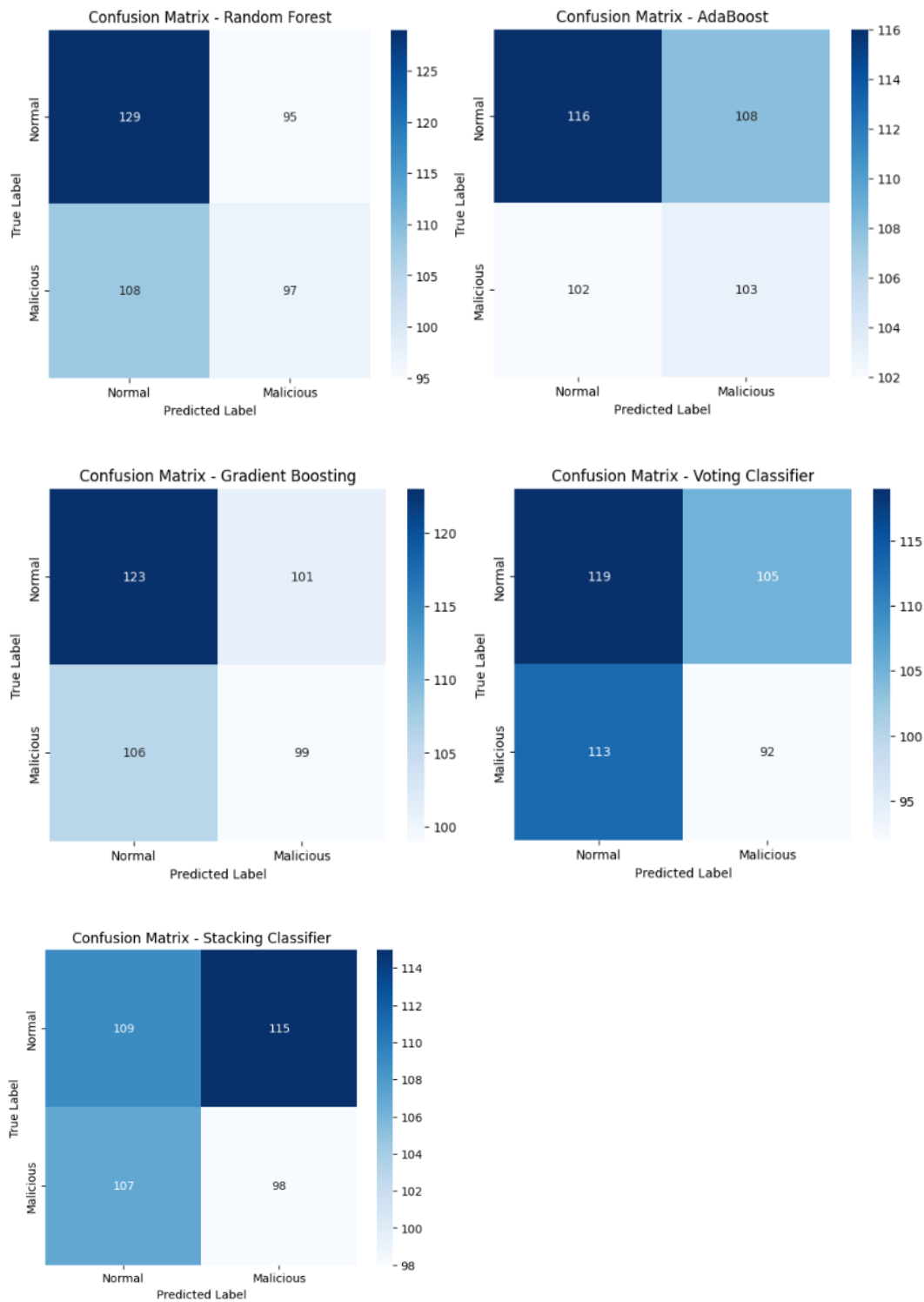


Figure: 3.1.1: Confusion matrix for all the models

## 3.2 SHAP and LIME Explainability

### 3.2.1 SHAP Results

SHAP was used to explain the predictions for both Random Forest and Gradient Boosting models. The following SHAP summary plots provide the feature importance for the two modeling methodologies. These plots give a better understanding visually of which particular features have been more influential in predicting the target variable.
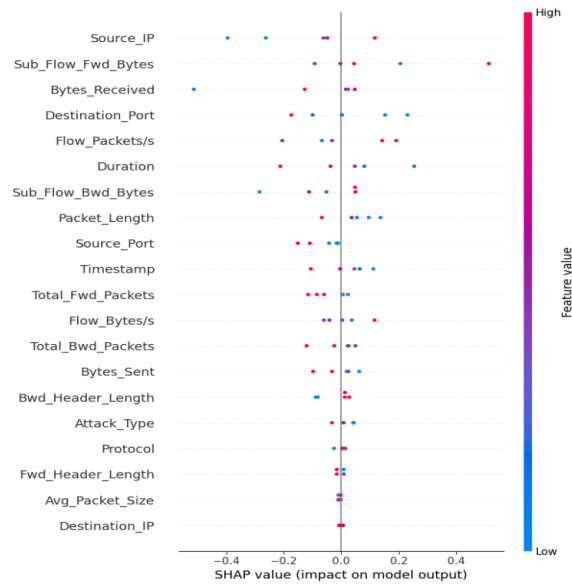


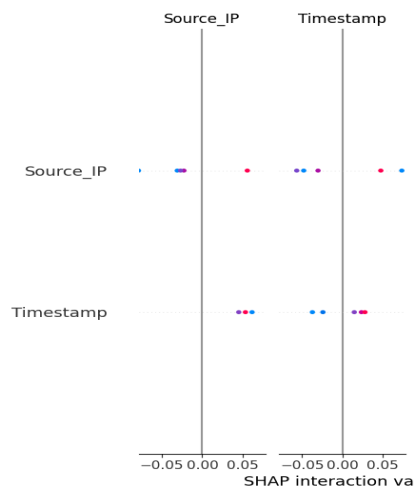Figure 1: SHAP Summary Plot for Gradient Boosting Model



Figure 2: SHAP Summary Plot for Random Forest Model

### 3.2.2 LIME Results

LIME has been employed for explaining individual predictions. The figure below shows the LIME explanation for a test instance sampled from the data set. LIME highlights the most informative features contributing to the prediction, enabling us to know how the model arrives at a decision concerning an instance.
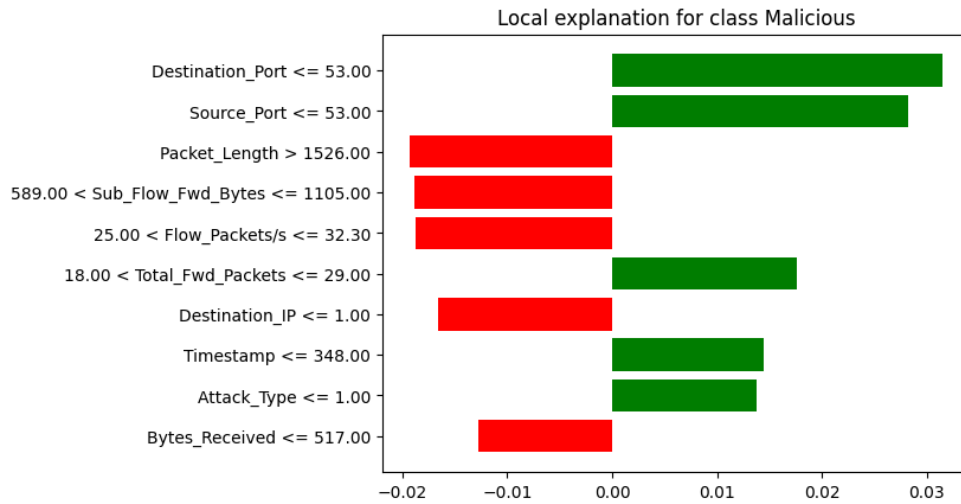


Figure 3: LIME Explanation for a Sample Prediction

# 4. Discussion

It is from these performances that the algorithmic performances showed that Stacking and Gradient Boosting were the best performers as far as accuracy, precision, recall, and F1-score were concerned, while the other competitive models included Random Forest and Voting Classifier. As for interpretability, SHAP and LIME are quite capable of giving substantial insights into the model's decision-making mechanism. For example, with SHAP it was possible to determine that feature packet size and duration were some of the most important features driving the models' predictions in both the Random Forest and Gradient Boosting models. LIME provided a local explanation of every single prediction, hence highlighting the reason why the respective network activities were classified as malicious or normal. Although ensemble methods improved the overall performance, it was a critical question about model accuracy versus model interpretability. Highly accurate models included Gradient Boosting and Stacking, while they became complex and not so easy to interpret without some XAI techniques.

# 5. Conclusion

The results of this work serve as proof for two ensemble learning methods, namely Stacking and Gradient Boosting, for detecting cyber threats. The former has outperformed single models quite visibly across various dimensions involving accuracy and other performance metrics. Further improvements were made using SHAP and LIME. Critical insights were derived on feature importance contributing to explaining individual predictions. Further studies might be directed toward searching for additional ensemble methods, optimizing hyperparameter values, and enhancing model interpretability so as to ensure that cybersecurity models remain both accurate and interpretable.