



IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 12 Issue: XI Month of publication: November 2024 DOI: https://doi.org/10.22214/ijraset.2024.65559

www.ijraset.com

Call: 🕥 08813907089 🔰 E-mail ID: ijraset@gmail.com



Detection of DDoS Attacks Using XGBoost-Based Feature Selection and Deep Learning

Arjun R Nair

I. INTRODUCTION

In the rapidly evolving landscape of cybersecurity, network intrusion detection systems (IDS) have become essential tools for safeguarding information assets against malicious activities. As cyber threats grow in complexity and frequency, particularly Distributed Denial of Service (DDoS) attacks, there is an increasing need for advanced detection methodologies that can accurately differentiate be- tween benign and malicious network traffic. Machine learning and deep learning techniques have emerged as powerful approaches to improve the efficacy of IDS by learning intricate patterns from network data.

The Canadian Institute for Cybersecurity's Intrusion Detection System 2017 (CICIDS2017) dataset has been established as a benchmark for evaluating intrusion detection algorithms. This dataset contains a rich mix of benign traffic and the most up-to-date common attacks, which closely resemble real-world network environments through Packet Capture (PCAP) files. It includes the results of network traffic analysis using CICFlowMeter, offering labeled flow data based on timestamps, source and destination IP addresses, ports, protocols, and attack types in the CSV format. A significant aspect of the CICIDS2017 dataset was the generation of realistic background traffic. Utilizing the B-Profile system proposed by Sharafaldin et al. (2016), the dataset profiles the abstract behavior of human interactions to generate naturalistic benign traffic. This approach models the behavior of 25 users across various protocols including HTTP, HTTPS, FTP, SSH, and email, thereby creating a comprehensive and realistic dataset for intrusion detection research.

Despite the availability of such detailed datasets, accurately detecting anoma- lies and outliers, especially DDoS attacks—remains a challenge. Factors such as high dimensionality, redundant features, and class imbalance in the data can hinder the performance of traditional detection methods. Therefore, there is a pressing need to develop advanced models that not only improve detection accuracy but also optimize computational efficiency and resource utilization.

A. Problem Statement

This research addresses the challenge of detecting anomalies and DDoS attacks within network traffic by improving intrusion detection accuracy using machine learning algorithms. Specifically, it focuses on enhancing the performance of IDS by leveraging deep learning techniques to classify network flow samples into benign or malicious categories, considering both binary and multiclass classifi- cation scenarios.

B. Objectives

The primary objective of this study is to develop a robust intrusion detection model that achieves high accuracy while maintaining computational efficiency. The specific goals include:

- *1)* Feature Selection: Implement feature selection using the Extreme Gra- dient Boosting (XGBoost) algorithm to identify the most significant fea- tures, reduce data dimensionality, and optimize memory usage without compromising detection performance.
- 2) Handling Class Imbalance: Address the issue of class imbalance in- herent in the CICIDS2017 dataset to ensure that minority classes, such as specific types of attacks, are accurately detected. Techniques such as resampling or cost-sensitive learning may be employed to mitigate this challenge.
- 3) Deep Learning Model Development: Leverage Multilayer Percep- tron (MLP) architectures through the Keras framework to construct deep learning models capable of learning complex patterns in network traffic data. The models will be configured for both binary classification (distin- guishing between benign and malicious traffic) and multiclass classification (identifying specific attack types).
- 4) Performance Optimization: Achieve state-of-the-art results with an accuracy of approximately 99% or higher. This involves fine-tuning the deep learning models and optimizing hyperparameters to enhance detec- tion capabilities while minimizing computational overhead.



5) Validation and Evaluation: Rigorously evaluate the proposed models using appropriate metrics such as accuracy, precision, recall, F1-score, and area under the Receiver Operating Characteristic (ROC) curve. Cross- validation and testing on unseen data will be conducted to ensure the generalizability of the results.

II. LITERATURE REVIEW

Intrusion Detection Systems (IDS) are critical components in modern network security infrastructure, designed to monitor network traffic and identify poten-tial threats or unauthorized access [2]. Over the years, IDS have evolved from signature-based detection methods to incorporate advanced machine learning (ML) and deep learning techniques, enabling them to detect novel and sophisticated cyber-attacks [3].

A. Intrusion Detection Systems and Machine Learning

Traditional IDS rely on predefined signatures or anomaly detection rules, which can be insufficient against zero-day attacks or sophisticated evasion techniques [4]. Machine learning offers a dynamic and adaptive approach, where models learn patterns from data to detect anomalies [4]. Studies have demonstrated that ML-based IDS can significantly improve detection rates and reduce false positives [1].

B. Feature Engineering in IDS

Feature engineering is a pivotal process in developing ML models for IDS. It involves selecting and transforming variables to improve model performance [5]. In network intrusion detection, features can include various network flow char- acteristics such as packet sizes, durations, protocols, and statistical measures [1]. Effective feature selection reduces dimensionality, mitigates the curse of dimensionality, and enhances model interpretability [6]. Techniques like prin- cipal component analysis (PCA) and autoencoders have been used to extract significant features [7].

Sharafaldin et al. [1] emphasized the importance of comprehensive feature selection and proposed a new approach for generating reliable datasets, leading to the creation of the CICIDS2017 dataset.

C. CICIDS2017 Dataset in IDS Research

The CICIDS2017 dataset has become a standard benchmark for evaluating IDS models due to its realistic representation of network traffic and inclusion of contemporary attack types [1]. Several studies have utilized this dataset to assess the effectiveness of various detection techniques.

For instance, Maciá-Fernández et al. [8] investigated the impact of metadata features on machine-learned IDS models using the CICIDS2017 dataset. They found that including certain metadata can contaminate the model, leading to overfitting and reduced generalizability. Their study highlights the necessity of careful feature selection to avoid the inclusion of non-representative data that could skew model performance.

Other researchers have employed the dataset to explore deep learning ap- proaches. Yin et al. [9] proposed a deep learning framework using recurrent neural networks (RNN) to achieve high detection rates. Similarly, Ullah and Mahmoud [10] developed a hybrid model combining convolutional neural net- works (CNN) and long short-term memory (LSTM) networks, demonstrating improved performance on the CICIDS2017 dataset.

D. Machine Learning Approaches for IDS

Among machine learning algorithms, ensemble methods like Extreme Gradient Boosting (XGBoost) have gained significant attention for IDS applications [11]. XGBoost is known for its scalability, efficiency, and high predictive accuracy, especially on large and complex datasets.

XGBoost operates by building an ensemble of weak learners, typically de- cision trees, in a sequential manner where each new tree focuses on correcting the errors of the previous ones [11]. This approach makes it powerful for classi- fication tasks within IDS, particularly in anomaly detection where patterns are subtle and complex.

A key advantage of XGBoost is its ability to handle imbalanced datasets, a common issue in IDS where attack instances are rare compared to normal traffic [12]. Techniques such as weighting the loss function and using appropriate evaluation metrics help in addressing class imbalance. Furthermore, XGBoost's parallelization capability ensures it trains faster compared to other gradient boosting methods.



Comparisons with other algorithms have been explored in the literature. LightGBM, developed by Microsoft, offers faster training and lower memory us- age by using histogram-based algorithms and leaf-wise tree growth [13]. It often achieves higher accuracy than level-wise growth used in many other algorithms, particularly when dealing with large datasets or high-dimensional data [14].

CatBoost, another gradient boosting algorithm, is designed to handle cate- gorical features efficiently and reduce overfitting [15]. It incorporates techniques to address the prediction shift caused by target leakage, which is beneficial in IDS applications where data integrity is crucial.

Random Forest, an ensemble method using bagging and feature randomness, constructs multiple decision trees and aggregates their results, offering robust- ness and ease of interpretation [16]. Although it may not achieve the same level of accuracy as boosted models, Random Forest is valuable for its robustness to overfitting and its performance on a wide range of datasets [17].

In the context of IDS, the choice between these algorithms depends on spe- cific dataset characteristics and computational constraints. For instance, Light- GBM might be preferred for very large datasets due to its speed and efficiency [14], while Random Forest may be suitable for problems where model inter- pretability is a priority.

E. Feature Interpretation and Explainability

Understanding model decisions is crucial in IDS to ensure trust and compliance. Techniques such as Shapley values provide insights into feature contributions to predictions [18]. Shapley values offer a unified approach to interpreting model outputs, helping analysts understand the importance of each feature in the detection process [19].

Moreover, tools like Pandas Profiling (now known as ydata-profiling) facili- tate exploratory data analysis by generating comprehensive reports on dataset features [20]. Such tools help in identifying data quality issues, understanding feature distributions, and uncovering patterns that may influence model perfor- mance.

Quantile functions are also used in statistical analysis to understand the distribution of features, which is valuable in preprocessing and normalizing data for ML models [21]. By analyzing quantiles, researchers can detect outliers and better prepare data for training robust IDS models.

F. Summary

The integration of advanced machine learning techniques in IDS has shown sig- nificant promise in enhancing detection capabilities. The CICIDS2017 dataset serves as a robust benchmark for evaluating these methods. XGBoost, among other algorithms, has demonstrated high accuracy and efficiency in handling the complexities of intrusion detection. However, the choice of algorithm should be tailored to the specific characteristics of the dataset and the requirements of the deployment environment. Feature engineering and model interpretability remain critical components in developing effective IDS solutions.

III. METHODOLOGY

A. Dataset Description

The Canadian Institute for Cybersecurity's Intrusion Detection System 2017 (**CICIDS2017**) dataset is utilized in this research as the primary source of network traffic data for intrusion detection analysis. This dataset is renowned for its comprehensive and realistic representation of modern network traffic pat- terns, including both benign activities and a variety of malicious attacks [1].

The CICIDS2017 dataset was created with the goal of resembling true real- world data, incorporating Packet Capture (PCAP) files that capture network traffic over a period of five consecutive days. The data collection commenced at 9 a.m. on Monday, July 3, 2017, and concluded at 5 p.m. on Friday, July 7, 2017, covering both normal and attack traffic periods.

Key features of the dataset include:

- 1) Benign Traffic: The dataset includes naturalistic benign background traffic. To generate realistic background traffic, the B-Profile system pro- posed by Sharafaldin et al. [1] was employed. This system profiles the abstract behavior of human interactions, simulating the activity of 25 users based on protocols such as HTTP, HTTPS, FTP, SSH, and email.
- 2) Attack Representations: The dataset encompasses the most up-to-date common attacks, executed during specific periods to emulate real-world scenarios. The attacks included are:
- Brute Force FTP: Unauthorized access attempts targeting FTP services.
- Brute Force SSH: Repeated login attempts to breach SSH security.
- DoS (Denial of Service): Attacks aiming to render network re- sources unavailable to legitimate users.
- Heartbleed: Exploitation of the Heartbleed vulnerability in OpenSSL.



International Journal for Research in Applied Science & Engineering Technology (IJRASET)

ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 12 Issue XI Nov 2024- Available at www.ijraset.com

- Web Attacks: Including SQL injection, cross-site scripting, and other web-based attacks.
- Infiltration: Unauthorized access and compromise of network sys- tems.
- Botnet: Activities related to botnet communication and coordina- tion.
- DDoS (Distributed Denial of Service): Coordinated attacks from multiple sources to overwhelm network resources.
- 3) Data Format: The dataset includes labeled flow-based features extracted using CICFlowMeter, providing detailed information such as timestamps, source and destination IP addresses, ports, protocols, and attack labels. The data is stored in CSV files, facilitating ease of use for machine learning applications.

The inclusion of both benign and malicious activities, along with comprehen- sive feature representation, makes the CICIDS2017 dataset suitable for develop- ing and evaluating intrusion detection models under realistic network conditions.

B. Data Preprocessing

Effective data preprocessing is essential to ensure the quality and reliability of the machine learning models developed for intrusion detection. The following steps were undertaken to prepare the dataset for analysis:

1) Data Consolidation:

• The dataset, originally divided into multiple CSV files represent- ing different days and attack types, was consolidated into a single DataFrame for uniform processing.

2) Handling Null Values:

- The dataset was checked for the presence of null values using the isnull() function.
- Null Value Detection:

- Null Value Removal:
- > All records with null values were removed to prevent inconsis- tencies in model training.

df_data.dropna(inplace=True)

3) Handling Duplicate Entries:

- Duplicate records can bias the model by over-representing certain instances.
- Duplicate Detection:

duplicate_count = df_data.duplicated().sum() print(f"{
 duplicate_count} duplicate entries have
 been found in the dataset\n")

- Duplicate Removal:
- > All duplicate rows were removed to ensure each record is unique.

df_data.drop_duplicates(inplace=True)

- Index Resetting:
- After removal operations, the DataFrame indices were reset for consistency.

df_data.reset_index(drop=True, inplace=True)

4) Data Type Inspection:

• Identification of categorical and numerical features was performed.



categorical_columns = df_data.select_dtypes(include =['object']).columns.tolist() print("Categorical columns:", categorical_columns, ' \n')

- The primary categorical column identified was the 'Label' column, indicating the class of each record.
- 5) Feature and Target Separation:
- Features (X): All columns except 'Label' were considered as fea- tures.

X = df_data.drop('Label', axis=1)

• Target (y): The 'Label' column was extracted as the target variable.

y = df_data['Label'].copy()

C. Transformation into Binary and MultiClass Classifi- cation

The nature of intrusion detection problems can vary based on the specific ob- jectives. In this research, the problem was approached from both binary and multiclass classification perspectives.

1) Justification for Transforming the Problem Binary Classification

- Objective: To distinguish between normal (benign) and abnormal (ma-licious) network traffic.
- Rationale:
- > Simplifies the classification task, enabling the model to focus on de- tecting any form of intrusion.
- > Suitable for scenarios where the primary concern is to flag potential threats for further investigation.
- > Helps to address the general imbalance between normal and malicious traffic.

Multiclass Classification

- Objective: To identify and categorize specific types of network attacks.
- Rationale:
- Provides granular insights into the nature of the detected intrusions.
- > Enables tailored response strategies for different attack types.
- Addresses the complex class imbalance across multiple attack cate- gories.
- ▶ Useful for comprehensive intrusion detection systems that prioritize detailed threat analysis.

2) Implementation of Classification Configurations Binary Classification Configuration

- Label Binarization:
- > The 'Label' column was transformed to represent two classes:
- * 0 for 'Benign' traffic.
- * 1 for all types of attacks.
- Code Implementation:

Binarize labels: Map 'Benign' to 0 and all other attack labels to 1
y_b = y.map({'Benign': 0}).fillna(1)

- Class Distribution Analysis:
- > Understanding the ratio of benign to malicious traffic is crucial due to class imbalance concerns.
- Minimum Baseline Accuracy:

A naive model predicting the majority class would achieve an accuracy equal to the proportion of the dominant class (e.g., if benign traffic is 84.92%, the baseline accuracy is 84.92%).



International Journal for Research in Applied Science & Engineering Technology (IJRASET) ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538

Volume 12 Issue XI Nov 2024- Available at www.ijraset.com

The model must outperform this baseline to be considered effec- tive.

Multiclass Classification Configuration

- Label Encoding:
- All unique attack types, along with 'Benign', were encoded into in- teger labels using Label Encoding.
- Code Implementation:

from sklearn.preprocessing import LabelEncoder label_encoder =
LabelEncoder()
y_encoded = label_encoder.fit_transform(y) y_m = pd.
Series(y_encoded)

- Class Distribution Analysis:
- > The distribution of each class was examined to identify minority classes.
- Class Imbalance Consideration:
- * Some attack types constitute a very small fraction of the dataset.

* Special techniques (e.g., resampling, class weighting) may be nec- essary to ensure these classes are adequately represented during model training.

3) Addressing Class Imbalance

- Techniques Employed:
- Resampling Methods:
- * Oversampling minority classes or undersampling majority classes to balance the dataset.



Algorithmic Approaches:

* Using models that incorporate class weights (e.g., XGBoost) to penalize misclassification of minority classes more heavily.

Evaluation Metrics:

* Relying on metrics beyond accuracy, such as precision, recall, F1-score, and area under the ROC curve, to assess model performance on imbalanced data.

4) Benefits of Dual Approach

- Comprehensive Evaluation:
- > Assessing model performance across different levels of classification complexity.
- > Extracting insights on feature importance and model generalization capabilities.
- Practical Applicability:
- Catering to different operational needs, from simple intrusion detec- tion to detailed attack classification.
- > Enhancing the adaptability of the IDS in various network security scenarios.

D. Conclusion on Transformation

Transforming the intrusion detection problem into both binary and multiclass classification tasks allows for a robust evaluation of machine learning models.



It addresses different practical needs in network security, from the quick detection of any intrusion to the precise identification of attack types. By meticulously preprocessing the data and thoughtfully configuring the classification tasks, the research sets a solid foundation for developing effective and efficient intrusion detection models.

E. Feature Engineering

Feature selection is a crucial step in developing efficient and effective machine learning models, especially in high-dimensional data scenarios common in intru- sion detection systems. In this research, we performed feature selection using the Extreme Gradient Boosting (XGBoost) algorithm. XGBoost is a tree- based ensemble learning algorithm that inherently performs feature selection during its training process. This embedded method helps in selecting the most informative features while reducing the dimensionality of the input space.

1) Introduction to Feature Selection with XGBoost

XGBoost is widely used for classification and regression tasks due to its scala- bility and speed [11]. One of its key advantages is the ability to automatically compute feature importance scores during model training. These scores indicate the contribution of each feature to the predictive performance of the model.

Explanation of Feature Importance Calculation In decision tree-based algorithms like XGBoost, feature importance scores are derived from the impact of features on reducing impurity in decision trees. One common measure of impurity used in classification tasks is the Gini impurity.

Gini Impurity

Gini impurity is a measure of how often a randomly chosen element from the set would be incorrectly labeled if it were randomly labeled according to the distribution of labels in the subset. For a binary classification problem with classes 0 and 1, the Gini impurity G for a node with N samples is calculated as:

$$G = 1 - \frac{1}{i} p$$
, (1)
 $i=0$

where p_i is the probability of class *i* in the node.

Feature Importance Calculation

During the construction of decision trees, at each node, possible splits on each feature are evaluated, and the split that maximally reduces the Gini impurity is selected. The feature importance score is then calculated based on the total reduction in impurity achieved by splitting on that feature across all trees in the ensemble.

Features that contribute more to reducing the impurity (i.e., lead to greater reduction in Gini impurity when used for splitting nodes) are considered more important by the algorithm. These important features can be used for further analysis or as input to other machine learning models, such as neural networks.

Relationship Between Gini Index and Feature Importance In XG- Boost, the Gini index serves as the default metric for assessing impurity when constructing decision trees. For a node with class probabilities p_1 and p_2 (where $p_1 + p_2 = 1$), the Gini index is calculated as:

$$GI = 1 - (p^2 + p^2).$$
 (2)

Lower values of the Gini index indicate a higher degree of class separation, signifying improved model performance. During tree construction, the algo- rithm selects splits that maximally reduce the Gini index. Feature importance is determined by the total reduction in the Gini index achieved by splits involv- ing that feature across all nodes of all trees in the ensemble.



ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 12 Issue XI Nov 2024- Available at www.ijraset.com

2) Implementation of Feature Selection

The feature selection process was implemented using the following code.

```
def get_feature_importances(original_data_clf_params,printing=False):
     Calculate the feature importances using an XGBoost classifier
         trained on the provided data.
     Args:

    original_data (list): A list containing X_train, X_val.

         , X test, y train, y val, and y test data.

    clf_params (dict): Parameters for the XGBoost classifier.

    printing (bool, optional): If True, prints the feature names alongside

          their importance scores.
     Returns:
      feature_importances (array): An array containing the importance
         scores of each feature.
     .....
     # Extract original data
     X_train, X_val, X_test, y_train, y_val, y_test = original_data
     # Create XGBoost classifier instance xgb_clf =
     XGBClassifier(**clf_params)
     # Train the classifier
     xgb_clf.fit(X_train, y_train, eval_set=[(X_val, y_val)], verbose=False)
     # Get feature importances
     feature_importances = xgb_clf_feature_importances_
     if printing:
         # Print feature importances
         max_feature_name_length = max(len(name) for name inX_train.
              columns)
         print(f"{'Feature Name':<{max_feature_name_length}}\ t</pre>
              Importance")
         for name, importance in zip(X, train.columns,
                                                  ): print(f"{name:<{
              feature importances
              max_feature_name_length}}\t{
                   importance*100:.2f}%")
   return feature importances
```

International Journal for Research in Applied Science & Engineering Technology (IJRASET)



d

ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 12 Issue XI Nov 2024- Available at www.ijraset.com

The function above calculates the feature importances using an XGBoost classifier trained on the provided training data. It can optionally print out the feature names alongside their importance scores for inspection.

<pre>Select features based on a given importance threshold. Ares: original_data (list): Contains X_train, X_val, X_test, v_train, v_val, v_test data. cill_natams(dict): Parameters for the XGBoost classifier. threshold(float): Importance threshold above which features will be selected. cilling (bod), optional): if True, prints the names and importances of the selected features. A_train_selected (DataFrame): Subset of X_train containing selected features. X_train_selected (DataFrame): Subset of X_train containing selected features. X_train_selected (DataFrame): Subset of X_train containing selected features. X_train_selected (DataFrame): Subset of X_test containing selected features. X_train_tuples(list):List of tuples containing feature names and their Importances. if Extract data X_train_X_val_X_test_v_train_v_val_v_test = original_data # Extract data X_train_X_val_X_test_v_train_v_val_v_test = original_data. sif_params) # Get feature importances feature_importances = get_feature_importances(original_data. sif_params) # Get feature names feature_names = X_train_columns # Select features based onthreshold selected_features = no_where(feature_importances > thresholdllO # Create list of feature names and importances feature_tuples = [lname. importance], in _zip(feature_names[selected_features], Teature_importances[selected_features]]) # Sort features by importance sorted_feature_tuples = sorted(feature_tuples, key=lambda.xx[1], reverse =True) sorted_feature_names = [tup[0] for tup in sorted_feature_tuples] # Subset data with selected features X_train_selected features X_train_selected features x_train_selected features max[isp(name)] for name in sorted_feature_names] X_test</pre>	ef select_features_by_threshold (original_data, clf_params, threshold=0.01, printing= False): """
<pre>Arss: • original_data (list): Contains X_train, X_val, X_test, v_train, v_val, v_test data. • clf_params (dict): Importance threshold above which features will be selected. • printing (bod), optional): if True, prints the names and importances of the selected features. Returns: • X_train_selected (DataFrame): Subset of X_train containing selected features. • Insportances. • Insportances. • Insportances. • Insportances = set_feature_entraining feature names and their importances = set_feature_importances(original_data) # Get feature names feature_names = X_train_columns # Select features based onthreshold selected_features = np.where(feature_importances ></pre>	Select features based on a given importance threshold.
<pre>Returns: X_train_selected (DataFrame): Subset of X_train containing selected features. X_test_selected (DataFrame): Subset of X_test containing selected features. Icature_tuples(I)st): List of tuples containing feature names and their importances. # Extract data X_train_X_val_X_test_v_train_v_val_v_test = original_data # Get feature importances feature_importances = get_feature_importances(original_data. cll_params) # Get feature names feature_names = X_train_columns # Select feature based onthreshold selected_features = np.where(feature_importances > threshold!00 # Create list of feature names and importances [eature_tuples = [iname, importance] of name, importance], ip 2ip(feature_names[selected_features], feature_importances[selected_features])) # Sort features by importance sorted_feature_s bill of tup in sorted_feature_tuples] # Subset data with selected features X_train_selected = X_train.loc(); sorted_feature_names] X_val_selected = X_val.loc(; sorted_feature_names] X_test_selected = X_test.loc(); sorted_feature_names] # Subset data with selected features X_train_selected features X_train_selected features X_train_selected features X_train_selected features X_train_selected features X_train_selected features X_train_selected features X_train_selected features max_feature_names]X_test_selected = X_test.loc(); sorted_feature_names] y printing: # Print selected features max_feature_names]X_test_selected = X_test.loc(); sorted_feature_names] printing: # Print selected features max_feature_names] = top() for name in sorted_feature_names) = top() for name in sorted_feature_names) = top() for name in sorted_feature_names) = top() for name in sorted_feature_names) = top() for name in sorted_feature_names) = top() for name in sorted_features) = top() for name in sorted_feature_names) = top</pre>	 Ares: original_data (list): Contains X_train, X_val, X_test, v_train, v_val, v_test data. Clf_params (dict): Parameters for the XGB cost classifier. threshold (float): Importance threshold above which features will be selected. printing (bool, optional): If True, prints the names and importances of the selected features.
<pre>## Extract data # Extract data X_train_X_val, X_test v_train_v_val, v_test = original_data # Get feature importances feature_importances = get_feature_importances(original_data,</pre>	 Returns; X_train_selected (DataFrame): Subset of X_train containing selected features. X_val_selected (DataFrame): Subset of X_val containing selected features. X_test_selected (DataFrame): Subset of X_test containing selected features. feature tuples (list): List of tuples containing feature names and their importances.
<pre># Get feature importances feature_importances = get_feature_importances[original_data,</pre>	# Extract data X_train, X_val, X_test, v_train, v_val, v_test = original_data
<pre># Get feature names feature_names = X_train_columns # Select features based onthreshold selected_features = np.where(feature_importances > threshold)[0] # Create list of feature names and importances feature_tuples = [[name, importance] for name, importance_int zip(feature_names[selected_features], feature_importances[selected_features]]] # Sort features by importance sorted_feature_tuples = sorted[feature_tuples, key=lambda x:x[1], reverse =True) sorted_feature_names = [tup[0] for tup in sorted_feature_tuples] # Subset data with selected features X_train_selected = X_train.loc[: sorted_feature_names] X_val_selected = X_val.loc[: , sorted_feature_names] X_test_selected = X_test.loc[: sorted_feature_names] if printing: # Print selected features max_feature_name_length = max[len(name) for name in sorted_feature_names] </pre>	<pre># Get feature importances feature_importances(original_data,</pre>
<pre>feature_names = X_train_columns # Select features based onthreshold selected_features = np.where(feature_importances > threshold[0] # Create list of feature names and importances feature_tuples = [(name, importancein_zip(feature_names[selected_features], feature_importances[selected_features])) # Sort features by importance sorted_feature_tuples = sorted(feature_tuples, key=lambda x:x[1], reverse =True) sorted_feature_names = [tup[0] for tup in sorted_feature_tuples] # Subset data with selected features X_train_selected = X_train_loci; sorted_feature_names] X_val_selected = X_val_loci; , sorted_feature_names] X_test_selected = X_test_loci; sorted_feature_names] if printing: # Print_selected_features max_feature_names_length = max[len(name) for name in sorted_feature_names] </pre>	# Get feature names
<pre># Select features based on threshold selected_features = np.where(feature_importances > threshold)[0] # Create list of feature names and importances feature_tuples = [[name, importance] for name, importancein_zip(feature_names[selected_features], feature_importances[selected_features]]] # Sort feature_tuples = sorted[feature_tuples, key=lambda x:x[1], reverse =True) sorted_feature_names = [tup[0] for tup in sorted_feature_tuples] # Subset data with selected features X_train_selected = X_train.loc[:, sorted_feature_names] X_val_selected = X_val.loc[:, sorted_feature_names] X_test_selected = X_test.loc[:, sorted_feature_names] if printing: # Print_selected_features max_feature_name_length = max(len(name) for name in sorted_feature_names) </pre>	feature_names = X_train.columns
<pre># Create list of feature names and importances feature_tuples = [(name, importancein_zip(feature_names[selected_features], feature_importances[selected_features])] # Sort features by importance sorted_feature_tuples = sorted(feature_tuples, key= lambda x:x[1], reverse =True) sorted_feature_names = [tup[0] for tup in sorted_feature_tuples] # Subset data with selected features X_train_selected = X_train.loc[:, sorted_feature_names] X_val_selected = X_val.loc[: , sorted_feature_names] X_test_selected = X_test_loc[:, sorted_feature_names] if printing: # Print_selected_features max_feature_name_length = max(len(name) for name in</pre>	# Select features based onthreshold <u>selected_features</u> = np.where(feature_importances > threshold)[0]
<pre>feature_importances[selected_features])] # Sort features by importance sorted_feature_tuples = sorted(feature_tuples, key=lambda x:x[1], reverse =True) sorted_feature_names = [tup[0] for tup in sorted_feature_tuples] # Subset data with selected features X train_selected = X train_loc[:, sorted_feature_names] X val_selected = X val_loc[:, sorted_feature_names] X test_selected = X test_loc[:, sorted_feature_names] # Print_selected_features max_feature_name_length = max(len(name) for name in</pre>	<pre># Create list of feature names and importances feature_tuples = [(name, importance) for name, importanceinzip(feature_names[selected_features],</pre>
<pre># Sort features by importance sorted_feature_tuples = sorted(feature_tuples, key=lambda x:x[1], reverse =True) sorted_feature_names = [tup[0] for tup in sorted_feature_tuples] # Subset data with selected features X_train_selected = X_train_loc[:, sorted_feature_names] X_val_selected = X_val_loc[: , sorted_feature_names] X_test_selected = X_test_loc[:, sorted_feature_names] if printing: # Print_selected_features max_feature_name_length = max(len(name) for name in sorted_feature_names)</pre>	feature_importances[selected_features])]
<pre>sorted_teature_names = [tup[0] for tup in sorted_teature_tuples] # Subset data with selected features X_train_selected = X_train_loc[:, sorted_feature_names] X_val_selected = X_val_loc[:, sorted_feature_names] X_test_selected = X_test_loc[:, sorted_feature_names] if printing: # Print selected features max_feature_name_length = max(len(name) for name in sorted_feature_names) print("The number of selected features in the selected feature_names) </pre>	<pre># Sort features by importance sorted_feature_tuples, key=lambda x:x[1], reverse =True)</pre>
<pre># Subset data with selected features X_train_selected = X_train_loc[:, sorted_feature_names] X_val_selected = X_val_loc[:, , sorted_feature_names] X_test_selected = X_test_loc[:, sorted_feature_names] if printing: # Print selected features max_feature_name_length = max(len(name) for name in sorted_feature_names) print("The number of selected features in the feature_names) </pre>	sorted_feature_names = [tup[0] for tup in sorted_feature_tuples]
<pre>if printing: # Print selected features max_feature_name_length = max(len(name) for name in sorted_feature_names) print("The number of selected features are lies(served features)</pre>	# Subset data with selected features X_train_selected = X_train_loc[:, sorted_feature_names] X_val_selected = X_val_loc[: , sorted_feature_names] X_test_selected = X_test_loc[:, sorted_feature_names]
print(f"{'Feature Name':< <u>{max_feature_name_length}}</u> \ timportance\n")	<pre>if printing: # Print selected features max_feature_name_length = max(len(name) for name in sorted_feature_names) print(f"The number of selected features are: {len(sorted_feature_names)}\n") print(f"{'Feature_Name':<{max_feature_name_length}}} tImportance\n")</pre>



This function selects features whose importance scores exceed a specified threshold.

3) Feature Selection Under Binary and Multiclass Configurations

We performed feature selection separately under binary and multiclass configu- rations.

Binary Configuration We defined the classifier parameters and executed the feature selection process:

<pre># Define classifier parameters clf_params_b = dict(objective='binaty:logistic',</pre>	
n_estimators=50, eval_metric='	
learning_rate=0.1, subsample=0.8,	
colsample_bytree=0.8, reg_alpha =0.1_early_stopping_rounds=10	
random_state=seed_pi	
)	
<pre>#Perform feature selection print("Under Binary Configuration:\n")</pre>	
X_train_selected_b, X_val_selected_b, X_test_selected_b, feature_tuples_b = select_features_by_threshold (original_	data_b,
=0.01, printing=True)	

Multiclass Configuration Similarly, for the multiclass classification task:

I

```
# Define number of classes num_classes = len(
class_counts_m)
# Define classifier parameters clf_params_m
= dict(
     objective='multi:softmax',
     num_class=num_classes,
     n_estimators=50, eval_metric='
     mlogloss', max_depth=3,
     learning_rate=0.1, subsample=0.8
     , colsample_bytree=0.8,
     reg_alpha=0.1,
     early_stopping_rounds =10,
     random_state=seed_pi
# Perform feature selection
print("Under Multiclass Configuration:\n") X_train_selected_m,
X val selected m, X test selected m,
    feature_tuples_m = select_features_by_threshold ( original_data_m,
     clf_params=clf_params_m, threshold
         =0.01, printing=True)
```



ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538 Volume 12 Issue XI Nov 2024- Available at www.ijraset.com

4) Visualization of Feature Importances

We plotted the feature importances for both configurations to visualize the significance of each feature.

```
# Define figure and subplots
fig, axs = plt.subplots(1, 2, figsize=(16, 8))
# Binary configuration
feature_names_b = [tup[0] for tup in feature_tuples_b] significance_scores_b
= [tup[1] for tup in feature_tuples_b] axs[0].barh(feature_names_b,
significance_scores_b, color="
    skyblue
                               ') axs[0].set_xlabel
('Feature Importance') axs[0].set_ylabel('Features')
axs[0].set_title("Feature Importance under Binary Classification")
axs[0].invert_yaxis()
# Multiclass configuration
feature_names_m = [tup[0] for tup in feature_tuples_m]
significance_scores_m = [tup[1] for tup in feature_tuples_m] axs[1].barh(
feature_names_m, significance_scores_m, color='
    red')
axs [1]. set_xlabel (' Feature
                                   Importance ') axs
                                                         [1].
                               ') axs[1].set_title("Feature
set_ylabel (' Features
Importance under Multiclass
    Classification") axs[1].
invert_yaxis()
# Adjust layout and display plt .
tight_layout () plt.show()
```

5) Analysis of Unique Features

We identified features that were unique to each configuration:

Features unique to binary configuration
unique_in_b = set(feature_names_b) - set(feature_names_m)
Features unique to multiclass configuration
unique_in_m = set(feature_names_m) - set(feature_names_b)

6) Creation of New Features

While the primary focus was on selecting existing features based on their impor- tance scores, we also explored the creation of new features that could potentially enhance model performance. This involved combining existing features or com- puting statistical measures that may capture underlying patterns in the data.







The intended impact of creating new features was to provide the model with additional information that could improve its ability to distinguish between benign and malicious traffic, as well as between different types of attacks in the multiclass configuration.

F. Conclusion on Feature Engineering

Feature selection using XGBoost allowed us to identify and retain the most infor- mative features, effectively reducing the dimensionality of the dataset without compromising model performance. The reduction in dimensionality leads to lower computational costs and potentially improves the generalization of the model. The differences in selected features between the binary and multiclass configurations highlight the importance of context when performing feature se- lection.

G. Outlier Analysis

Note that outlier analysis is performed only for the **Binary Configuration**. We avoided this procedure for the multiclass setup to save time, as we expect analogous results. Since there is no outlier-handling procedure in the multiclass setting, this analysis serves to enhance our understanding of the problem.

1) Definition of Outlier

An outlier in statistics is a data point that deviates significantly from the over- all pattern of the remaining data. Mathematically, there isn't a universally accepted definition, but two common approaches are often used:

Z-scores We can define an outlier based on its Z-score, which measures how many standard deviations a data point is away from the mean:



Where:

• *x* is the data point in question,

Interquartile Range (IQR)

- μ is the population mean (often estimated by the sample mean), and
- σ is the population standard deviation (often estimated by the sample standard deviation).

We can then define a threshold for Z-scores (e.g., ± 3 standard deviations). Data points exceeding this threshold in absolute value can be considered outliers.

This approach utilizes the quartiles of the data distribution:

- Q1 (First quartile): represents the 25th percentile (value below which 25% of the data lies).
- Q3 (Third quartile): represents the 75th percentile (value below which 75% of the data lies).

We calculate the Interquartile Range (IQR):

$$IQR = Q3 - Q1$$

The lower and upper bounds are defined as:

Lower bound = $Q1 - 1.5 \times IQR$ Upper bound = $Q3 + 1.5 \times IQR$

(4)

Outliers can be defined as data points falling outside these bounds. Data points below the lower bound or exceeding the upper bound can be considered potential outliers.

(5)

Since the IQR method is robust against skewed distributions, while Z-scores are better for normally distributed data, the IQR approach will be adopted in this analysis.



2) Implementation of Outlier Detection

```
if.'X.final.b' not in globals():
      # Concatenate 'X_train_final_b', 'X_val_final_b', and 'X_test_final_b'
      # into a DataFrame named X. final. b
      X_final_b = pd.concat([X_train_final_b, X_val_final_b, X_test_final_b])
# Define function for IQR outlier detection def
detect_outliers_igr(df):
      Detects outliers using the Interquartile Range (JQR) method.
      Args:
           df (pandas.DataFrame): The dataframe containing the features.
      Returns:
           pandas DataFrame: A new dataframe with additional columns
                indicating outliers.
      # Extract columns
      columns = df.columns.tolist()
      # Create a new DataFrame to store the outliers outliers = pd.Data
      Frame(columns=[col + '_outliers' for
          col in columns) for col
      in columns:
           Q1 = df[col].guantile(0.25)Q3 = df[
           col].guantile(0.75)
           IQR = Q3 - Q1
           lower_bound = Q1-1.5 * IQR upper_bound
           = Q3 + 1.5 * IQR
           # If they fall outside of boundaries, they are defined as
                outliers
           outliers[col + '_outliers'] = df[col].apply(lambda x
               :
                                                1 if (x < lower, bound or x >
                                                    upper_bound) else 0)
      return outliers
# Detect outliers
df_with_outliers = detect_outliers_igr(X_final_b.copy()) # .sum() will
create the requested result
print("Number of outlier samples for every column:") print("[for a
concatenation of X train final b,
     X_val_final_b, X_test_final_b dataframes] \n")
print(df_with_outliers.sum())
```

We implemented the outlier detection using the IQR method on our dataset. The following code demonstrates how outliers are detected:



3) Outlier Boxplot Visualization

To visualize the distribution of features and the presence of outliers, we used boxplots. The following function generates boxplots for the specified features:

```
def visualize_boxplots(df, features, figsize=(15, 28)): """
     Visualizes boxplots for the specified features using a grid layout.
     Args:
           df (pandas.DataFrame): The dataframe containing the features.
           features (list): A list of feature names to visualize.
           figsize (tuple, optional): The size of the figure.
                Defaults to (15, 28).
     .....
     n_features = len(features)
                                            # Example: 20 features
     rows = int(np.ceil(n_features / 3))
                                                        # Calculate number of
          rows for a 3-column grid
     fig, axes = plt.subplots(rows, 3, figsize=figsize)
     # Flatten the axes array for easy iteration axes_flat = axes.ravel
     ()
     for i, feature in enumerate(features): axes_flat[i].boxplot(
           df[feature]) axes_flat[i].set_title(f"{feature} Boxplot")
          axes_flat[i].set_ylabel("Value")
     plt.tight_layout() plt.show
# Visualize boxplots
visualize_boxplots(X_final_b, X_final_b.columns.tolist())
```

Interpretation of Boxplots The boxplot elements provide insights into the distribution of data along the y-axis:

- Whiskers: Extend from the first quartile (Q1) to the third quartile (Q3), indicating the range where most data points lie (within 1.5 times the IQR).
- Box: Represents the Interquartile Range (IQR), emphasizing the central 50% of the data distribution.
- Median: Depicted as a horizontal line within the box, indicates the value that separates the lower and upper halves of the data.
- Outliers: Individual data points plotted beyond the whiskers, may signify values significantly deviating from the primary distribution.



© IJRASET: All Rights are Reserved | SJ Impact Factor 7.538 | ISRA Journal Impact Factor 7.894 |



International Journal for Research in Applied Science & Engineering Technology (IJRASET) ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538

Volume 12 Issue XI Nov 2024- Available at www.ijraset.com

4) Conclusion

Our statistical analysis identified a significant number of outliers, about 25% of the data. To address these outliers, efficient data scaling techniques will be preferred over traditional methods like winsorization, in order to minimize potential information loss. While other outlier detection methods like Isolation Forest and Local Outlier Factor exist, their high computational costs make them impractical for the large size of the current dataset.

REFERENCES

- [1] Sharafaldin, I., Lashkari, A. H., & Ghorbani, A. A. (2018). Toward generat- ing a new intrusion detection dataset and intrusion traffic characterization. In *Proceedings* of the 4th International Conference on Information Systems Security and Privacy (ICISSP) (pp. 108–116).
- [2] Scarfone, K., & Mell, P. (2007). Guide to Intrusion Detection and Preven- tion Systems (IDPS). NIST Special Publication, 800(2007), 94.
- [3] Ahmed, M., Mahmood, A. N., & Hu, J. (2016). A survey of network anomaly detection techniques. *Journal of Network and Computer Appli- cations*, 60, 19–31.
- [4] Sommer, R., & Paxson, V. (2010). Outside the closed world: On using machine learning for network intrusion detection. In 2010 IEEE Symposium on Security and Privacy (pp. 305–316). IEEE.
- [5] Hall, M. A. (1999). Correlation-based feature selection for machine learning. PhD Thesis, University of Waikato.
- [6] Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. Journal of Machine Learning Research, 3(Mar), 1157–1182.
- [7] Wang, W., Sheng, Y., Wang, J., Zeng, X., Ye, X., & Huang, Y. (2018). HAST-IDS: Learning hierarchical spatial-temporal features using deep neu- ral networks to improve intrusion detection. *IEEE Access*, 6, 1792–1806.
- [8] Maciá-Fernández, G., Garc'ia-Teodoro, P., & Mirsky, Y. (2022). Establish- ing the Contaminating Effect of Metadata Feature Inclusion in Machine- Learned Network Intrusion Detection Models. In *Engineering Secure Soft- ware and Systems* (pp. 21–37). Springer. Retrieved from https://link. springer.com/chapter/10.1007/978-3-031-09484-2_2
- [9] Yin, C., Zhu, Y., Fei, J., & He, X. (2017). A deep learning approach for intrusion detection using recurrent neural networks. *IEEE Access*, 5, 21954–21961.
- [10] Ullah, I., & Mahmoud, Q. H. (2020). A hybrid model for anomaly-based intrusion detection in software-defined networks. *Journal of Network and Computer Applications*, 157, 102563.
- [11] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting sys- tem. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 785–794).
- [12] Sun, Y., Kamel, M. S., Wong, A. K., & Wang, Y. (2007). Cost-sensitive boosting for classification of imbalanced data. Pattern Recognition, 40(12), 3358–3378.
- [13] Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., ... & Liu, T. Y. (2017). LightGBM: A highly efficient gradient boosting decision tree. In Advances in Neural Information Processing Systems, 30.
- [14] Wang, C., Pang, Y., Li, Y., & Yuan, F. (2020). LightGBM: A novel en- semble boosting model for accurate soccer result prediction based on book- maker odds. *Entropy*, 22(4), 437.
- [15] Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., & Gulin, A. (2018). CatBoost: Unbiased boosting with categorical features. In Advances in Neural Information Processing Systems, 31.
- [16] Breiman, L. (2001). Random forests. Machine Learning, 45(1), 5-32.
- [17] Liaw, A., & Wiener, M. (2002). Classification and regression by random- Forest. R News, 2(3), 18–22.
- [18] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. In Advances in Neural Information Processing Systems, 30.
- [19] Molnar, C. (2020). Interpretable Machine Learning. Retrieved from https://christophm.github.io/interpretable-ml-book/shapley.html
- [20] DataCamp. (2023). Pandas Profiling (ydata-profiling) in Python: A Guide for Beginners. Retrieved from https://www.datacamp.com/tutorial/ pandas-profilingydata-profiling-in-python-guide
- [21] Ross, K. D. (2021). Quantile functions. In *Probability and Statistics* (pp. 111–119). Open Educational Resource. Retrieved from https://bookdown. org/kevin_davisross/probsim-book/quantile-functions.html











45.98



IMPACT FACTOR: 7.129







INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 🕓 (24*7 Support on Whatsapp)