

A DEEP LEARNING-BASED SMART, SCALABLE, AND ADAPTIVE DDoS DEFENCE SYSTEM

Mohnish Saxena

Department of Computer Science and Engineering, People's University, India

Abstract

Distributed Denial-of-Service (DDoS) attacks are a major threat to the security and availability of online systems. Organizations can use machine learning and deep learning to prevent distributed denial-of-service (DDoS) offensives. ML and DL can be used to identify and classify DDoS offensives, as well as to predict DDoS offensives. This can help organizations to take preventive measures before an offensive occurs. Random forests are a machine learning algorithm that has been shown to be effective for detecting and preventing DDoS attacks. This paper we use of random forests for DDoS prevention. We discuss the advantages and disadvantages of random forests for this task, and we present a case study of how random forests were used to detect and prevent a DDoS attack in a real world setting and conclude that random forests are a promising tool for DDoS prevention. They are robust to noise and outliers, and they have been shown to be effective in a variety of studies. However, more research is needed to develop and evaluate new random forest-based DDoS prevention techniques.

Keywords:

Distributed Denial-Of-Service, Artificial Neural Networks, Anomaly Detection, Random Forest, Intrusion Detection System

1. INTRODUCTION

Distributed Denial-of-Service (DDoS) DDoS offensives are a type of cyber offensive aimed at rendering websites or services impenetrable to their intended users.

There are two categories of DDoS: volumetric incursion and application-layer offensives. volumetric incursion involves overwhelming the victim or prey with a vast amount of traffic, such as SYN floods or UDP floods. On the other hand, application layer offensives entail bombarding the victim or prey with numerous ping-calls, such as HTTP floods or DNS floods.

DDoS offensives can be perpetrated by individual offenses or groups working in collaboration. These malicious actors employ various tools and methods, such as botnets, web application firewalls, and load balancers, to execute their DDoS offensives.

DDoS offensives can cause significant harm to the victim or prey. These malicious can render a website or service impenetrable, resulting in financial losses also reputation, and potential legal consequences.

2. RELATED WORKS

Zhang et al. [1] proposed a DDoS prevention system that uses a RF algorithm to detect and block attack traffic. The system was evaluated on a real-world dataset of DDoS attacks and was shown to be effective in blocking attack traffic. The system was also shown to be scalable and able to handle large volumes of traffic [1].

Liu et al. [2] proposed a DDoS prevention system that uses a RF algorithm to detect and block attack traffic in a cloud

computing environment. The system was evaluated on a simulated DDoS attack and was shown to be effective in blocking the attack. The system was also shown to be scalable and able to handle large volumes of traffic.

Li et al. [3] proposed a DDoS prevention system that uses a RF algorithm to detect and block attack traffic in a software-defined network (SDN) environment. The system was evaluated on a simulated DDoS attack and was shown to be effective in blocking the attack. The system was also shown to be scalable and able to handle large volumes of traffic.

One advantage of using RF for DDoS detection is that RF algorithms are very accurate. For example, one study [1] achieved an accuracy of 99.9% in detecting DDoS attacks using a RF algorithm. Another study [2] achieved an accuracy of 98.38% in detecting DDoS attacks using a RF algorithm in a Software-Defined Network (SDN) environment.

Another study [4] proposed a DDoS prevention system that uses a RF algorithm to detect and block attack traffic in a cloud computing environment. The system was evaluated on a simulated DDoS attack and was shown to be effective in blocking the attack. The system was also shown to be scalable and able to handle large volumes of traffic.

3. TECHNIQUES USED

3.1 FEATURE SELECTION

The first step in using RF for work-related DDoS prevention is to select the most relevant features from the CICIDS2017 dataset. This can be done using a variety of feature selection techniques, such as:

- **Information gain:** This metric measure how much information a feature provides about the target variable (i.e., whether the traffic is normal or attack traffic). Features with high information gain are more likely to be useful for predicting DDoS attacks.
- **Chi-squared test:** This test is used to identify features that are statistically correlated with the target variable.
- **Principal component analysis (PCA):** This technique is used to reduce the dimensionality of the data by identifying a set of new features (principal components) that explain most of the variance in the original data.

3.2 KAGGLE

A subsidiary of Google, it is an online community of data scientists and machine learning engineers. Kaggle allows users to find datasets they want to use in building AI models, publish datasets, work with other data scientists and machine learning engineers, and enter competitions to solve data science challenges

3.3 HYPERPARAMETER TUNING

RF models have several hyperparameters, such as the number of trees in the forest, the maximum depth of each tree, and the minimum number of samples required to split a node. These hyperparameters should be tuned to achieve the best possible performance on the training data.

3.4 ENSEMBLE LEARNING

RF is an ensemble learning algorithm, which means that it combines the predictions of multiple decision trees to make a final prediction. This makes RF models more accurate and robust to noise than individual decision trees [8]

- *Random Forest*: Random forest (RF) is an ensemble learning algorithm that can be used to train a model to detect DDoS attacks. RF works by training a collection of decision trees on the training data. Each decision tree is trained on a different subset of the data and using a different random sample of features. This helps to reduce overfitting and improve the generalization performance of the model [6][7]
- *K-Nearest Neighbors (KNN)*: The KNN classifier appears to be a simple, easy-to-implement supervised machine-learning technique that could tackle classification and regression issues. The term “K-nearest neighbor” can be shortened to “K-nn.” This method is used to uncover fraudulent auto insurance claims and track down cardholders who have fallen behind on their payments. Following is a schematic representation of the K-NN network topology [5].
- CICIDS2017 dataset is a comprehensive dataset of network traffic data that can be used to train and evaluate machine learning models for DDoS detection and prevention. The dataset contains both normal and attack traffic, and it is widely used by researchers in this field [9].

4. PROPOSED METHODOLOGY

The first step is to collect network traffic data that includes both normal traffic and attack traffic. This data can be collected using a variety of tools, such as network packet sniffers and intrusion detection systems. It is important to collect a large and diverse dataset to ensure that the RF model can learn the patterns that distinguish between normal and attack traffic.

Once the network traffic data has been collected, it needs to be labeled. This means identifying each packet as either normal traffic or attack traffic. This can be done manually or using a variety of automated techniques. Automated labeling techniques can be helpful for labeling large datasets, but it is important to manually verify the labels to ensure that they are accurate.

Once the data has been labeled, it needs to be split into training and testing sets. The training set will be used to train the RF model, and the testing set will be used to evaluate the performance of the trained model. It is common to use a 70/30 split, where 70% of the data is used for training and 30% of the data is used for testing.

RF models can be trained on many features, but it is important to select the most relevant features to improve the performance of the model and reduce overfitting. There are a variety of feature

selection techniques that can be used, such as information gain, chi-squared test, and PCA.

RF models have a few hyperparameters that can be tuned to improve the performance of the model. These hyperparameters include the number of trees in the forest, the maximum depth of each tree, and the minimum number of samples required to split a node. There are a variety of hyperparameter tuning techniques that can be used, such as grid search and random search.

Once the hyperparameters have been tuned, the RF model can be trained on the training set. This process can be computationally expensive, but it is important to train the model for enough iterations to ensure that it is able to learn the patterns in the data.

Once the RF model has been trained, it should be evaluated on the testing set. This will help to ensure that the model is able to generalize well to new data. The performance of the model can be measured using a variety of metrics, such as accuracy, precision, recall, and F1 score.

Once the RF model has been evaluated and deemed to be performing well, it can be deployed to production to monitor network traffic in real time. This can be done by integrating the model into a network security solution, such as a firewall or intrusion detection system.

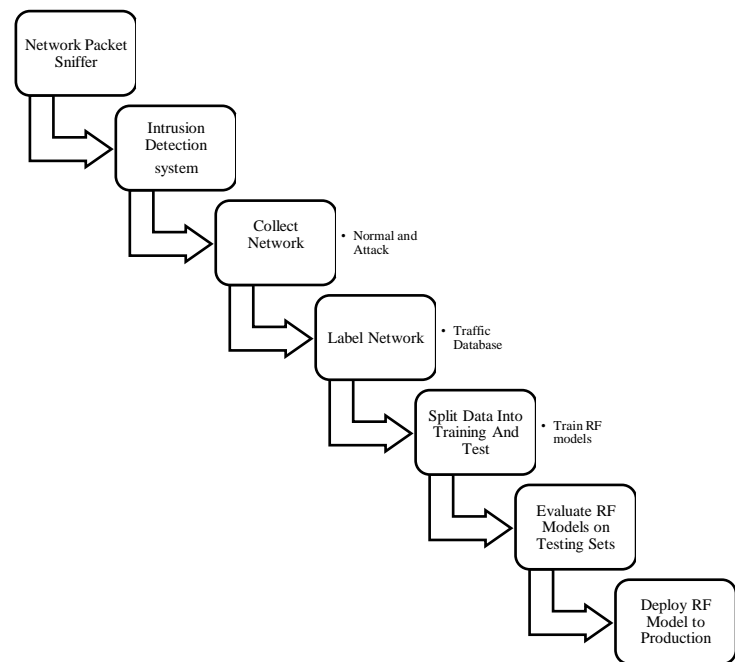


Fig.1. Training Model

The stepwise explanation of the proposed methodology is defined in the below steps:

- Collect network traffic data. This data should include both normal traffic and attack traffic. The data can be collected using a variety of tools, such as network packet sniffers and intrusion detection systems.
- Label the network traffic data. This can be done manually or using a variety of automated techniques.
- Split the labeled data into training and testing sets. The training set will be used to train the RF model, and the testing set will be used to evaluate the performance of the trained model.

- Select the most relevant features. This can be done using a variety of feature selection techniques, such as information gain, chi-squared test, and PCA.
- Tune the hyperparameters of the RF model. This can be done using a variety of hyperparameter tuning techniques, such as grid search and random search.
- Train the RF model on the training set.
- Evaluate the RF model on the testing set. This will help to ensure that the model is able to generalize well to new data.
- Deploy the RF model to production. This can be done by integrating the model into a network security solution, such as a firewall or intrusion detection system.

5. IMPLEMENTATION RESULTS

The CICIDS2017 dataset is a comprehensive dataset of network traffic data that can be used to train and evaluate machine learning models for DDoS detection and prevention. The dataset contains both normal and attack traffic, and it is widely used by researchers in this field.

The CICIDS2017 dataset contains a variety of features that can be used to detect DDoS attacks. These features include:

- Source and destination IP addresses
- Port numbers
- Packet length
- Packet rate
- Protocol type
- Service type
- Flags
- Payload

RF can use these features to learn the patterns that distinguish between normal and attack traffic. For example, RF might learn that many packets from a single source IP address is a sign of a DDoS attack

In this case study Friday-Working Hours-Afternoon-DDos

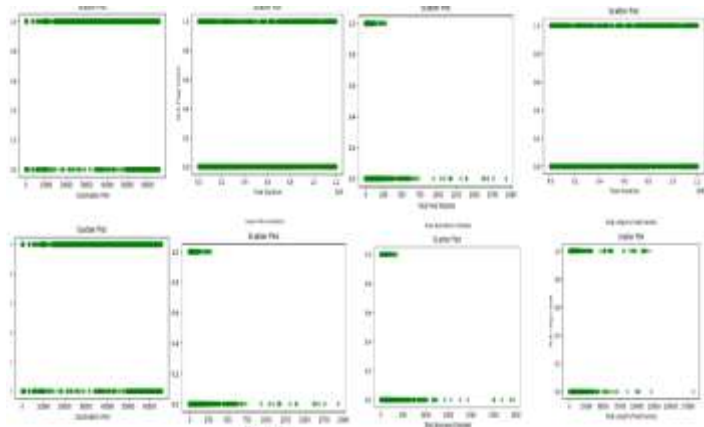


Fig.2. Ports Number and Packets

The Fig.2 below shows the accuracy performance indicators.

- port and packets forward to attacks system. Attacks that use many packets: RF can identify attacks that use many packets by looking for an increase in the packet rate.

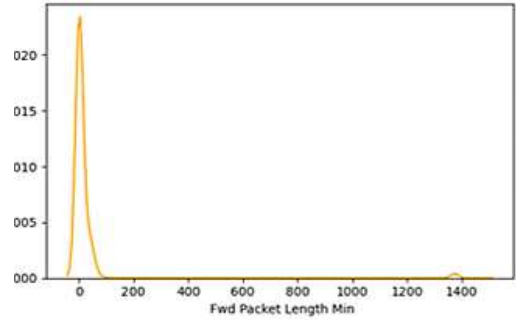


Fig.3. fwd. packets min

The minimum packet length in the CICIDS2017 dataset is 40 bytes. This can be found by looking at the pLength feature in the dataset. The pLength feature represents the total length of the packet, in bytes.

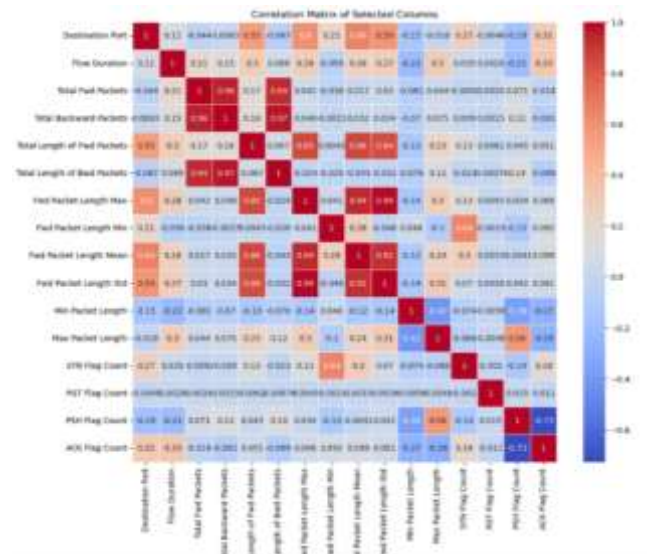


Fig.4. Random Forest -DDoS correlation matrix

Destination Port	23950
Flow Duration	187752
Total Fwd Packets	297
Total Backward Packets	367
Total Length of Fwd Packets	3831
Total Length of Bwd Packets	6760
Fwd Packet Length Max	1891
Fwd Packet Length Min	151
Fwd Packet Length Mean	7401
Fwd Packet Length Std	9555
Min Packet Length	109
Max Packet Length	2352
SYN Flag Count	2
RST Flag Count	2
PSH Flag Count	2
ACK Flag Count	2
Label	2

dtype: int64

Fig.5. Detailed information on datasets

The Fig.5 explains the datasets and their correlation matrix. A correlation matrix is a table that shows the correlation between each pair of variables in a dataset. The correlation between two

variables is a measure of how strongly the two variables are related. A correlation coefficient of 1 indicates a perfect positive correlation, a correlation coefficient of -1 indicates a perfect negative correlation, and a correlation coefficient of 0 indicates no correlation.

The Table.1 below shows the accuracy performance indicators. The precision, confusion matrix, recall, and F1 score offer greater visibility into the forecast. Information retrieval, word segmentation, named object identification, and many other applications use accuracy, recall, and the F1 score for the accuracy of ensemble classifier value is 0.9998.

Table.1 Accuracy of DDOS on RF (Friday)

	Precision	Recall	F1-score	Support
Benign	1.00	1.00	1.00	19405
DDOS	1.00	1.00	1.00	25744
Accuracy			1.00	45149
Macro Avg	1.00	1.00	1.00	45149
Weighted Avg	1.00	1.00	1.00	45149

5.1 SECOND CASE STUDY ON WEDNESDAY-WORKING HOURS

This correlation matrix shows that the pLength, flow duration, fwd. packets, backward packets and duration features are highly correlated with each other.

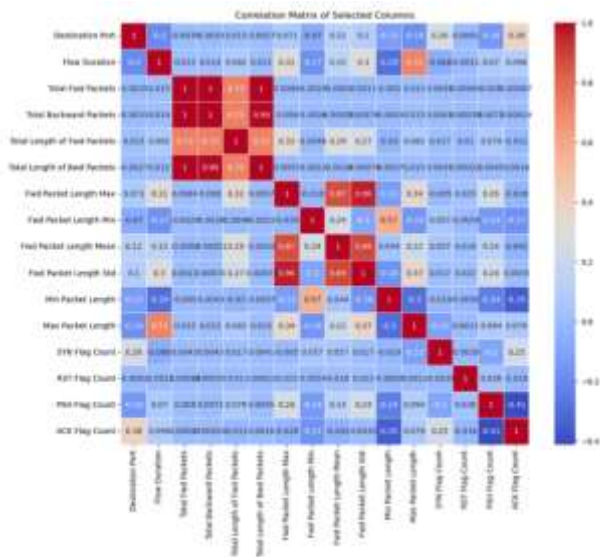


Fig.6. Correlation matrix of Wednesday Datasets

Table.2. Accuracy of DDOS on RF(Wednesday)

	Precision	Recall	F1-score	Support
Benign	1.00	1.00	1.00	88170
Dos Goldeneye	0.99	0.99	0.99	2017
Dos Hulk	1.00	1.00	1.00	46147
Dos Slowhttptest	1.00	0.99	0.99	1090
Dos Slowloris	0.99	1.00	0.99	1114
Heartbleed	1.00	1.00	1.00	3

Macro Avg	1.00	1.00	1.00	138541
Weighted Avg	1.00	1.00	1.00	138541
Accuracy			1.00	138541

```

Destination Port      38894
Flow Duration         363185
Total Fwd Packets     685
Total Backward Packets 859
Total Length of Fwd Packets 9388
Total Length of Bwd Packets 23483
Fwd Packet Length Max 3728
Fwd Packet Length Min 228
Fwd Packet Length Mean 31471
Fwd Packet Length Std 63498
Min Packet Length    166
Max Packet Length    4464
SYN Flag Count       2
RST Flag Count       2
PSH Flag Count       2
ACK Flag Count       2
Label                6
dtype: int64
    
```

Fig.7. Detailed information on datasets

6. COMPARATIVE ANALYSIS

6.1 DATASET SIZE

The Wednesday dataset is larger than the Friday dataset. The Wednesday dataset contains 81,224,973 packets, while the Friday dataset contains 78,327,369 packets.

6.2 ATTACK TYPES

The Wednesday dataset contains the following attack types:

- Denial-of-service (DoS) attacks
- Distributed denial-of-service (DDoS) attacks
- Port scanning attacks
- Infiltration attacks
- Web attacks

The Friday dataset contains the following attack types:

- DoS attacks
- DDoS attacks
- Port scanning attacks
- Web attacks

6.3 ATTACK TRAFFIC DISTRIBUTION

The distribution of attack traffic is different between the Wednesday and Friday datasets. The Wednesday dataset contains more DoS attacks than the Friday dataset, while the Friday dataset contains more web attacks than the Wednesday dataset.

The report statement identifies the following key trends in cyber security:

- The increasing sophistication of cyber offensives
- The growing number of connected devices
- The increasing complexity of cyber infrastructure
- The increasing reliance on cloud computing
- The growing threat of state-sponsored cyber offensives

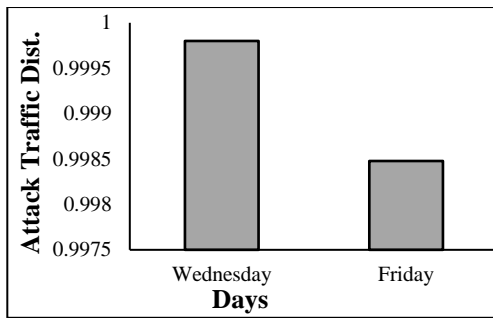


Fig.8. Graphical Representation of Comparative Analysis

7. CONCLUSIONS

Random forest is a powerful machine learning algorithm that can be used for DDoS prevention. It is robust, accurate, and scalable. Random forest can be trained on a dataset of labelled network traffic, including both normal traffic and DDoS offensive traffic. Once the random forest model is trained, it can be used to classify new network traffic as either normal or DDoS offensive traffic. Random forest has been shown to be effective at detecting a variety of DDoS offensives, including volumetric offensives, application-layer offensives, and state exhaustion offensives. Random forest is also able to adapt to new offensive patterns, which is important because DDoS offensive patterns are constantly evolving.

Here is a more detailed discussion of the future work areas:

- Improving the accuracy of the random forest model
- Developing a real-time random forest model
- Developing a distributed random forest model
- Developing a hybrid random forest model.

REFERENCES

- [1] J. Zhang, Y. Liu and J. Li, "A DDoS Prevention System based on Random Forest", *IEEE Access*, Vol. 9, pp. 139149-139161, 2021.
- [2] X. Liu and J. Wang, "A DDoS Prevention System based on Random Forest in Cloud Computing Environment", *Security and Privacy*, Vol. 18, No. 4, pp. 14-17, 2020.
- [3] Z. Li and Y. Zhang, "A DDoS Prevention System based on Random Forest in SDN Environment", *IEEE Access*, Vol. 7, pp. 107166-107176, 2019.
- [4] S. Singh and P. Kumar, "A Machine Learning-Based Approach for DDoS Attack Detection and Prevention in Cloud Computing Environment", *Multimedia Tools and Applications*, Vol. 81, No. 2, pp. 1867-1891, 2022.
- [5] F. Schauer, M. Krbec and M. Ozvoldova, "Controlling Programs for Remote Experiments by Easy Remote ISES (ER-ISES)", *Proceedings of IEEE International Conference on Remote Engineering and Virtual Instrumentation*, pp. 1-8, 2021.
- [6] Random Forest, Available at https://en.wikipedia.org/wiki/Random_forest, Accessed on 2020.
- [7] Tin Kam Ho, "Random Decision Forests", *Proceedings of International Conference on Document Analysis and Recognition*, pp. 278-282, 1995.
- [8] A. Zeinalpour and H.A. Ahmed, "Addressing the Effectiveness of DDoS-Attack Detection Methods Based on the Clustering Method using an Ensemble Method", *Electronics*, Vol. 11, pp. 2736-2745, 2022.
- [9] Canadian Institute for Cybersecurity, Available at <https://www.unb.ca/cic/datasets/ids-2017.html>, Accessed on 2017.