# SYN Flood Attack From TCP/IP Using Statistical Analysis of Machine Learning Techniques

[1]**Vankayalapati Nagaraju,**

*Research Scholar, Department of ECE, VISTAS, India*

[2]**Dr. Arun Raaza, Director,**

*CARD, Department of ECE, VISTAS, India*

[3]**Dr. V.Rajendran,**
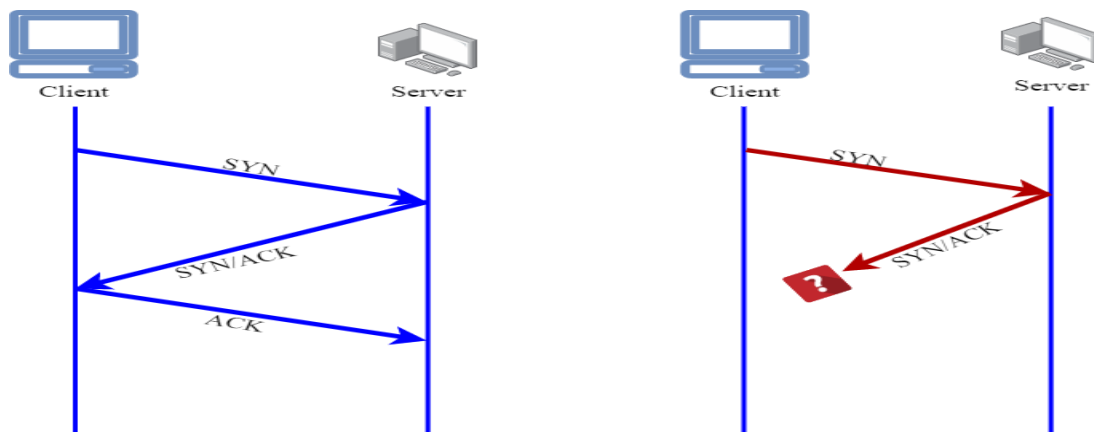
*Director, Department of ECE, VISTAS, India*

*Abstract*

*Most of the network management affected by SYN flood attack needed to secure the server since being harassed by malevolent attackers. The Transmission control protocol synchronize (SYN) flood malicious attack happens whilst the hacker floods the network system with requests in order to overcome the intention and make it unable to respond to new real connection requests. This makes many of the intention server's communications ports into half-open state. If this kind of attack established in big data analysis, Artificial Intelligence, as well as Internet of Things, then SYN flood has to be found. So we proposed novel model for detecting the attack through statistical analysis in machine learning techniques such as decision tree classifier, ensemble gradient boosting classifier, MLP classifier for enhancing the model performance by evaluating the metrics such as accuracy, F-Score, FNR.*

*Keywords: Synchronize (SYN), Machine Learning Techniques, Decision Tree classifier, Ensemble gradient boosting classifier, Multi-Layer Perceptron (MLP) Classifier, Accuracy, False Negative Rate (FNR), F-Score.*
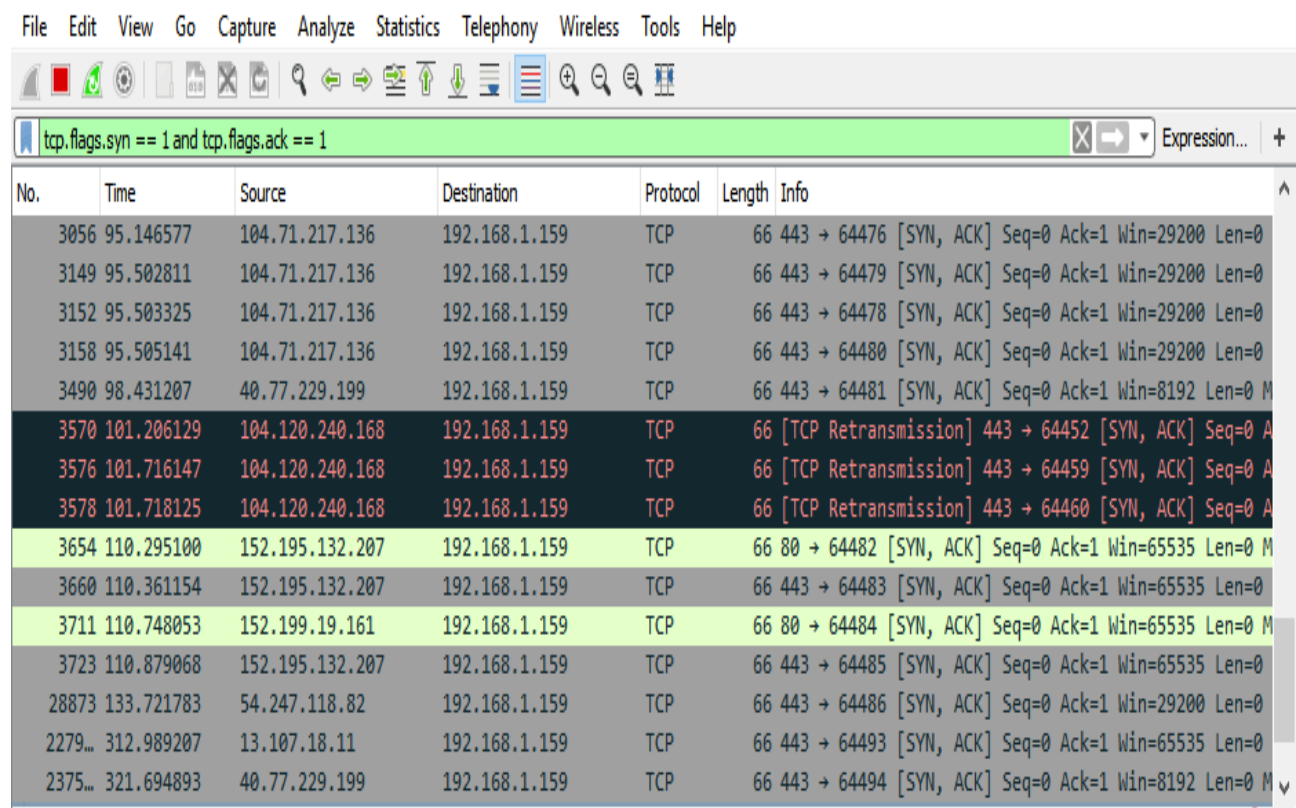
## 1. Introduction

The working principle of SYN flood attack is to undertake the handshaking (client-server communication) process of Transmission control Protocol (TCP). Initially, the client sends the data by SYN packets to the server, then server respond to the same with acknowledgement as SYN/ACK for improving communication between client and server.

This kind of malicious attack is simple to secure through placing a trouble-free firewall rule to obstruct packets with the attacker's source IP address that makes the attack shutdown after any attack found. The way of preventing SYN flood attack as follows:

a. Filtering

b. Increase in backlog

c. Half open TCP

d. Firewalls and Proxies

e. Reducing SYN-Received Timer

f. SYN collection

g. Reprocessing the oldest half-open TCP

[1] Introduced novel method based on Software Defined Network in machine learning algorithm for detecting TCP based SYN flood attack in the network. The final accuracy achieved around 96% deal with the exchange between accuracy and capacity of the device which makes secure for increase and improvement in the performance. Attackers hurriedly send SYN packets exclusive of spoofing their Internet Protocol source address in SYN flood attack.



Here, every packet source IP address varies with destination HTTP port 80, length 120, size of the window as 64 ($2^5$). The protocol used in this attack is Transmission Control Protocol, the value of tcp.flags.syn==tcp.flags.ack==1, and the number of SYN/ACK is moderately very less.
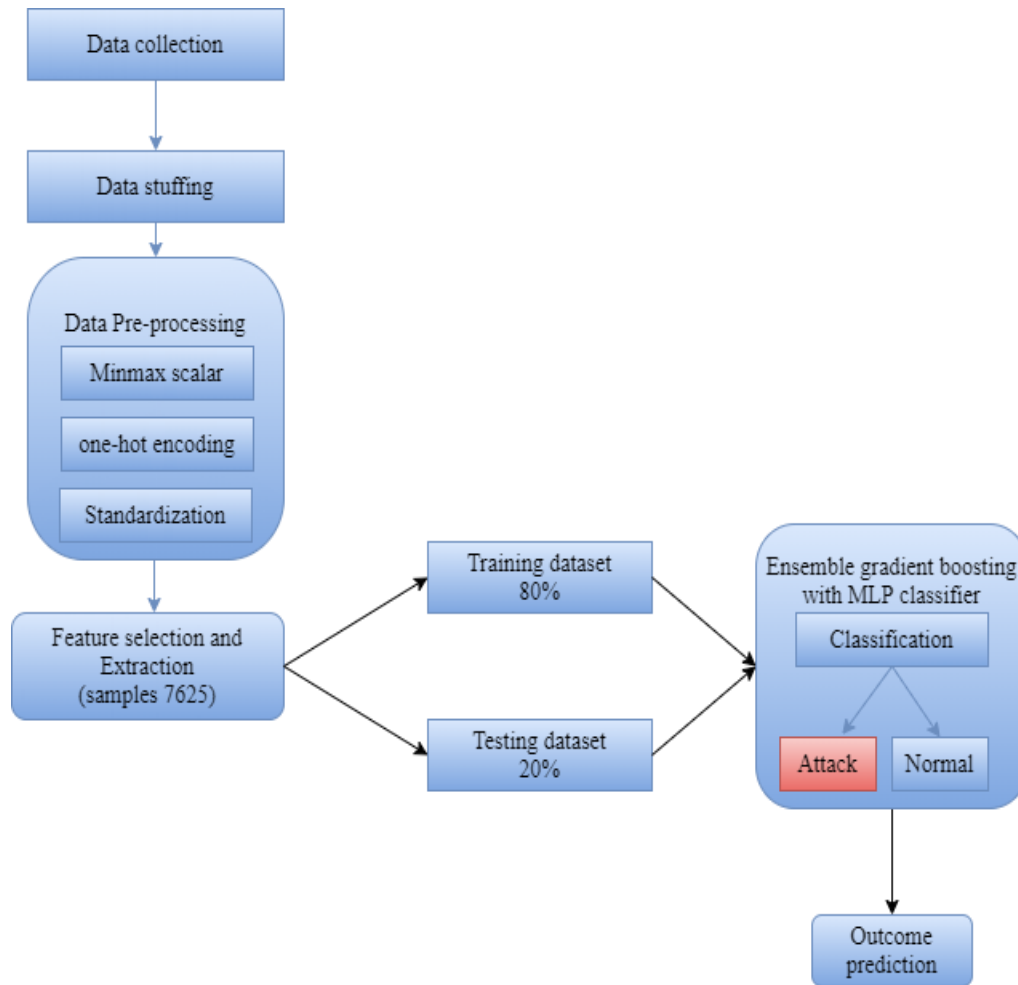
## 2. Background

 [2] Designed methodology for perceiving LDDoS attacks depends upon the individuality malevolent transmission control protocol (TCP) flows which is a kind of DDoS attack. Here, two different arrangements of datasets one produce from simulated network, the other from publically available CIC DoS dataset. The experimental outcome shows the accuracy around 99.99% with remarkably low false rate, low FNR in detection of LDDoS. [3] Detecting DDoS attack could be found by analyzing flow of packets in the network. The novel method introduced telemetry usage from the cloud through machine learning techniques such as KNN and Decision Tree algorithm (CART) to find DDoS attack. [4] Proposed five rule-based machine learning algorithm namely JRip, decision table, PART, OneR, and ZeroR. Among all algorithms, Slow loris achieves high accuracy of 99.7% using PART classifier for realizing the detection of ICMP attack as well as HTTP flood attack. [5] Suggested the improvement of new defense technologies for differentiating the attack from normal using machine learning method based on DDoS detection system. The investigational outcomes reveals that detection rate accuracy as 96% with high precision, low FAR, through sampling rate with 20% of traffic in the network. [6] Introduced synchronize flood attack in cloud based technology and finding attacks extorted from TCP/IP header based on features in the datasets. The final output has been predicted for model performance by identifying and classifying attacks in the network with reference to some dataset features. [8] Imposes literature survey by comparing various researchers regarding detection of finding attacks in the system through development of protocol, enumeration, selection strategy, and synthesizing using machine learning algorithms. [9] Utilized binary classification and two optimizers for detecting the accuracy of attacks in the system via data pre-processing methods, tuning hyper parameters, and Neural Network architectures. [10] [15] Intended to detect the DDoS system based on C.45 decision based random forest algorithm for classifying the network as normal or attack. The novel introduced algorithm attached with signature based detection techniques produces decision tree to achieve automatic, effective detection of signature attacks for flooding attacks. Machine learning algorithm used for validates the model to enhance the better performance. [12] Recommended the test data calculation based on IRE values evaluated during testing phase via monitoring the incoming and outgoing traffic of the web server for detecting the SYN flood attack in the network device. The auto encoder found the attack with delay of 19.2 seconds after scrutinizing the traffic in the network. [13] exposed the detection of SYN flood attack in the network due to IP header through payload and impractical field which makes the detection faster, and highly effective. Finally the alarming section gives the alert if any abnormal behavior occurs in the network. [14] focused on different security attacks in machine learning algorithms and cloud can be used for identifying attacks in the network. Various machine learning techniques like Naïve Bayes, SVM, Logistic Regression, and ensemble methods can be used for finding attack such as authentication attack, Man-in-the-Middle attack, malware attack, DoS attack. [16] Used DARPA, KDD Cup, and CONFICKER datasets for categorizing attack or normal through evaluating metrics such as accuracy, average delay, cost, loss of packets, overhead, packet delivery ratio as well as throughput. The simulation results shows 99% of high accuracy, less false reduction, less overhead, less packet loss and increase in throughput and packet delivery ratio.

## 3. Proposed work

### 3.1 Workflow for proposed model

The flow of work includes organizing and repetitive prototype of any action facilitate by some association resources into processes which converts the resources, services or any information about process.
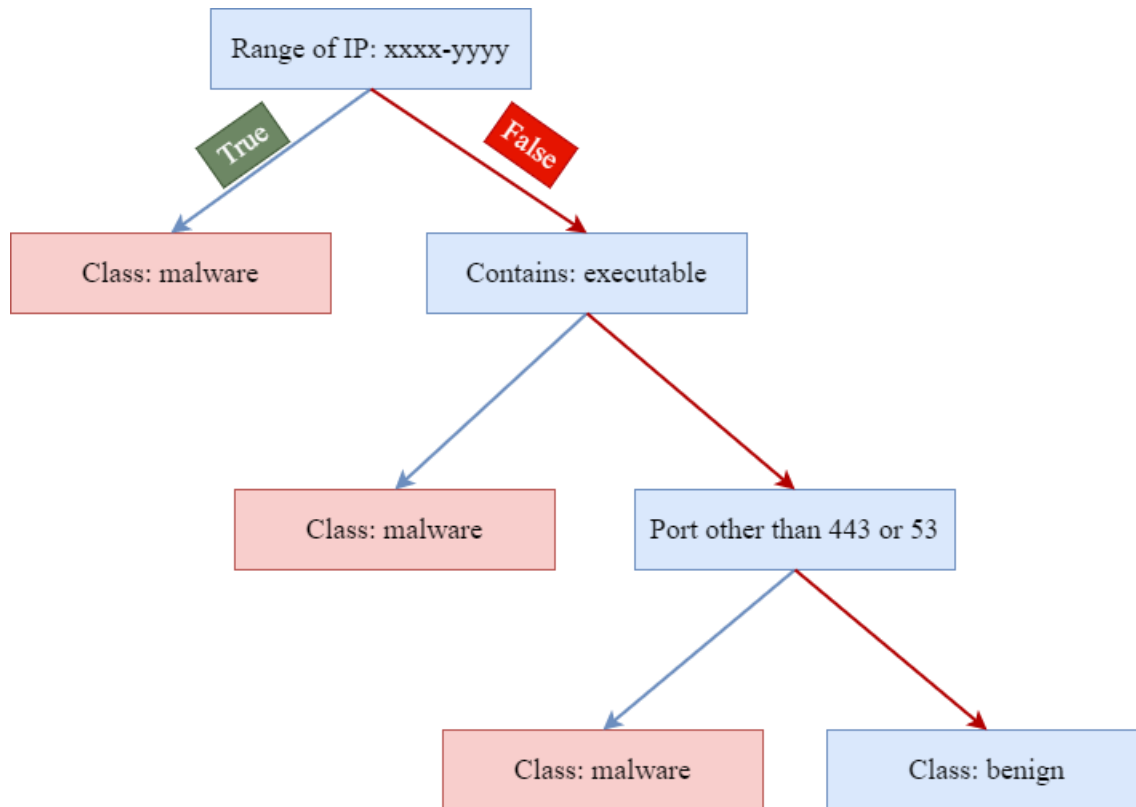
The first step is to collect the data from corresponding repository, then loading the data for further processing leads to data pre-processing in which min-max scalar initialized, then apply it to the features for further transformation, thus scaling applied features can be shown including percentage of attacks. Hence, calculate the feature-wise distribution for every feature with training and testing ratio as 80:20 as well as choosing classifiers such as ensemble gradient boosting along with MLP classifier for categorizing the network as malicious or normal. Thus the overall performance model can be predicted through the final outcome prediction via classifier algorithms in machine learning techniques.

### 3.2 Proposed work algorithm

In the proposed method, finding SYN flood attack through machine learning techniques such as decision tree classifier, ensemble gradient boosting classifier, Multi Layer Perceptron classifier.

*A. Decision tree classifier:* Decision tree classifier comes under supervised machine learning for categorizing the SYN flood based on Transmission Control Protocol into either malware or benign. The figure shows that if IP range specified as xxxx-yyyy, then the class found to be malware otherwise IP address is executable. If IP address is executable, then class is specified as malware or else referring some port address. If it is mentioned as port then the class is malware or found the class to be benign.
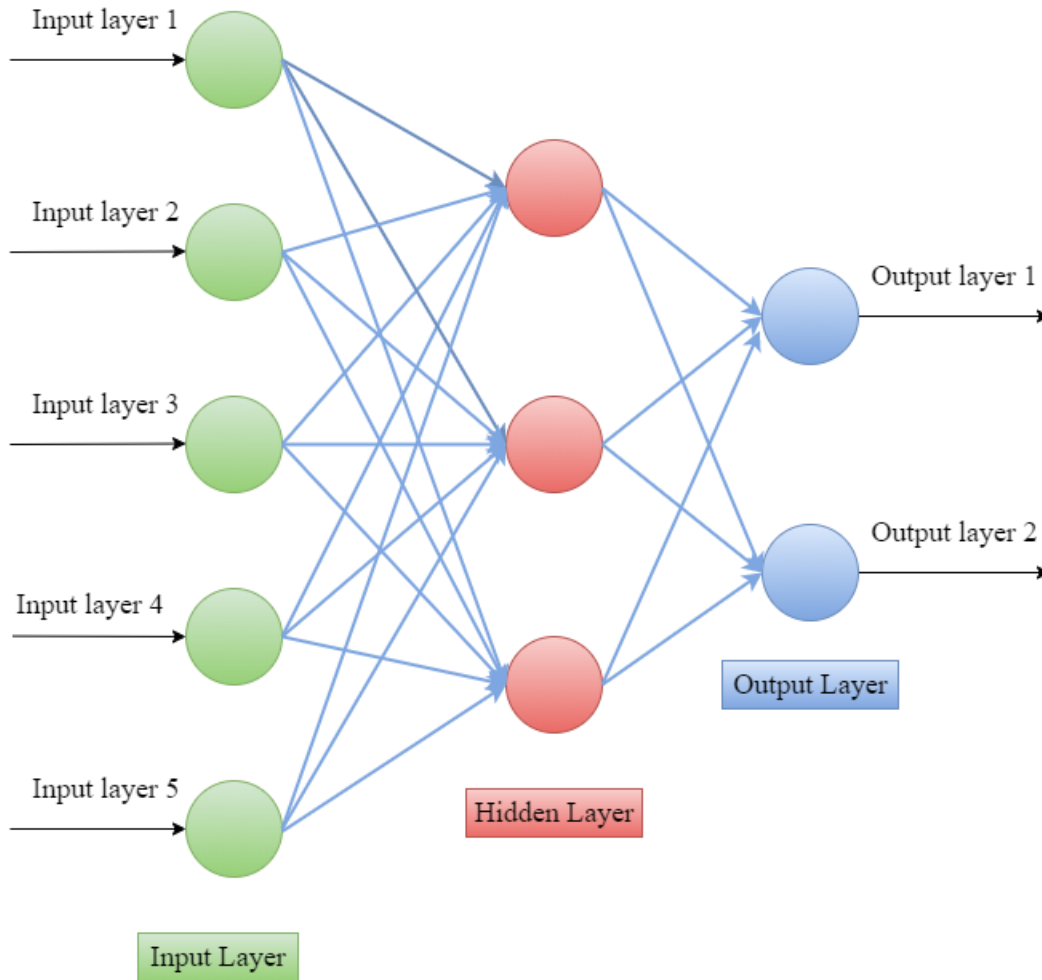
*B. Ensemble Gradient Boosting Classifier*: An ensemble gradient boosting classifier is one of the machine learning methods which can be utilized in case of both classification as well as regression problems that generates a prophecy model in the form of an ensemble of weak model normally decision trees. [11] Proposed support vector machine for categorizing the datasets as normal or attack in the network processed under four phases namely attack simulation, data collection, feature selection and classification. Finally, the process generates the classification accuracy as 100% and True Positive Rate as 100% reveals dazzling performance which supply as valuable asset in the field of network security. In our proposed model, the detection accuracy achieved as 100% which outstands the performance of the model through ensemble gradient boosting algorithm along with decision tree.

*C. Multi-Layer Perceptron:* Mainly MLP relies on intrinsic neural network to complete the undertaking of classification problems. Importing and initializing MLP classifier using Python code.

```
#importing MLP classifier
From sklearn. neural network import MLPClassifier
# initializing the MLPClassifier
Clf_C= MLPClassifier (solver='lbfgs', alpha=1e-5, hidden_layer_sizes = (5, 2), random_state=1)
```

Now, the solver = lbfgs or gradient descent parameter in neural network optimizes the log loss function, alpha= 1e-5 parameter for regularization that assist in evade over fitting by chastising weights with huge scale, hidden layer size parameter permit the number of nodes and layers in neural network classifier, random_state parameter allocate to set seed for generating the output.

The above diagram shows the input layer as 5 and the output layer as 2 with hidden layers (i.e)
hidden_layers_size= (5, 2)

## 4. Dataset description

The dataset is collected from Kaggle resource datasets which has 7167 samples with 82 features that
performs splitting phase which has 5733 samples for training sets with 80 percent and the remaining 1434
samples allotted for testing sets with 20 percent as well. Herein, the number of normal records is 2166,
the number of attacks refer 5001 attacks, the percentage of attacks in SYN flood found to be 69.78%. This
kind of attack has to be found via statistical analysis of some techniques in machine learning.

| Total number of records | Number of attacks | Number of normal records | Percentage of attacks |
|---|---|---|---|
| 7167 | 5001 | 2166 | 69.78% |

Some of the features namely source IP, source port, destination port, destination IP, protocol, timestamp,
flow duration, total forward packets, total backward packets, active std, active min max, idle mean, idle
std, similar HTTP, inbound, and type of the attack. The features used in detecting SYN flood attack with
its description are summarized.

1749

| Features | Description |
|---|---|
| Source IP | The IP address of the device sending IP packet |
| Source port | The subsequent offered data assigned to the user by TCP/IP |
| Destination IP | The machine to which the packet is being used |
| Destination port | Well-known ports (Server application) |
| Protocol | Data transmission between dissimilar devices in similar net |
| Timestamp | Current time of an incident saved by a system |
| Flow duration | Time duration for flow of packets |
| Total forward packets | Number of received packets from interface which is linked to a particular system |

| Metrics | Description | Formula |
|---|---|---|

| Total backward packets | Reverse of forward packets |
|---|---|
| Active standard | Node performs certain operation in the network |
| Active minimum | Minimum flow of packets |
| Active Maximum | Maximum flow of packets |
| Idle mean | Finding mean value |
| Idle Standard | Finding standard value in the network |
| HTTP | Hypertext Transfer Protocol |
| Inbound | Inbound firewall defend the system beside referral traffic from internet |
| Attack type | Finding type of the attack |

4.1 Metrics Evaluation in proposed model

The overall model performance can be evaluated through metrics such as accuracy, prediction time, F-Score, False Negative Rate (FNR), for both training and testing samples. The number of incorrect and correct predictions are summarized with count values and broken down by every class. Confusion matrix frequently used to illustrate classification model performance.

|  | Class 1 Predicted value | Class 2 Actual value |
|---|---|---|
| Class 1 Actual value | TP | FN |
| Class 2 Predicted value | FP | TN |

| Accuracy | Correctly classified percentage | $Accuracy = \dfrac{TP + TN}{TP + TN + FP + FN}$ |
|---|---|---|
| F-Score | Harmonic mean of precision and recall | $F-score = \dfrac{2 * recall * precision}{recall + precision}$ |
| FNR | Miss rate | $FNR = \dfrac{FN}{FN + TP}$ |
| Precision | Measure of result relevancy | $Precision = \dfrac{TP}{TP + FP}$ |
| Recall | Measure of success prediction when classes are imbalanced. | $Recall = \dfrac{TP}{TP + FN}$ |
| Prediction time | Time to predict the future value | - |

Based on the metrics evaluation such as accuracy, F-score, FNR, Precision, Recall, prediction time, the overall model performance can be calculated through statistical analysis of metrics.
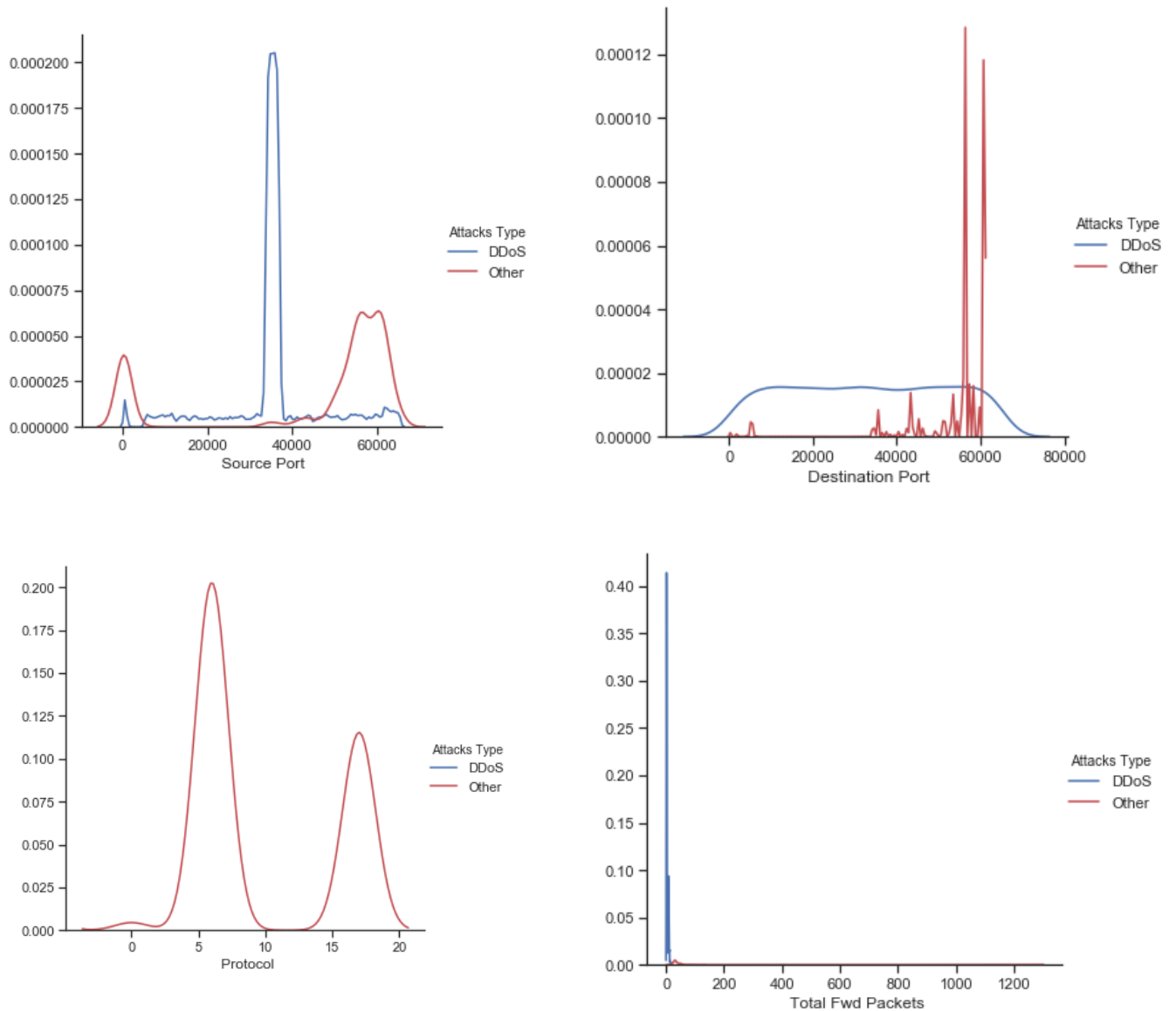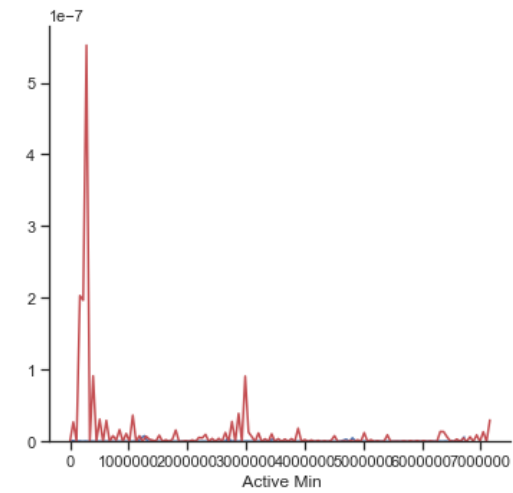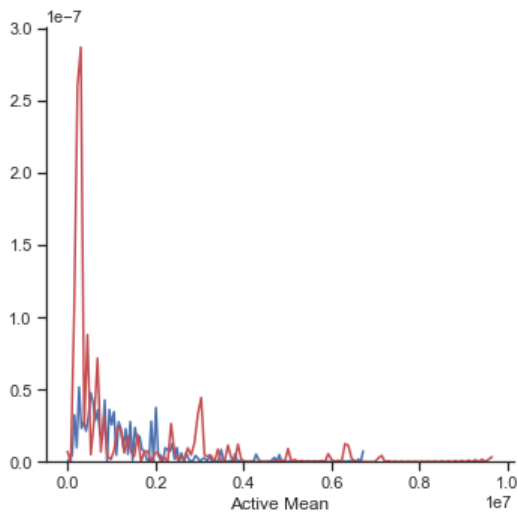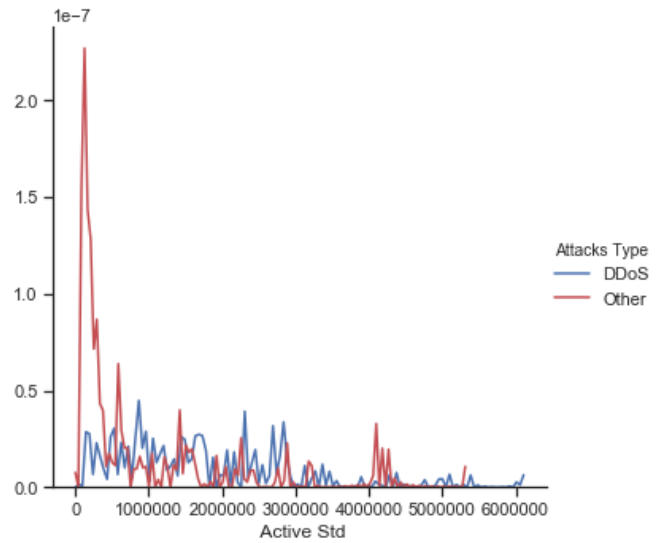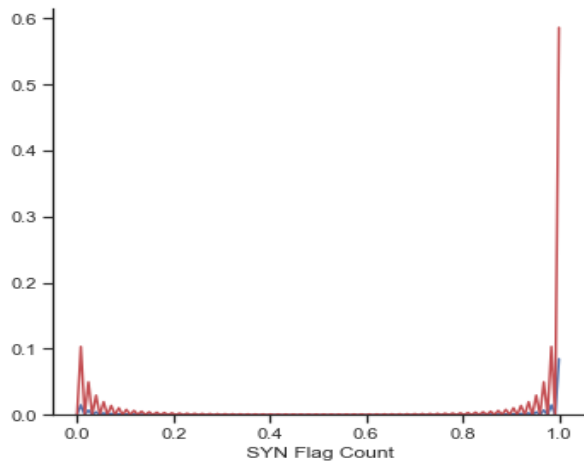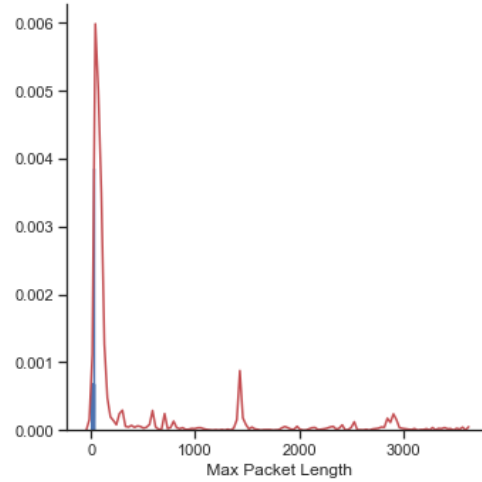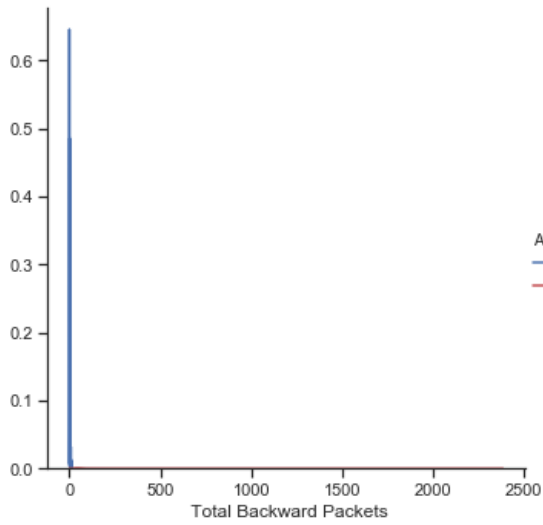
### 4.2 Feature calculation for mean, median, mode

The below table reveals that finding mean, median, mode for both attack or normal mentioned some features of synchronize attack based on Transmission control protocol.
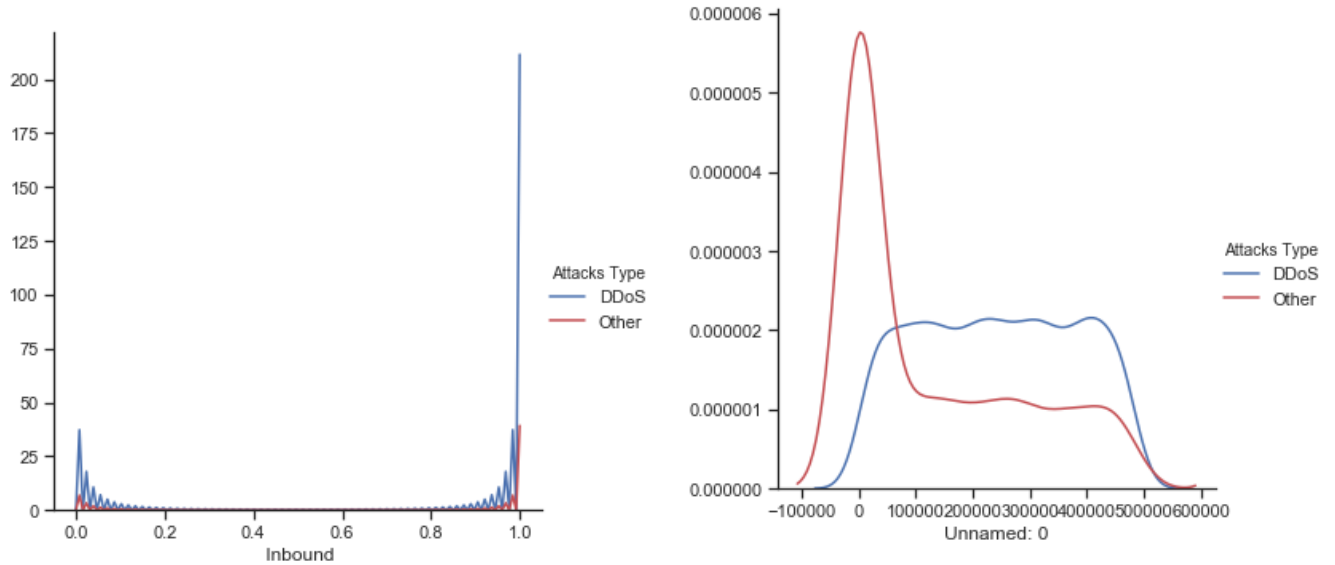
| Features | Mode attack | Median attack | Mean attack | Mode normal | Median normal | Mean attack |
|---|---|---|---|---|---|---|
| Source port | 35363.5 | 35461.000000 | 3.545055e+04 | 443.0 | 56221.500000 | 4.566191e+04 |
| Destination port | 19951.0 | 32793.000000 | 3.305446e+04 | 53.0 | 443.000000 | 1.017277e+04 |
| Protocol | 6.0 | 6.000000 | 6.000000e+00 | 6.0 | 6.000000 | 9.847645e+00 |
| Total Fwd Packets | 2.0 | 2.000000 | 4.188562e+00 | 2.0 | 2.000000 | 9.573869e+00 |
| Total Backward Packets | 0.0 | 2.000000 | 1.708458e+00 | 2.0 | 2.000000 | 1.040859e+01 |
| Total Length of Fwd Packets | 12.0 | 12.000000 | 2.514897e+01 | 6.0 | 64.000000 | 8.994307e+02 |
| Total Length of Bwd Packets | 0.0 | 12.000000 | 1.024115e+01 | 0.0 | 82.000000 | 8.007978e+03 |
| Fwd Packet Length Max | 6.0 | 6.000000 | 6.008798e+00 | 6.0 | 32.000000 | 1.426385e+02 |
| Fwd Packet Length Min | 6.0 | 6.000000 | 6.004399e+00 | 6.0 | 6.000000 | 1.687073e+01 |
| Fwd Packet Length Mean | 6.0 | 6.000000 | 6.005279e+00 | 6.0 | 31.000000 | 3.884630e+01 |
| Fwd Packet Length Std | 0.0 | 0.000000 | 1.854832e-03 | 0.0 | 0.000000 | 3.865311e+01 |
| Bwd Packet Length Max | 6.0 | 6.000000 | 3.212158e+00 | 0.0 | 41.000000 | 2.445300e+02 |
| Bwd Packet Length Min | 6.0 | 6.000000 | 3.212158e+00 | 0.0 | 6.000000 | 3.311450e+01 |
| Bwd Packet Length Mean | 6.0 | 6.000000 | 3.212158e+00 | 0.0 | 32.666667 | 8.485964e+01 |
| Flow IAT Mean | 1.0 | 35.333333 | 1.472787e+06 | 1.0 | 6963.166666 | 3.987781e+05 |

## 5. Statistical Analysis

Suppose if we wish to concern in some other techniques such as big data analytics, Artificial Intelligence, Internet of Things which causes major issues to the network of recognize attacks. Thus identifying and classifying the network as malicious or standard is necessary in recent technology trends. So, the novel proposed method expose to do favor to the network by implementing statistical analysis of such classifiers in machine learning models for categorizing malicious or normal. [7] Anticipated machine learning methods based on cloud which is on source side for detecting the attack namely Denial of Service in the network. This paper illustrates the model performance through statistical analysis from user as a virtual machine and cloud server to check both inward and outward packets in the network. Hence, the accuracy achieved as 99.7% in detecting DoS attack via experimental results.

## 6. Results and Discussion

### 6.1 Algorithm evaluation with full feature set

| | 1% training samples | | | 10% training samples | | | 100% training samples | | |
|---|---|---|---|---|---|---|---|---|---|
| | GB | DT | MLP | GB | DT | MLP | GB | DT | MLP |
| Train_time (s) | 0.031404 | 0.003999 | 0.015275 | 0.140637 | 0.005999 | 0.028002 | 1.264150 | 0.056145 | 0.102025 |
| Prediction_time | 0.000000 | 0.006002 | 0.007003 | 0.000000 | 0.005007 | 0.007006 | 0.008994 | 0.000000 | 0.004998 |
| Accuracy_train | 0.976667 | 0.973333 | 0.596667 | 1.000000 | 1.000000 | 0.590000 | 1.000000 | 1.000000 | 0.590000 |
| Accuracy_test | 0.987448 | 0.982566 | 0.598326 | 0.996513 | 0.997211 | 0.597629 | 1.000000 | 1.000000 | 0.597629 |
| F_train | 0.983683 | 0.981395 | 0.699752 | 1.000000 | 1.000000 | 0.693267 | 1.000000 | 1.000000 | 0.693267 |
| F_test | 0.991045 | 0.987605 | 0.685590 | 0.997484 | 0.997996 | 0.684182 | 1.000000 | 1.000000 | 0.684182 |
| FNR_train | 0.000000 | 0.000000 | 0.331754 | 0.000000 | 0.000000 | 0.341232 | 0.000000 | 0.000000 | 0.341232 |
| FNR_test | 0.000000 | 0.000000 | 0.369478 | 0.005020 | 0.000000 | 0.372490 | 0.000000 | 0.000000 | 0.372490 |

## 6.2 Algorithm evaluation with statistical feature set

| | 1% training samples | | | 10% training samples | | | 100% training samples | | |
|---|---|---|---|---|---|---|---|---|---|
| | GB | DT | MLP | GB | DT | MLP | GB | DT | MLP |
| Train_time (s) | 0.059633 | 0.009479 | 0.011116 | 0.152682 | 0.152682 | 0.018001 | 1.189550 | 0.052131 | 0.049473 |
| Prediction_time | 0.011002 | 0.004001 | 0.005886 | 0.015635 | 0.015635 | 0.007006 | 0.015626 | 0.004996 | 0.002510 |
| Accuracy_train | 0.976667 | 0.970000 | 0.296667 | 1.000000 | 1.000000 | 0.296667 | 1.000000 | 1.000000 | 0.296667 |
| Accuracy_test | 0.987448 | 0.978382 | 0.305439 | 0.993724 | 0.993724 | 0.305439 | 1.000000 | 1.000000 | 0.305439 |
| F_train | 0.983683 | 0.979118 | 0.000000 | 1.000000 | 1.000000 | 0.000000 | 1.000000 | 1.000000 | 0.000000 |
| F_test | 0.991045 | 0.984676 | 0.000000 | 0.995461 | 0.995461 | 0.000000 | 1.000000 | 1.000000 | 0.000000 |
| FNR_train | 0.000000 | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 1.000000 | 0.000000 | 0.000000 | 1.000000 |
| FNR_test | 0.000000 | 0.000000 | 1.000000 | 0.009036 | 0.009036 | 1.000000 | 0.000000 | 0.000000 | 1.000000 |

The table specifies the algorithm evaluation for full feature set as well as statistical feature set with 1%, 10%, 100% training samples for all classifiers in the sequence of 57, 573, 5733 samples on training datasets.

### 6.3 Feature importance determined by t-test

The t-test score value can be calculated for determining the importance of features in synchronizes dataset for classifying the attack as break up.

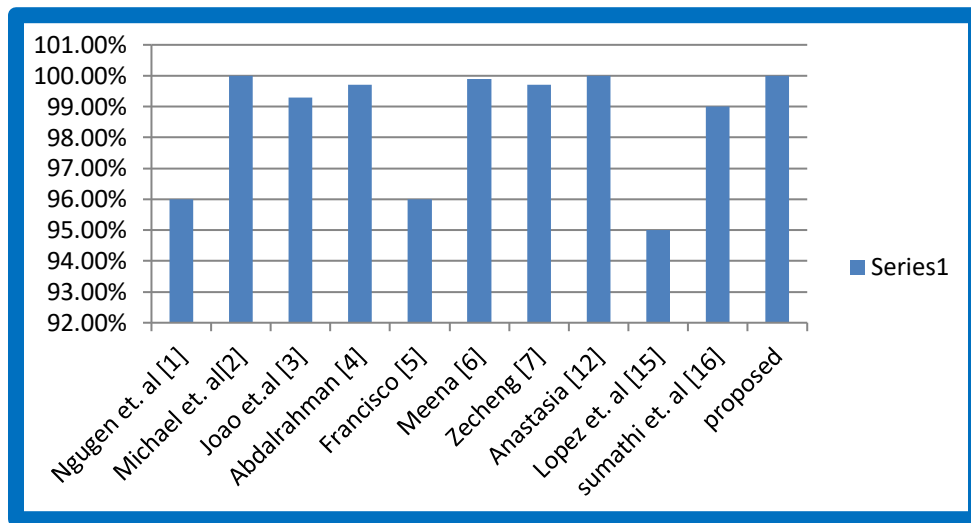| | ACK Flag Count | Inbound | URG Flag Count | Protocol | Destination Port | Bwd Packet Length Min | CWE Flag Count | Average Packet Size | Avg Fwd Segment Size | Fwd Packet Length Mean | Subflow Bwd Bytes | Idle Min | Active Mean | Flow IAT Min | Fwd IAT Min | SYN Flag Count | Active Std | Bwd IAT Std |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| t-Score (squared) | 22960.54668 | 21861.03077 | 4078.03611 | 2558.23402 | 2035.87731 | 1638.49086 | 1325.82968 | 1268.59244 | 1209.33493 | 1209.33493 | 26.79566 | 25.63777 | 23.95676 | 17.27751 | 10.06518 | 7.62092 | 7.11081 | 4.77502 |
| P-Value | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000003 | 0.000152 | 0.000578 | 0.000768 | 0.028891 |

From analyzing the algorithm evaluation statistically with three different classifiers, the outcomes could be predicted for finding the model performance. In this proposed model, among three classifiers in MLT, ensemble gradient boosting and decision tree classifier shows better accuracy as 100 %. Hence, by referring the percentage of accuracy in validation stage this novel proposed model is best for categorizing the net as attack or normal in SYN flood attack based upon TCP.

## 6.4 Comparative Analysis

Comparative Analysis for accuracy detection determined by various authors is shown below:

| Authors | Techniques used | Detection of attack | Accuracy detection |
|---|---|---|---|
| [1] | MLP | TCP-SYN | 96% |
| [2] | DT, KNN | LDDoS | 99.99% |
| [3] | KNN | DoS | 99.3% |
| | CART | | 91.07% |
| [4] | PART | HTTP flood attack | 99.7% |
| [5] | ML | DoS | 96% |
| [6] | NN | TCP/ IP syn flood | 99.9% |
| | NB | | 99.1% |
| | KNN | | 98.2% |
| | DT | | 99.9% |

| | | | |
|---|---|---|---|
| **[7]** | Cloud based ML | DoS | 99.7% |
| **[8]** | ML and DL | DDoS | |
| **[9]** | BNN | DDoS | |
| | LSTM RNN | | |
| **[10]** | C4.5 Decision Tree | DDoS | |
| **[11]** | SVM | DoS | 100% |
| **[12]** | Auto encoder | DoS, Code injection | |
| **[13]** | Packet filtering algorithm | TCP SYN flood | |
| **[14]** | SVM, NB,DT, ensemble methods | | ✓ |
| **[15]** | Decision Tree based RF | DDoS | |
| | KNN | | 95% |
| **[16]** | Deep Learning NN classifier | DoS | 99% |
| **Proposed** | Machine Learning classifier | SYN flood attack | 100% |



The above chart shows the classification accuracy detection for all authors classifies attack or normal with 100% accuracy using decision tree classifier. Thus our proposed work shows better outcome among various authors.

## 7. Conclusion

In this novel proposed method, the SYN flood attack could be recognized and categorizing the datasets as attack or normal using statistical analysis in machine learning techniques like Decision Tree, ensemble gradient boosting classifier, as well as Multi-Layer Perceptron classifier by extracting features acquired from TCP/IP header protocol. Hence, the efficiency of IDS in distinguishing SYN flood attack depends on extracting specific features using MLT. The investigational outcomes demonstrate that ensemble gradient boosting and decision tree algorithm generates 100% accuracy with least prediction time as 0.01 seconds in detection of SYN flood attack. Based on accuracy prediction evaluated through various classifiers, the performance of the model can be acknowledged. The future scope is to assess the feature extraction implemented in machine learning techniques that concerned it in real world application through indispensable alteration which may improve the performance of the model.

**References**

[1] Nguyen Ngoc Tuan, Pham Huy Hung, Nguyen Danh Nghia, Nguyen Van Tho, Trung V. Phan, Nguyen Huu Thanh "A Robust TCP-SYN Flood Mitigation Scheme Using Machine Learning Based on SDN", IEEE, ICTC 2019.

[2] Michael Siracusano, Stavros Shiaeles, Bogdan Ghita "Detection of LDDoS Attacks Based on TCP Connection Parameters", IEEE,

[3] João Henrique Corrêa, Patrick Marques Ciarelli, Moises R. N. Ribeiro, Rodolfo da Silva Villaca "On Machine Learning DoS Attack Identification from Cloud Computing Telemetry", LANCOMM'19, 2019.

[4] Abdalrahman Hwoij, Mouhammd Al-kasassbeh, Mustafa Al-Fayoumi "Detecting Network Anomalies using Rule-based machine learning within SNMP-MIB dataset".

[5] Francisco Sales de Lima Filho, Frederico A. F. Silveira, Agostinho de Medeiros Brito Junior, Genoveva Vargas-Solar, and Luiz F. Silveira "Smart Detection: An Online Approach for DoS/DDoS Attack Detection Using Machine Learning", Security and Communication Networks, 2019.

[6] Muna Sulieman AL-Hawawreh, "SYN Flood Attack Detection in Cloud Environment Based on TCP/IP Header Statistical Features", International Conference of Information Technology (ICIT), 2017.

[7] Zecheng He, Tianwei Zhang, Ruby B. Lee, "Machine Learning Based DDoS Attack Detection from Source Side in Cloud", IEEE 4th International Conference on Cyber Security and Cloud Computing, 2017.

[8] Saritha, B. RamaSubba Reddy, A Suresh Babu "Prediction of DDoS Attacksusing Machine Learning and Deep Learning Algorithms", International Journal of Recent Technology and Engineering (IJRTE), 2019.

[9] Meejoung Kim "Supervised learning-based DDoS attacks detection: Tuning hyper parameters", ETRI Journal WILEY, 2019.

[10] Marwane Zekri, Said El Kafhali, Noureddine Aboutabit and Youssef Saadi "DDoS Attack Detection using Machine Learning Techniques in Cloud Computing Environments" IEEE, 2018.

[11] Zerina Mašetić, Dino Kečo, Nejdet Doğru, Kemal Hajdarević "SYN Flood Attack Detection in Cloud Computing using Support Vector Machine", TEM Journal, 2018.

[12] Anastasia Gurina, and Vladimir Eliseev "Anomaly-Based Method for Detecting Multiple Classes of Network Attacks", MDPI, Information 2019.

[13] Haris, S.H.C.; Ahmad, R.B.; Ghani, M.A.H.A. Detecting TCP SYN flood attack based on anomaly detection. In Proceedings of the 2010 Second International Conference on Network Applications, Protocols and Services, Kedah, Malaysia, 22–23 September 2010; pp. 240–244.

[14] Dhivya R, Dharshana R, Divya V "Security Attacks Detection in Cloud using Machine Learning Algorithms", International Research Journal of Engineering and Technology, 2019.

[15] Lopez, Alma D.; Mohan, Asha P.; and Nair, Sukumaran (2019) "Network Traffic Behavioral Analytics for Detection of DDoS Attacks," SMU Data Science Review: Vol. 2 : No. 1 , Article 14.

[16] S. Sumathi, N. Karthikeyan "Detection of distributed denial of service using deep learning neural network" Journal of Ambient Intelligence and Humanized Computing, 2020.