

PROJECT REPORT
OF
DATA MINING AND DATA WAREHOUSING
(CSPC-308)



Submitted To:
Dr. Nonita Sharma
CSE Department

**DR. B.R. AMBEDKAR NATIONAL INSTITUTE OF
TECHNOLOGY JALANDHAR**

TEAM INFORMATION

<u>SNO</u>	<u>NAME</u>	<u>ROLL NO</u>	<u>CONTACT</u>
1.	RAMAVATH SAIKUMAR	18103077	ramavaths.cs.18@nitj.ac.in
2.	ROHIT MITTAL	18103081	rohitm.cs.18@nitj.ac.in
3.	SHOBHIT TEWARI	18103086	shobhitt.cs.18@nitj.ac.in
4.	VIKAS	18103094	vikas.cs.18@nitj.ac.in

DETECTION AND CLASSIFICATION **OF** **DDoS ATTACKS**

Abstract

In recent years, with increase in internet technology and users, the advanced threats against Cyber–Physical Systems (CPSs), attacks, are also increasing. And it can be observed that DDoS attacks are one of the most hazardous attacks and contributes a big part overall network attacks. These attacks are highly dangerous because it is very difficult for the networks to distinguish between normal and malicious flows. The testing and mitigation of DDoS attacks is very difficult because of many reasons like rigidity, complexity, cost and architecture of equipment's and protocols used are different for different vendors.

A lot of work has been done and is being done to detect and classify the DDoS attacks by using Machine Learning techniques but still the question is open for the best algorithm for this task.. This paper is motivated by the question: “Which supervised learning algorithm will give the best outcomes to detect DDoS attacks”? We used various ML techniques (Decision Tree, Random Forest, Naïve Bayes, K Nearest Neighbor and Gradient Boost) on the dataset and are able to get the maximum accuracy of 98.84%.

Keywords: Classification, Supervised Learning, DDoS detection, DDoS attack, security, Network Threats, KNN, Gradient Boost, Random Forest Classifier, Naïve Bayes, TCP SYN, UDP, MSSQL

TABLE OF CONTENTS

S No.	Contents	Page No.
1	Introduction	6
2	DDoS Attacks	
	2.1 Overview	7
	2.2 Types of DDoS attacks	7-8
	2.3 Common DDoS attacks	8-9
3	Methods	
	3.1 Decision Tree	10-11
	3.2 Naïve Bayes	11-12
	3.3 K Nearest Neighbor	12-13
	3.4 Random Forest	13-14
	3.5 Gradient Boost	14-15
4	Experimentation	
	4.1 Overview	16
	4.2 Data Preparation	17-20
	4.3 Training	21-26
	4.4 Testing	26
5	Results and Discussions	27-29
6	Conclusions and Future Work	
	6.1 Conclusions	30
	6.2 Future Work	30
7	APPENDIX	31

1. INTRODUCTION

In today's world all the important organizations such as defense, research, government organizations and nuclear plants depends highly on computer networks. Distributed Denial of Service (DDoS) attack is a very big problem to network security which aims at wear out the target networks with malicious and waste traffic. DDoS attack the attacker uses the Botnet-based Computers which are infected by Malware and fully controlled by the attacker to not only choke the system resources but also available bandwidth and hence causing the networks to halt for some time.

This can cause loss of billions to service providing companies (It has been seen that a DDoS attack can cause a loss in millions of dollars every hour for any service providing organizations along with loss in trust of its customers) and even can cause country's defense huge threats and problems. The most dangerous thing about DDoS attacks is that these can't be detected with the traditional strategies such as Signature-Based Intrusion as here the attacker uses the botnet computers to attack.

These reasons have forced us to adopt such strategies which can deal with such situations and losses. The testing and mitigation of DDoS attacks is very difficult because of many reasons like rigidity, complexity, cost and architecture of equipment's and protocols used are different for different vendors. Hence application of Machine Learning comes into play which is the aim of our paper.

In this paper we tried to analyze different machine learning approaches namely Decision Tree, Random Forest, Gradient Boost, Naïve Bayes, K-Nearest Neighbor to find which approach works the best in detecting and classifying the different DDoS attacks.

The dataset used for this paper is [CICDDoS2019](#) which contains benign and the most up-to-date common DDoS attacks. The dataset mainly contains 7 types of DDoS attacks (SYN Flood, UDP Flood, UDP Lag, MSSQL, NetBIOS, LDAP, and PortMap) and benign i.e. normal flow. The dataset is collected by **Canadian Institute for Cyber Security** for studying the different attacks in 2019. Here in features extraction process from the raw data, CICFlowMeter-V3 has been used with the help of which more than 80 traffic features are extracted.

This paper is basically divided into six sections. The next section presents a small review on DDoS attacks and its types. In the third section, under Methods, different machine learning algorithms are explained. The fourth section outlines the details of the design and implementation of the various algorithms which is followed by the comparison and discussions on the results of various algorithms. After that conclusions are made and also future scopes are discussed.

2. DDoS Attacks:

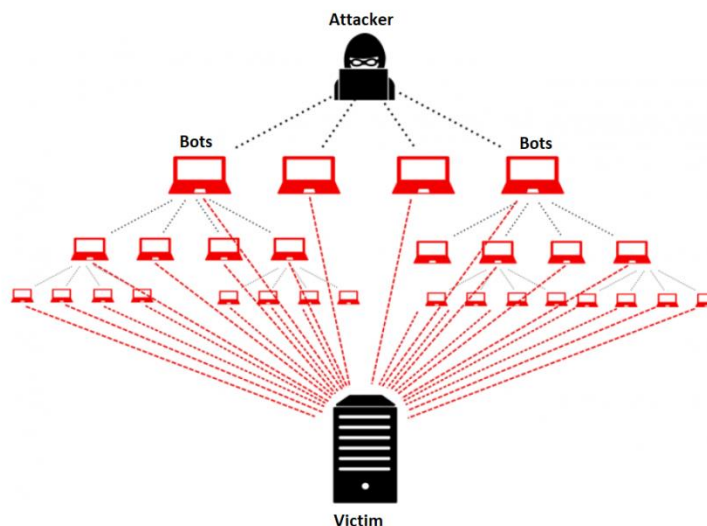
2.1 Overview:

For further discussions we must know what DDoS attacks are and what their main types are. DDoS attack refers to Distributed Denial of Service attack. It is one of the most commonly used Denial of Service attack as it is very difficult to classify.

It is essentially an effort to make a service provider's web service inaccessible to customers, normally by temporarily interrupting or halting the hosting server's services. The attacker sends the data in such a high volume that either the bandwidth completely consumes or the victims' resources get exhausted.

The attacker uses a botnet or a network of compromised computers to carry out this attack. They are referred to as a "botnet." These are used to saturate specific websites, servers, and networks with data that they can't handle.

The botnets can submit more link requests than a server can manage, or they may send massive quantities of data that surpass the targeted victim's bandwidth capabilities. Botnets may have tens of thousands to millions of computers under cybercriminals' influence. There are high chances that the owner of system acting as bot for the attacker doesn't have any idea about this.



2.2 Types of DDoS Attacks:

The DDoS attacks can mainly be classified as:

2.2.1 Volume Based Attacks:

In this type of attack the attacker attempts to saturate the bandwidth of the victim network, rendering the facilities inaccessible to the victims' customers. It includes UDP and ICMP floods, among other items.

2.2.2 Protocol Attacks

In this form of attack, the attacker tries to exploit a flaw in the authentication protocol used by the victim computer, such as TCP's three-way handshake. Other examples of this attack are Floods, Ping of Death, and other SYN events.

2.2.3 Application Layer Attacks

This form of attack focuses primarily on web applications and is considered the most complex and dangerous. The attacker's goal is to bring the web server down. for example GET/POST floods.

2.3 Common DDoS attacks:

The most common attacks are as follows:

2.3.1 SYN Flood

It takes advantage of flaws in the TCP three-way handshake link chain. To begin the "handshake," the host machine receives a synchronized (SYN) message. And also allocates the resources for the request and acknowledges the message by sending an acknowledgement (ACK) flag to the initial request side, but the third packet is never sent and hence the connection is never closed which leads to complete exhaust of resources and which then closes the connection. However, in a SYN flood, spoofed messages are sent and the link is not closed, effectively shutting down operation.

2.3.2 UDP Flood

A UDP flood uses UDP packets to flood random ports on a device or network. The host looks for an application that is listening on those ports, but none is found. The attacker keeps sending the packets and hence exhausting the bandwidth and hence making the services unavailable for the victim's customers.

2.3.3 UDP Lag Attack:

A UDP Lag attack is also volumetric distributed denial-of-service (DDoS) attack using the User Datagram Protocol (UDP) whose main aim is to slow down the services of the server. It is mainly used for lagging the live streaming apps and games so that users experience goes down and hence victim server suffers huge losses.

2.3.4 PortMap Attack:

The PortMap service guides the client to the correct port number so that it can connect with the RPC service that was requested. As with many other UDP-based services (DNS, NTP), it is being used by attackers to conceal the source of the attack and increase its scale. The attacker sends a huge amount of requests to PortMap and hence the real client has to wait an infinite time to get the service of PortMap and hence leading to crash of servers.

2.3.5 LDAP Attack:

Web-based applications that create LDAP statements based on user feedback are vulnerable to LDAP Injection. When an application fails to properly sanitise user input, a local proxy may be used to alter LDAP statements. As a consequence, arbitrary commands, such as granting permissions to unauthorised queries and content alteration within the LDAP tree, may be executed.

2.3.6 MSSQL Attack:

When connecting to a database server or cluster of several database instances, MC-SQLR allows clients to define the database instance with which they are attempting to interact. The SQL Resolution Protocol can be used if a client requires information on network-configured MS SQL servers. A list of instances is returned to the client by the server.

2.3.7 NetBIOS Attack:

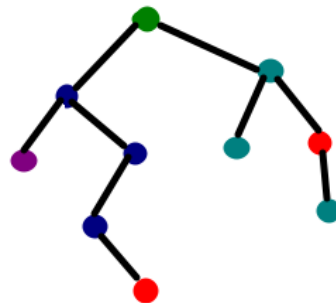
NetBIOS' main objective is to enable applications on different computers to communicate and create sessions in order to access shared resources and locate each other over a local area network. The intruder sends a large number of requests to NetBIOS, causing actual users to wait longer and the server's services to halt.

3. Methods:

There are various supervised machine learning algorithms that can be used for classification of various DDoS attacks. **Some of these are explained below:**

3.1 Decision Tree

The Decision Tree algorithm is a supervised learning algorithm that can be used for both classification and regression. Its main objective is to build a training model (a set of rules) from training data and predicts the target variable for given conditions using this model. It uses a tree structure to build classification or regression models. It incrementally breaks down a data set into smaller and smaller subsets based on attribute values, while also developing an associated decision tree. A tree with decision nodes and leaf nodes is the end product. A classification or decision is represented by a leaf node. So the important point is on what basis these subsets are made from the dataset.



There are mainly two types of criteria which are information gain and gini index which are explained as below:

3.1.1 Information Gain:

The amount of information obtained from a feature about the class is referred to as information gain. To put it another way, it calculates the entropy reduction caused by transforming a dataset based on a function. It can be determined by subtracting the entropy of the parent node from the weighted average entropy of the child nodes.

$$IG(S, A) = H(S) - H(S, A)$$

Where $H(S)$ is calculated as:

$$H(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

It is the most widely used criterion for building decision trees from a training dataset by assessing the information gain for each variable and choosing the variable that maximizes the information gain, which minimizes the entropy and better divides the dataset into groups for successful classification.

3.1.2 Gini Index:

When a function is randomly chosen, the Gini Index calculates the likelihood that it will be labelled incorrectly. As a result, for decision tree design, the function with the lowest gini index should be chosen. It ranges from 0 to 1, with

0 indicating that the classification is pure, i.e. all of the elements belong to a single class, and 1 indicating that only one class exists. The information is unadulterated. And 1 indicates the random distribution of elements across various classes. The value of 0.5 of the Gini Index shows an equal distribution of elements over some classes.

It can be calculated by deducting the sum of squared of probabilities of each class from one. i.e.

$$\text{Gini Index} = 1 - \sum_{i=1}^n (P_i)^2$$

Pseudo Code:

Algorithm Decision Tree

start

∀ attributes a_1, a_2, \dots, a_n

Find the attribute that best divides the training data using Information Gain and Gini Index

$a_best \leftarrow$ the attribute with highest information gain

Create a decision node that splits on a_best

Recurse on the sub-lists obtained by splitting on a_best and add those nodes as children of node

end

3.2 Naïve Bayes:

Naïve Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems. It is based on the assumption that all features are independent or unrelated. It is one of the simple and most effective Classification algorithms which help in building the fast machine learning models that can make quick predictions. It is probabilistic classifier, i.e. it calculates the probability of all the labels depending upon the given test case. The label with the highest probability is given as the result.

Bayes theorem provides a way of calculating posterior probability $P(c|x)$ from $P(c)$, $P(x)$ and $P(x|c)$. Look at the equation below:

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Above,

- $P(c|x)$ is the posterior probability of *class* (c , *target*) given *predictor* (x , *attributes*).
- $P(c)$ is the prior probability of *class*.
- $P(x|c)$ is the likelihood which is the probability of *predictor* given *class*.
- $P(x)$ is the prior probability of *predictor*.

Pseudo code:

Algorithm Naïve Bayes

start

Let $S = \{a_1, a_2, \dots, a_n\}$, where S = training set and a = articles:

Calculate the probability of the classes $P(C)$

Calculate likelihood of attribute A for each class $P(A/C)$

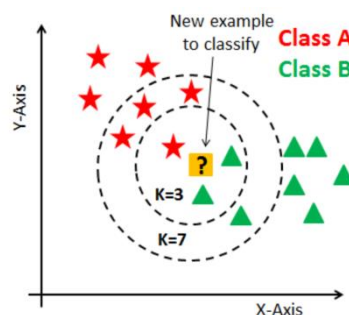
Calculate the conditional probability $P(C|A)$

Assign the class with the highest probability

end

3.3 K-Nearest Neighbor:

The K-Nearest Neighbor is one of the simplest algorithms used for classification. It is an instance-based classifier. When the k-NN is used, instances within a dataset are contained in a dimensional space, where a new instance is labeled based on its similarity with other instances. These instances are referred to as neighbors. A new instance is labeled x , if x is the most similar class for the neighboring observations. A distance function is applied to determine the similarity between instances. There are many distance functions e.g. Manhattan Distance, Euclidean Distance, Minkowsky Distance etc.



For the purpose of this study, the distance function employed is Euclidean. The Euclidean function is a relatively common method as it reflects the human perception of distance.

Pseudo code:

Algorithm k-Nearest Neighbour

start

Let $S = \{a_1, a_2, \dots, a_n\}$, where S represents the training set and a represents article documents

$k \leftarrow$ the desired number of nearest neighbours

Compute $d(x, y)$ between new instance i and all $a \in S$

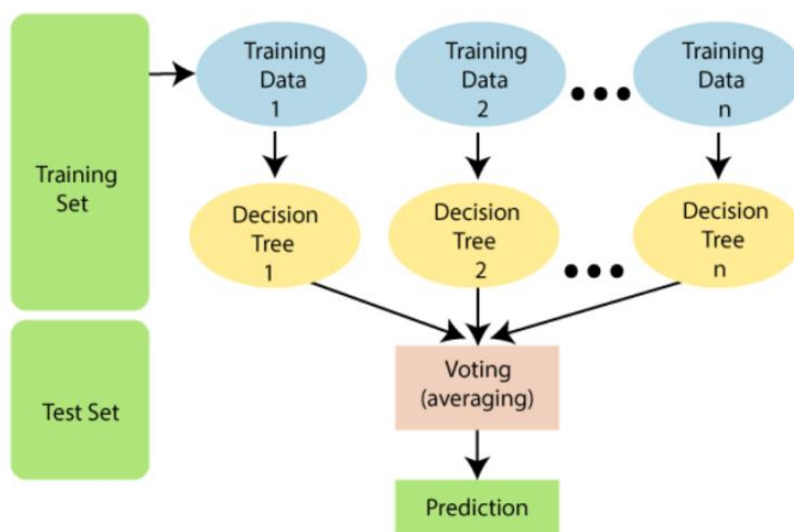
Select the k closest training samples to i

$Class_i \leftarrow$ best voted class

end

3.4 Random Forest:

Random forest is a supervised learning algorithm. The "forest" it builds, is a collection of decision trees, usually trained with the "bagging" method. The general idea of the bagging method is that a combination of learning models increases the overall result. So, it is based on the principle that instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output. The different trees are made using random data from the datasets and using all these trees the predictions are made.



Pseudo code:

Algorithm Random Forest

Require *IDT* (a decision tree inducer), *T* (the number of iterations), *S* (the training set), μ (the subsample size), *N* (the number of attributes used in each node)

start

$t \leftarrow 1$

repeat

$S_t \leftarrow$ Sample μ instances from *S* with replacement.

 Build classifier M_t using *IDT*(*N*) on S_t

$t++$

until $t > T$

end

3.5 Gradient Boost:

Gradient Boosting is a word made by combining both Gradient Descent and Boosting. Gradient boosting is a type of machine learning algorithm which is based on the principle that the best possible next model can be made if we combine previous models and try to minimize the overall prediction error. This happens by optimizing the loss function. So the difference between random forest and gradient boost is that in random forest all the trees are independent from each other while here the next tree is made based upon the predecessor. It mainly consists of 3 parts:

3.5.1 Weak Learner/Naïve Model

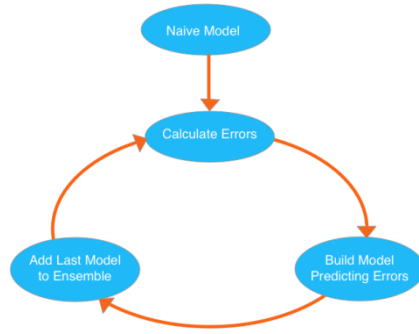
Decision trees are used as the weak learner in gradient boosting, also known as the Nave Model. Trees are built greedily, with the best split points chosen based on Gini scores or to mitigate the loss. It is collective to make the weak learners in unique ways, such as with the most layers, nodes, breaks, or leaf nodes possible. This is to ensure that the learners stay fragile but can still be greedily built.

3.5.2 Loss Function

The loss function used depends on the type of problem which is to be solved. It must be differentiable, but many standard loss functions are supported and we can define our own.

3.5.3 Additive Model

Current trees in the model are not updated, and new trees are introduced one at a time. When adding trees, a gradient descent technique is used to minimise the loss. Gradient descent has traditionally been used to reduce a number of parameters, such as the coefficients in a y on x equation or the weights in a neural network. After calculating error or loss, weights are updated to minimise that error. After calculating the loss, to perform the gradient descent procedure, we must add a tree to the model that reduces the lost (i.e. follow the gradient).



Pseudo Code:

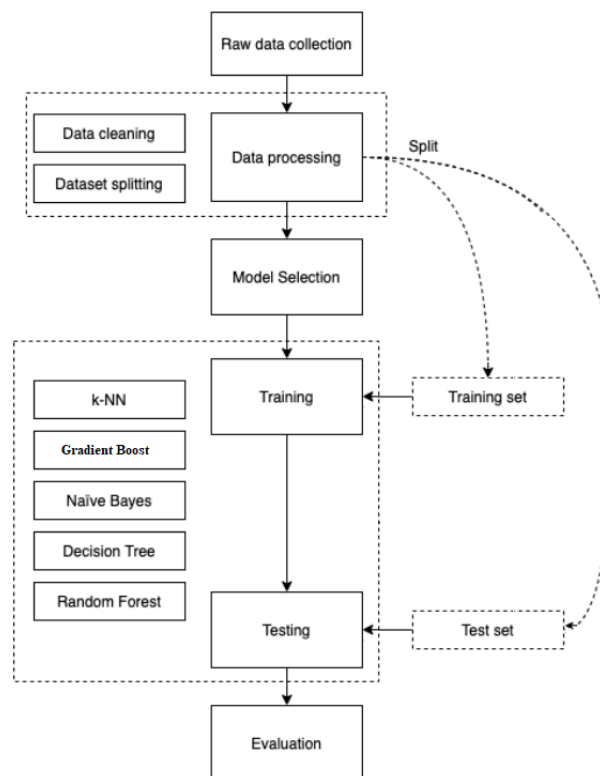
Algorithm	Gradient Tree Boosting Algorithm.
<ol style="list-style-type: none"> 1. Initialize $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$. 2. For $m = 1$ to M: <ol style="list-style-type: none"> (a) For $i = 1, 2, \dots, N$ compute $r_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}.$ (b) Fit a regression tree to the targets r_{im} giving terminal regions R_{jm}, $j = 1, 2, \dots, J_m$. (c) For $j = 1, 2, \dots, J_m$ compute $\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma).$ (d) Update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$. 3. Output $\hat{f}(x) = f_M(x)$. 	

4. Experimentation

This section mainly explains the details of the design and implementation of the proposed solution. The solution is implemented in Python 3. Firstly, an overview of the solution is presented, briefly describing the phases of this implementation. Section 4.2 describes the data preparation process, including details on data cleaning and transformation, and dataset splitting. The next section presents the modeling process, with a detailed account of the training and testing processes. And the last section contains details about testing phase.

4.1 Overview

Figure presents a flow chart of the supervised learning process adopted in this study,



as part of the proposed solution.

Firstly, DDoS raw data is gathered from the [link](#). After collection, data is processed to construct the final datasets for modeling. Data processing includes data cleaning, transformation of data types, and dataset splitting which is explained in section. This is followed by the model selection process. The models are trained using five different algorithms; k-nearest neighbor, Gradient Boost, Naïve Bayes, Decision tree, Random forest. Finally, the model is tested with unseen data. The results are evaluated using several performance metrics, as described in Section.

4.2 Data Preparation

4.2.1 Data Cleaning and Transformation

4.2.1.1 Missing data: Handling missing data is vital in machine learning, as it could lead to incorrect predictions for any model. Accordingly, null values are eliminated by propagating the last valid observation forward along the column axis.

```
In [18]: df.shape
```

```
Out[18]: (3894072, 87)
```

```
In [19]: df1 = df[df.isna().any(axis=1)]  
df1.shape
```

```
Out[19]: (27, 87)
```

As the number of records having missing data are very less. So, removing those records will not make any considerable difference.

```
In [20]: df=df.dropna()
```

```
In [21]: df.shape
```

```
Out[21]: (3894045, 87)
```

4.2.1.2 Infinite and large data: Handling infinite values is as important as null values as it hampers the working of algorithms. Also very large values also cause problems in the data.

```
In [27]: count = np.isinf(df).values.sum()  
count
```

```
Out[27]: 396926
```

```
In [28]: col_name = df.columns.to_series()[np.isinf(df).any()]  
print(col_name)
```

```
Flow Bytes/s      Flow Bytes/s  
Flow Packets/s    Flow Packets/s  
dtype: object
```

```
In [29]: print(df['Flow Bytes/s'].value_counts())
```

```
1.200000e+07    659587
4.580000e+08    261201
inf            198463
2.944000e+09    170937
8.020000e+08     69737
...
1.148087e+00         1
1.219178e+00         1
1.136842e+07         1
1.526512e+00         1
1.257952e+04         1
Name: Flow Bytes/s, Length: 199034, dtype: int64
```

```
In [30]: print(df['Flow Packets/s'].value_counts())
```

```
2.000000e+06    1898111
1.000000e+06     225619
inf            198463
4.166667e+04     117175
4.081633e+04      95690
...
2.979368e+01         1
3.507681e-01         1
1.525107e+01         1
1.201880e+00         1
2.034270e-01         1
Name: Flow Packets/s, Length: 175505, dtype: int64
```

As infinite values are very large so removing these can cause problems to the real model. Also it can be seen that values are very large. So can normalize it and also the infinite values be given finite value which will be larger than all the values so that its credibility doesn't go down will help.

```
In [31]: df['Flow Bytes/s']=df['Flow Bytes/s']/1000000 #for basic normalization we divided the data by 10^6
mm = df.loc[df['Flow Bytes/s'] != np.inf, 'Flow Bytes/s'].max()
mm=mm+200
df['Flow Bytes/s'].replace(np.inf,mm,inplace=True)
print(df['Flow Bytes/s'].value_counts())
```

```
1.200000e+01    659587
4.580000e+02    261201
3.144000e+03    198463
2.944000e+03    170937
8.020000e+02     69737
...
1.734579e-06         1
2.757282e+00         1
1.292287e-06         1
1.845178e-06         1
9.321774e-07         1
Name: Flow Bytes/s, Length: 199032, dtype: int64
```

```
In [32]: df['Flow Packets/s']=df['Flow Packets/s']/1000000
m = df.loc[df['Flow Packets/s'] != np.inf, 'Flow Packets/s'].max()
m=m+200
df['Flow Packets/s'].replace(np.inf,m,inplace=True)
print(df['Flow Packets/s'].value_counts())
```

```
2.000000e+00    1898111
1.000000e+00    225619
2.040000e+02    198463
4.166667e-02    117175
4.081633e-02     95690
...
5.673759e-03         1
3.340000e-07         1
1.917285e-07         1
2.165298e-07         1
2.526832e-07         1
Name: Flow Packets/s, Length: 175502, dtype: int64
```

4.2.1.3 Removing Unnecessary Features: The dataset has 87 columns out of which columns Timestamp, Source IP, Source Port, Destination IP, Destination Port, Flow ID are removed as these columns will affect the credibility of data because Timestamp and Flow ID are different for all records and Source IP, Destination IP, Source Port, Destination Port have just specific values as the dataset creators

```
In [33]: df.drop(columns=['Flow ID'], inplace=True)
df.drop(columns=['Timestamp'], inplace=True)
df.drop(columns=['Source IP'], inplace=True)
df.drop(columns=['Source Port'], inplace=True)
df.drop(columns=['Destination IP'], inplace=True)
df.drop(columns=['Destination Port'], inplace=True)
df
```

have used.

4.2.1.4 Transformation: The format of the collected data might not be suitable for modeling. In such cases, data and data types need to be transformed so that the data can then be fed into the models. Accordingly, some data features were transformed into numeric or float, since models do not perform well with strings, or do not perform at all.

```
In [25]: #SimillarHTTP is a object type so has to be converted to string
df['SimillarHTTP'] = df['SimillarHTTP'].astype('S')

#Label Encoders for converting the string data to integer values
encoder1=preprocessing.LabelEncoder()
df['SimillarHTTP']=encoder1.fit_transform(df['SimillarHTTP'])
df['Label']=encoder1.fit_transform(df['Label'])
```

4.2.1.5 Class Labels: Each dataset instance represents a snapshot of the network traffic at a given point in time. These instances are labeled according to the nature of the traffic i.e. whether the benign or some DDoS attack and in DDoS attack which type is it of (SYN Flood, UDP Flood, UDP Lag, MSSQL, NetBIOS, LDAP, and PortMap).So Classification is Multinomial.

```
print(df['Label'].value_counts())
```

```
Syn          1801189
MSSQL         694593
UDP           688602
NetBIOS       438261
LDAP          230013
Portmap       22480
BENIGN        14958
UDPLag        3976
Name: Label, dtype: int64
```

4.2.1.6 Splitting Datasets: A key characteristic of a good learning model is its ability to generalize to new or unseen, data. A model which is too close to a particular set of data is described as over fit, and therefore, will not perform well with unseen data. A generalized model requires exposure to multiple variations of input samples. Primarily, models require two sets of data, one to train and another to test. The training data is the set of instances that the model trains on, while the testing data is used to evaluate the generalizability of the model, that is, the performance of the model with unseen data.

```
#Dividing Data for training and testing
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.3, random_state=1) # 70% training and 30% test
print(X_train['Flow Packets/s'].value_counts())
#returning the training and testing data
return (X_train, X_test, Y_train, Y_test)
```

4.3 Training

During the training process, the selected algorithms are provided with training data to learn from to eventually create machine learning models. Accordingly, the training set is used. At this point in the process, the input data source needs to be provided and should contain the target attribute (class label). The training process involves finding patterns in the training set that map the input features with the target attribute. Based on the observed patterns, a model is produced which are as follows:

4.3.1 Decision Tree Classifier:

We built two decision tree classifiers with different criteria: One with Information Gain (Entropy) as criteria of splitting and other with gini index

4.3.1.1 With Entropy:

Function for creating and saving Decision Tree Classifier With Entropy Model

```
def decisiontree_entropy(X_train,Y_train):  
    from sklearn import tree  
    clf=tree.DecisionTreeClassifier(criterion='entropy')  
    clf.fit(X_train,Y_train)  
    save(clf, "c:/Users/Hp/PycharmProjects/pythonProject/decisiontree_entropy.mdl")  
    return True
```

```
decisiontree_entropy(X_train,Y_train)
```

```
saved
```

```
True
```

Feature Importance:

Also we can see the importance of every feature i.e. how much the classifier depends upon given feature.

```
importance_decision_entropy = decision_entropy .feature_importances_  
dictt=dict(zip(columns, importance_decision_entropy))  
importance_decision_entropy={k: v for k, v in sorted(dictt.items(), key=lambda item: item[1],reverse=True)}  
print ("{:<40} {:<10} ".format('Feature', 'Importance'))  
# print each data item.  
for key, value in importance_decision_entropy.items():  
    print ("{:<40} {:<10} ".format(key, value))
```

Feature	Importance
Min Packet Length	0.49640934178656027
Average Packet Size	0.24683216474284056
Max Packet Length	0.22378687853379026
ACK Flag Count	0.011530837511773704
Init_Win_bytes_forward	0.004629425168503786
Flow Bytes/s	0.004545919323931296
min_seg_size_forward	0.0017105906547327317
Total Length of Fwd Packets	0.0014545624668625171
Inbound	0.0014195048235718856
Fwd IAT Min	0.0009258695594770828
Avg Fwd Segment Size	0.0008873818205409081
Fwd Header Length.1	0.0008793732760365883
Fwd Header Length	0.0008312454151216615
Fwd Packet Length Max	0.0006344423431374986
Packet Length Mean	0.000629605856426146
Subflow Fwd Bytes	0.0006280177142926165
Fwd Packet Length Min	0.0005839164218679005
Flow IAT Mean	0.00024861068435902894
Fwd Packet Length Mean	0.00018802422036468425
Flow Duration	0.0001632982433701273
SYN Flag Count	0.0001256673339181236
Packet Length Variance	0.000116328450648585
Flow Packets/s	0.0001131596316423964
Flow IAT Min	0.00010350086971082603

4.3.1.2 With Gini Index:

Function for creating and saving Decision Tree Classifier With GINI INDEX Model

```
def decisiontree_gini(X_train,Y_train):  
    from sklearn import tree  
    clf=tree.DecisionTreeClassifier()  
    clf.fit(X_train,Y_train)  
    save(clf, "C:/Users/Hp/PycharmProjects/pythonProject/decisiontree_gini.mdl")  
    return True
```

```
decisiontree_gini(X_train,Y_train)
```

saved

True

Feature Importance:

```
importance_decision_gini = decision_gini.feature_importances_  
dictt=dict(zip(columns, importance_decision_gini))  
importance_decision_gini={k: v for k, v in sorted(dictt.items(), key=lambda item: item[1],reverse=True)}  
print("{:<40} {:<10}".format('Feature', 'Importance'))  
# print each data item.  
for key, value in importance_decision_gini.items():  
    print("{:<40} {:<10}".format(key, value))
```

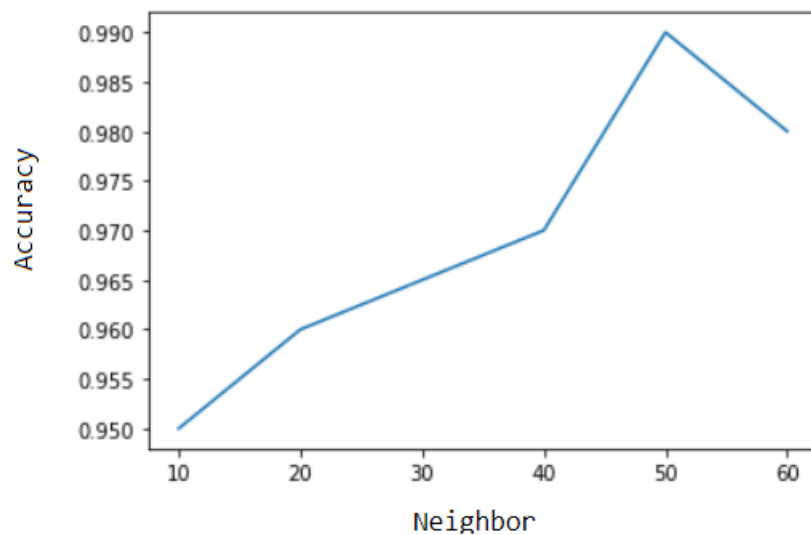
Feature	Importance
ACK Flag Count	0.4601559073777657
Average Packet Size	0.20715925332628973
Fwd Packet Length Max	0.19211514839534263
Max Packet Length	0.11650024082875801
Min Packet Length	0.008901870212226257
Inbound	0.0031265227968344887
Fwd Packet Length Mean	0.0022784878376733014
Flow Bytes/s	0.0015899952629221258
min_seg_size_forward	0.0014492606389429593
Total Length of Fwd Packets	0.0012729817128562264
Fwd Packet Length Min	0.0009164165072383531
Fwd Header Length	0.0007836354706048853
Packet Length Mean	0.0005776801781685499
Fwd Header Length.1	0.0004978169048885016
Subflow Fwd Bytes	0.0004023908303609716
Avg Fwd Segment Size	0.00034146230243139477
Fwd IAT Min	0.0002862950529003008
Flow Packets/s	0.00019118375038008531
Packet Length Std	0.00018711527333152895
Init_Win_bytes_forward	0.00018267462501565728
SYN Flag Count	0.0001237970701384819
Flow IAT Max	0.00011608703131644967

4.3.2 K Nearest Neighbor:

K nearest neighbor is tested for different values of Neighbor (with a very small part of dataset) whose accuracy is as follows:

Neighbor	Accuracy
10	0.95
20	0.96
30	0.965
40	0.97
50	0.99
60	0.98

This can be represented in form of graph as :



So we used Neighbor=50 in our model.

Function for creating and saving K Nearest Neighbor Classifier Model

```
def KNearestNeighbor(X_train,Y_train,neighbor,name):  
    model = KNeighborsClassifier(n_neighbors=neighbor)  
    model.fit(X_train,Y_train)  
    save(model, name)  
    return True
```

```
KNearestNeighbor(X_train,Y_train,50,"C:/Users/Hp/PycharmProjects/pythonProject/KN50.mdl")
```

saved

True

4.3.3 Naïve Bayes Classifier:

There are three types of Naïve Bayes classifier present in Sklearn.

- Multinomial Naïve Bayes
- Bernoulli Naïve Bayes
- Gaussian Naïve Bayes

For this dataset Multinomial Naïve Bayes can't be used as it is used for discrete counts and also Bernoulli Naïve Bayes can't be applied as for these features should be binary data but in our dataset it is not true.

So we used Gaussian Naïve Bayes classifier.

Function for creating and saving Naive Bayes Model

```
def Naive_Bayes(X_train,Y_train):  
    model = GaussianNB()  
    model.fit(X_train,Y_train)  
    save(model, "C:/Users/Hp/PycharmProjects/pythonProject/Naive_Bayes.mdl")  
    return True
```

```
Naive_Bayes(X_train,Y_train)
```

```
saved
```

```
True
```

4.3.4 Gradient Boost Classifier:

For this paper we used XGBoost Classifier because it is fast and also gives more accuracy i.e. it is basically updated and most popular Gradient Boost Classifier.

Function for creating and saving Gradient Boost Classifier Model

```
def GradientBoost(X_train,Y_train,lr,name):  
    model = xgb.XGBClassifier(learning_rate=lr)  
    model.fit(X_train,Y_train)  
    save(model, name)  
    return True
```

```
GradientBoost(X_train,Y_train,0.5,"C:/Users/Hp/PycharmProjects/pythonProject/GradientBoost05.mdl")
```

```
saved
```

```
True
```

Feature Importance:

```
importance_gd = clf_gd.feature_importances_  
dictt=dict(zip(columns, importance_gd))  
importance_gd={k: v for k, v in sorted(dictt.items(), key=lambda item: item[1],reverse=True)}  
print ("{:<40} {:<10} ".format('Feature', 'Importance'))  
# print each data item.  
for key, value in importance_gd.items():  
    print ("{:<40} {:<10} ".format(key, value))
```

Feature	Importance
ACK Flag Count	0.8490367531776428
Fwd Packet Length Std	0.07568186521530151
Min Packet Length	0.025686411187052727
Fwd Packet Length Max	0.012024401687085629
Average Packet Size	0.009461821056902409
Fwd Packet Length Min	0.008905714377760887
Fwd Packet Length Mean	0.007664439734071493
Max Packet Length	0.007374416571110487
URG Flag Count	0.0011450940510258079
Inbound	0.0009124670177698135
Init_Win_bytes_forward	0.0005865833954885602
Total Fwd Packets	0.00030265495297499
Total Length of Fwd Packets	0.00023466389393433928
Packet Length Mean	0.00019598338985815644
min_seg_size_forward	0.00010507513070479035
SYN Flag Count	8.282488124677911e-05

4.3.5 Random Forest Classifier:

In this paper we as in Decision Tree Classifier have implemented two Random Forest Classifier with different splitting criteria one using entropy and other gini index. We have used estimators=100 i.e. decision trees made are 100.

4.3.5.1 Using Entropy:

Function for creating and saving Random Forest Classifier With Entropy Model

```
def random_forest_entropy(X_train,Y_train):  
    model = RandomForestClassifier(n_estimators=100, criterion="entropy")  
    model.fit(X_train, Y_train)  
    save(model, "C:/Users/Hp/PycharmProjects/pythonProject/random-forest-entropy.mdl")  
    return True
```

```
random_forest_entropy(X_train,Y_train)
```

saved

True

Feature Importance:

```
importance_entropy = clf_entropy.feature_importances_  
dictt=dict(zip(columns, importance_entropy))  
importance_entropy={k: v for k, v in sorted(dictt.items(), key=lambda item: item[1],reverse=True)}  
print ("{:<40} {:<10} ".format('Feature', 'Importance'))  
# print each data item.  
for key, value in importance_entropy.items():  
    print ("{:<40} {:<10} ".format(key, value))
```

Feature	Importance
Min Packet Length	0.1516582390898866
Avg Fwd Segment Size	0.10319791081195219
Average Packet Size	0.10243897039431166
Fwd Packet Length Min	0.09237528551001967
Fwd Packet Length Max	0.08174898895884546
Packet Length Mean	0.07729604928884745
Fwd Packet Length Mean	0.07418742326418588
Max Packet Length	0.05945193694944073
Subflow Fwd Bytes	0.051239177161365404
Total Length of Fwd Packets	0.04550593225797074
Flow Bytes/s	0.03635540247212011
Protocol	0.020263168962096047
ACK Flag Count	0.018727358730789064
Init_Win_bytes_forward	0.015746603998996606
Packet Length Std	0.0072441169618582355
Packet Length Variance	0.006783049656778316
min_seg_size_forward	0.005141247071090835
Fwd Header Length	0.005064907387462579
Flow Duration	0.004364465270563021
Fwd Packet Length Std	0.004316034545065283
Flow IAT Std	0.00304284067249798
Fwd IAT Max	0.002607040638415359
act_data_pkt_fwd	0.0025693904794002174
Fwd IAT Total	0.002568781121420273

4.3.5.2 Using Gini Index:

Function for creating and saving Random Forest Classifier With GINI INDEX Model

```
def random_forest_gini(X_train,Y_train):  
    model = RandomForestClassifier(n_estimators=100, criterion="gini")#max_depth=15  
    model.fit(X_train, Y_train)  
    save(model, "C:/Users/Hp/PycharmProjects/pythonProject/random-forest-gini.mdl")  
    return True
```

```
random_forest_gini(X_train,Y_train)
```

```
saved
```

```
True
```

Feature Importance:

```
importance_gini = clf_gini.feature_importances_  
dictt=dict(zip(columns, importance_gini))  
importance_gini={k: v for k, v in sorted(dictt.items(), key=lambda item: item[1],reverse=True)}  
print ("{:<40} {:<10} ".format('Feature', 'Importance'))  
# print each data item.  
for key, value in importance_gini.items():  
    print ("{:<40} {:<10} ".format(key, value))
```

Feature	Importance
Max Packet Length	0.11129622920770219
Packet Length Mean	0.09949788110653787
Fwd Packet Length Mean	0.08978768361272463
Fwd Packet Length Min	0.082308915592799
Fwd Packet Length Max	0.07948099233888806
Average Packet Size	0.0714209534437103
Min Packet Length	0.06650453244657727
ACK Flag Count	0.06295935865501383
Avg Fwd Segment Size	0.058923114953712864
Total Length of Fwd Packets	0.053642089216641695
Subflow Fwd Bytes	0.050678226447501584
Init_win_bytes_forward	0.04385040786410594
Protocol	0.024315066262049524
Flow Bytes/s	0.01783189049094836
Fwd Header Length.1	0.009547008362278472
Fwd Packet Length Std	0.007639142147601071
min_seg_size_forward	0.007070725536660922
Packet Length Variance	0.006548380174441645

4.4 Testing

In the last stage of experimentation the models are tested with unseen data. The unseen data used at this stage is the resulting test set from the data split (20%). Testing is conducted to assess how a model represents data and how well it will perform in the future. This study ensured that any tweaks to the models were done prior to testing, so that the testing data is used only once. Various performance metrics were generated to be able to analyse the performance of the DDoS datasets, such as accuracy, precision, recall, and F-measure and confusion metrics. These are described in the next section.

5 Results and Discussions:

A crucial part of understanding the performance of a model is generating performance metrics. In this study, various metrics are generated. These are described below.

5.1 Confusion Matrix

A confusion matrix is a very important metric in analyzes of any machine learning model. It is basically a way to represent the summary of predicted values in tabular form. In this the columns represent the count of instances in a predictive class and the rows represent the count of instances in an actual class.

The confusion Matrix for all the classifiers used are as follows:

Decision Tree With Entropy

	pred:0	pred:1	pred:2	pred:3	pred:4	pred:5	pred:6	pred:7
true:0	239	0	0	0	0	0	0	0
true:1	0	3750	13	1	0	0	0	0
true:2	0	140	10974	0	0	1	57	0
true:3	0	0	0	7040	2	0	1	0
true:4	0	0	2	354	1	1	0	0
true:5	0	0	0	0	0	29034	3	8
true:6	0	0	145	2	0	1	10975	7
true:7	0	0	0	0	0	1	8	43

Decision Tree With Gini

	pred:0	pred:1	pred:2	pred:3	pred:4	pred:5	pred:6	pred:7
true:0	238	0	0	0	0	1	0	0
true:1	0	3749	15	0	0	0	0	0
true:2	0	139	10974	1	0	1	57	0
true:3	0	0	0	7040	2	0	1	0
true:4	0	0	2	354	1	1	0	0
true:5	0	0	0	0	0	29034	3	8
true:6	0	0	143	2	0	0	10978	7
true:7	0	0	0	0	0	1	8	43

KNN

	pred:0	pred:1	pred:2	pred:3	pred:4	pred:5	pred:6	pred:7
true:0	220	0	0	0	0	1	0	18
true:1	0	3453	9	0	0	0	302	0
true:2	0	28	10219	1	0	909	15	0
true:3	0	0	0	6482	557	3	1	0
true:4	0	0	1	104	253	0	0	0
true:5	0	0	2274	0	0	26771	0	0
true:6	0	902	26	1	0	4	10195	2
true:7	6	0	0	0	0	0	0	46

Gradient Boost

	pred:0	pred:1	pred:2	pred:3	pred:4	pred:5	pred:6	pred:7
true:0	239	0	0	0	0	0	0	0
true:1	0	3750	13	1	0	0	0	0
true:2	0	138	10991	0	0	0	43	0
true:3	0	0	0	7041	1	0	1	0
true:4	0	0	1	354	0	2	1	0
true:5	0	0	0	0	0	29034	3	8
true:6	0	0	142	1	0	0	10978	9
true:7	0	0	0	0	0	1	8	43

Random Forest with Entropy

	pred:0	pred:1	pred:2	pred:3	pred:4	pred:5	pred:6	pred:7
true:0	239	0	0	0	0	0	0	0
true:1	0	3753	11	0	0	0	0	0
true:2	0	139	10979	0	0	0	54	0
true:3	0	0	0	7041	1	0	1	0
true:4	0	0	2	354	1	1	0	0
true:5	0	0	0	0	0	29034	3	8
true:6	0	0	141	2	0	0	10980	7
true:7	0	0	0	0	0	1	8	43

Random Forest with Gini

	pred:0	pred:1	pred:2	pred:3	pred:4	pred:5	pred:6	pred:7
true:0	239	0	0	0	0	0	0	0
true:1	0	3751	13	0	0	0	0	0
true:2	0	139	10977	0	0	0	56	0
true:3	0	0	0	7041	1	0	1	0
true:4	0	0	1	354	1	2	0	0
true:5	0	0	0	0	0	29033	3	9
true:6	0	0	143	1	0	0	10979	7
true:7	0	0	0	0	0	1	8	43

Naive Bayes

	pred:0	pred:1	pred:2	pred:3	pred:4	pred:5	pred:6	pred:7
true:0	39	29	7	1	147	7	3	6
true:1	0	3086	0	0	678	0	0	0
true:2	1	8478	3	1	2682	7	0	0
true:3	8	4839	0	0	2196	0	0	0
true:4	0	272	0	0	86	0	0	0
true:5	347	10767	17	120	14834	2926	13	21
true:6	0	3645	190	0	2008	0	5285	2
true:7	1	0	0	0	50	1	0	0

5.2 Accuracy: Accuracy is defined as the ratio of how much the instances have been correctly identified to the total instances. It is one of the most important and simple metric for evaluating the classification models. Mathematically it can be represented as:

$$\text{Accuracy} = \frac{\text{Correctly classified instances}}{\text{Total instances}} \times 100\%$$

5.3 Precision: Although accuracy is a good metric to find whether the model of being trained correctly or not but it does not give information on detailed information on the specific application. Because of which other performance metrics are to be used, and one of this is precision. Precision can be defined as the ratio of correctly identified positives to the total positives predicted by the model.

$$\text{Precision} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

5.4 Recall: Recall is defined as the ratio of correctly identified positives to the total positives present in the test cases. It is very important metric for such datasets in which the prediction of true is very important like predicting any disease.

$$\text{Recall} = \frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$$

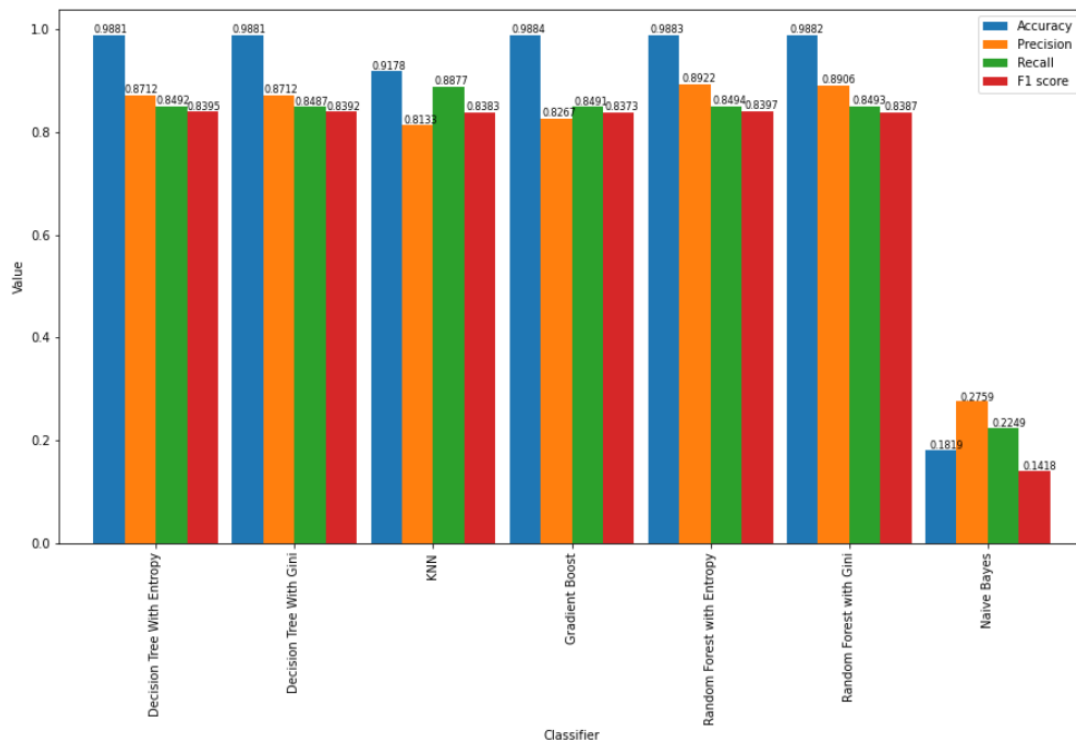
5.5 F-measure: The F-measure is a metric that provides an overall accuracy score for a model by combining both precision and recall. If F-measure score comes to be fairly highly then it can be concluded that the model has both low false negatives and low false positives.

$$F\text{-measure} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

The following Table presents the evaluation metrics of the machine learning models including accuracy, precision, recall, F-measure:

	Accuracy	Precision	Recall	F1 score
Decision Tree With Entropy	0.988106	0.871213	0.849193	0.839487
Decision Tree With Gini	0.988122	0.871243	0.848670	0.839242
KNN	0.917775	0.813344	0.887743	0.838256
Gradient Boost	0.988424	0.826734	0.849085	0.837309
Random Forest with Entropy	0.988329	0.892208	0.849422	0.839683
Random Forest with Gini	0.988233	0.890590	0.849318	0.838729
Naive Bayes	0.181918	0.275895	0.224891	0.141753

This can be represented in graphical way as:



The goal of this evaluation is to analyze the performance of different classifiers in terms of their capacity to detect and classify different DDoS attacks. For analyzing it should be kept in mind that the dataset is highly unbalanced. So instead of accuracy precision, recall and f1 score are more important.

From the results shown above, it shows that based upon accuracy Gradient Boost Performs best as its accuracy is 98.84% but in terms of precision, recall and F1 score Random Forest Classifier performs the best. On the other hand, the Naïve Bayes Classifier performs worst in all the parameters.

6 Conclusions and Future Work

6.1 Conclusions

With increase in complexity in DDoS attacks methods, the detection and classification of DDoS attacks is becoming difficult day by day. It has increased the risk for heavy losses to the service based companies and are also a threat to countries security.

Machine Learning models are anomaly detection techniques, which are accurate and practical methods to identify DDoS traffic from legitimate traffic. These ML-based classifiers can be trained and tested using a real-life Intrusion Detection datasets, for better performance in real scenarios.

Seeing the above, we have used five such ML algorithms. The results obtained in our analysis revealed that our study is able to address the main question raised in this paper. i.e. which machine learning techniques is best for detecting and classifying DDoS attacks. From the above studies it can be concluded that Random Forest Classifier performs the best with accuracy of 99.83% and F1 score of 83.97%.

6.2 Future Work

In this paper we have used various machine learning approaches but interesting area that could be explored is how we can do the same classification using various hybrid algorithms and artificial neural networks and further be analyzed with deep learning, once the datasets could be represented in non-structured forms of data. The above can also be seen as two separate problems, which the future study could expand on.

Also the study can be done on other datasets and can see how these datasets perform in these scenarios.

7 APPENDIX

The complete code and dataset can be seen from the link:

<https://github.com/rohit04445/Detection-and-classification-of-DDoS-attacks>