



Distributed Denial of Service Attack Prediction Using Novel Ensemble Model

Nikhil Anand Mahendrakar¹, Sanjay A S², Manoj M³

Siemens Healthineers, Bengaluru, India¹⁻³

Abstract: In the modern world, internet has become a major part in everyone's life. Most of businesses have now gone digital and internet is playing a major role in their success. Internet has enabled businesses to have larger market to which they can sell their products. One important component of the internet is the server. All the devices in the internet are broadly either a server or a client. Server stores the data that are then accessed by the clients. Distributed Denial-of-Service is a network attack on the servers. This attack targets the availability of the server by overwhelming them by a flood of internet traffic. This prevents the server from giving services to the legitimate users. DDoS attack causes significant losses to the organization. In this paper, we discuss an ensemble modelling based approach to mitigate this problem. Ensemble modelling is a process where multiple models which are different fundamentally are combined to make one classifier model. We try to predict the type of DDoS attack by developing an ensemble model by combining Naïve Bayes, Random Forest, Multilayer Perceptron, Stochastic Gradient Descent. The accuracy score of 98.62 percent has been achieved. It can be concluded that this system proves to be an effective deterrent for DDoS attacks.

Keywords: DDoS attack, http-flood, udp-flood, smurf, machine learning, ensemble model

I. INTRODUCTION

Advancement in communication and information technology further strengthens the role of internet in business. Internet is widely used in organization for marketing and promotion of products and services. The concept of the internet (sometimes also referred to as World Wide Web) was conceived by a British computer scientist Tim-Berners Lee in the year 1989. The number of people coming online since then has increased exponentially. It is also important to note here that internet is a very recent phenomenon spanning just decades ago. We are still in very early stages of the world wide web technology. The internet has indeed changed a lot of things but there is still scope of a lot of development left.

As discussed earlier, the increased popularity of the internet has brought various advantages as well as new challenges. One of the important challenge for data stored in the internet is to make the data always available to the user (Availability), there must be no illegal edits to the data (Integrity) and there should be privacy for the data (Confidentiality). This is also popularly called as the CIA triad.

Distributed Denial of Service attack is caused by various freely available compromised hosts which create a fake flood of internet traffic. These compromised hosts are also sometimes referred to as "zombies". Distributed Denial of Service is used to violate the foundational security principle of the internet – CIA (Confidentiality, Integrity, Availability).

It is usually not possible for an attacker to cause a Denial of Service attack using a single host. Therefore, the attacker attacks the various vulnerable devices present on the internet with malware. Once, these devices are compromised, then the attacker uses those devices for causing a flood of internet traffic to the server. Usually the owners of those devices are not even aware that they have been compromised. The identity of the attacker is also hidden since the attack is taking place through the compromised devices which seem as legitimate clients. The server is then serving the compromised devices and the legitimate clients are not able to access the server. This causes huge loss to the organization.

In this paper we try to solve the following problem using machine learning based ensemble modelling technique. Most of the datasets that are available publicly are not equipped with the latest DDoS attacks. The dataset used in this paper has covered the latest DDoS attacks that are being carried out. The dataset has around 2 million rows and 28 columns. Feature Selection on the attributes has been done to reduce the number of features for the prediction without affecting the accuracy. Feature Engineering (Hashing based approach discussed in detail in the Implementation section) on the values of the attributes has been done to convert from string representation to numerical representation for various computing purposes. An ensemble based modelling has been proposed between Naïve Bayes, Random Forest, Multilayer Perceptron and Stochastic Gradient Descent, for achieving a higher accuracy by training with much less data.

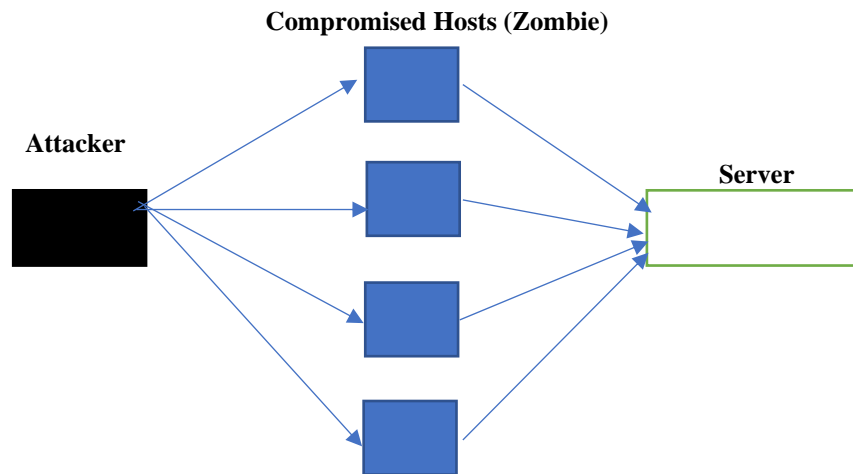


Fig. 1 Typical Structure of a DDoS attack

II. LITERATURE SURVEY

In paper [1], author has used Random Forest Classifier technique for classifying the packet into normal or anomalous packets. The dataset used is the famous NSL-KDD dataset which contains just 23000 odd rows. The number of features used in the dataset is 42. Similar accuracy can be achieved with much less attributes. More number of attributes increases the time for model building and also increases the complexity. Also, it is a bi-classification problem that means it classifies the packet into two classes normal and anomalous. Our paper presents a method to detect accurately which type of DDoS attack it is.

In paper [2], author has used the combination of both supervised and unsupervised machine learning technique. Agglomerative Clustering is used at the first step to form a cluster of packets which is similar to each other and then they are labelled manually. This labelled data, then acts as an input for building the model. Author has used k-Nearest Neighbors (kNN), Support Vector Machine (SVM) and Random Forest (RF) for building the model and achieved the accuracy of 95%, 92% and 96.6%. The accuracy achieved in this paper is around 98 percent.

In paper [3], author used the firewall based approach. A firewall is a system which acts as the first point of entrance for the packet to enter into a server or a computing device. Access control list (ACL) is used to check whether the packet is authentic or not. The IP addresses are used for differentiating anomalous and normal addresses. Firewalls do not provide efficient solution when the attack size increases.

In paper [4], author has used to classify the ICMPv6 DDoS attack. ICMP is an abbreviation for Internet Control Message Protocol which is used for sending control and error messages. They have used popular classification techniques like Neural Network, Decision Tree, Naïve Bayes, Support Vector Machine and K-Nearest Neighbors and compared the results.

In paper [5], the author uses Ensemble modelling technique, where 4 classifier models are trained and voting is done among the models for the final prediction. The dataset used in this research paper is the famous NSL-KDD dataset. The four classification algorithms used are MLP (NN), SMO (SVM), IBK (KNN) and J48 (DT-C4.5). Hard Voting is used to combine the prediction of different models to get the final outcome. However, Hard Voting ignores the fact that each model has different accuracy and gives equal weightage to each model while giving the final outcome.

In paper [6], The authors propose a RNN-based model, referred to as DeepDefense. Unlike previously discussed methods, the model does not use machine learning based techniques. Instead, the authors extracted 20 features from packet headers and applied sliding windows to separate continuous network traffic into sequences of network packets. For a given sequence of packets, the model classifies the last packet either as a legitimate or attack traffic. Their proposed deep learning model achieves lower error rate than traditional machine learning techniques. However, reliance on packet-by-packet.



III. METHODS AND MATERIALS

The following chapter is further divided into multiple sub-chapters where each sub-chapter describes the multiple processes for achieving the final result.

A. Data

The data is taken from the public dataset [7]. The data consists of the modern DDoS (Denial Distributed of Service) attacks. The attacks include the HTTP Flood, UDP Flood, SIDDOS, Normal, Smurf. These attacks do not cover the entire DDoS attacks but contain the most common ones. The data is present in the text format with the values for each features separated using comma. There are around 2 million rows and 28 features. We believe we have trained the model with all possible scenario. Some of the features which do not contribute in the enhancement of the accuracy have been removed.

```

3,24,3,389693,21,23,top,1540,-----,4,11339,16091,24780100,Switch1,Router,35,529786,35,529786,35,539909,0,328,240918,505490,1540,0,236321,0,35,519662,35,530032,1,50,02192,Normal
15,24,15,201196,23,24,top,1540,-----,16,6274,16092,24781700,Router,server1,20,176725,20,176725,20,186848,0,328,205808,505437,1540,0,236337,0,20,156478,20,186848,1,50,030211,Normal
24,15,15,61905,23,22,ack,55,-----,16,1930,16092,883060,Router,Switch2,7,049955,7,049955,7,059958,0,328,206042,18051,3,55,0,008441,0,7,039952,7,069962,1,030045,50,060221,UDP-Flood
24,9,9,443135,23,21,ack,55,-----,10,12670,16085,884675,Router,Switch1,39,62797,39,62797,39,637973,0,328,064183,18043,5,55,0,008437,0,39,617967,39,647976,1,030058,50,060098,Normal
24,8,8,157335,23,21,ack,55,-----,9,4901,16088,884840,Router,Switch1,16,039806,16,039806,16,049810,0,328,113525,18046,2,55,0,008438,0,16,029803,16,059813,1,030054,50,061864,Normal
24,1,1,219320,21,1,ack,55,-----,2,6837,16091,883005,Switch1,client-1,21,885768,21,885768,21,895771,0,328,297902,18056,4,55,0,008440,21,865762,21,895771,1,030016,50,043427,Normal
24,13,13,480053,24,23,ack,55,-----,14,13609,16103,885665,server1,Router,42,45032,42,45032,42,460323,0,328,460278,18065,3,55,0,008446,0,42,45032,42,48033,1,030032,50,055747,Normal
24,1,1,224,22,551227,24,23,ack,55,-----,3,15392,16091,883005,server1,Router,47,910078,47,910078,47,920081,0,328,26412,18054,5,55,0,008440,47,910078,47,940088,1,030022,50,048477,Normal
24,2,2,399941,23,22,ack,55,-----,3,11595,16091,24780100,Switch1,Router,36,314926,36,314926,36,325603,0,000554,328,26404,305526,1540,0,236321,0,001724,36,304803,36,336896,1,50,018467,Normal
9,2,24,51,33450,23,24,ack,1500,-----,10,1279,9108,13662000,router,server1,32,180177,32,181257,32,181257,0,00108,1016,496869,1524750,1500,0,130291,0,03252,32,158857,32,191377,1,9,960185,UDP-Flood
24,11,11,356924,23,22,ack,55,-----,12,10505,16103,885665,Router,Switch2,33,015317,33,015317,33,025230,0,328,522947,18068,8,55,0,008446,0,33,005314,33,035233,1,030019,50,046382,Normal
24,5,5,349309,21,5,ack,55,-----,6,10306,16091,883005,Switch1,client-5,49,037576,49,037576,49,047579,0,328,204851,18051,3,55,0,008440,49,01757,49,047579,1,030042,50,057349,Normal
6,1,24,26,629078,23,24,ack,1000,-----,27,5629,6250,6250000,Router,server1,70,05264,70,05312,70,0632,0,00048,124,942027,124942,1000,0,059605,0,00096,70,032,70,0632,25,75,0232,Normal
24,5,5,565692,21,5,ack,55,-----,6,15761,16091,883005,Switch1,client-5,49,037576,49,037576,49,047579,0,328,204851,18051,3,55,0,008440,49,01757,49,047579,1,030042,50,057349,Normal
24,1,1,478904,21,1,ack,55,-----,2,13574,16091,883005,Switch1,client-1,42,382176,42,382176,42,392179,0,328,297902,18056,4,55,0,008440,42,36217,42,392179,1,030016,50,043427,Normal
24,7,7,518117,21,7,ack,55,-----,3,14561,16090,884950,Switch1,client-7,45,38957,45,38957,45,399573,0,328,167767,18049,2,55,0,008440,45,369563,45,399573,1,030051,50,059851,Normal
24,6,6,13839,23,21,ack,55,-----,7,4223,16091,883005,Router,Switch1,2,524349,2,524349,2,534352,0,328,19309,18050,6,55,0,008440,2,514346,2,544355,1,030048,50,059112,Smurf
4,24,4,1,40886,4,21,top,1540,-----,5,4393,16091,24780100,client-4,Switch1,14,477384,14,477384,14,487307,0,328,217832,505455,1540,0,236321,0,14,477384,14,507754,1,50,025368,Normal
24,2,2,46,40094,2,1,ping,65535,-----,0,1,210,13762400,client-2,Switch1,44,4444,015243,0,23,188498,1519660,65535,0,2,131248,0,056214,44,4444,056214,0,9,056214,Smurf
24,11,11,337499,22,11,ack,55,-----,12,10014,16103,885665,Switch2,client-11,31,533911,31,533911,31,543914,0,328,522947,18068,8,55,0,008446,0,31,513905,31,543914,1,030019,50,046382,Normal
24,15,15,264811,22,15,ack,55,-----,16,8176,16092,883060,Switch2,client-15,25,957029,25,957029,25,967032,0,328,206042,18051,3,55,0,008441,0,25,957022,25,967032,1,030045,50,060221,Normal
24,13,13,201476,22,13,ack,55,-----,14,6278,16103,885665,Switch2,client-13,20,20168,20,20168,20,211683,0,328,460278,18065,3,55,0,008446,0,20,181674,20,211683,1,030032,50,055747,Normal
24,3,3,395290,21,3,ack,55,-----,4,11465,16091,883005,Switch1,client-3,35,968766,35,968766,35,978771,0,328,241051,18053,3,55,0,008440,35,94876,35,97877,1,030029,0,051929,Normal
14,24,14,186408,22,23,top,1540,-----,15,5812,16103,24798600,Switch2,Router,18,77443,18,77443,18,784540,0,328,431807,505785,1540,0,236498,0,18,764307,18,794677,1,50,029965,UDP-Flood
24,12,12,572328,23,22,ack,55,-----,13,15933,16103,885665,Router,Switch2,49,53893,49,53893,49,548933,0,328,491596,18067,5,50,008446,0,49,528927,49,538936,1,030026,50,051067,Normal
24,2,2,374184,24,23,ack,55,-----,3,10939,16091,883005,server1,Router,34,329077,34,329077,34,339081,0,328,26412,18054,5,55,0,008440,34,329077,34,339087,1,030022,50,048477,Normal
9,1,24,29,641250,21,23,ack,1000,-----,30,6233,6250,6250000,Switch1,Router,74,87408,74,8748,74,88488,0,00072,124,940828,124941,1000,0,059605,0,00144,74,864,74,89568,25,75,02368,Normal
9,24,9,452679,23,24,top,1540,-----,10,12912,16085,24779000,Router,server1,40,3718,40,3718,40,381923,0,328,063862,505218,1540,0,2362340,40,331554,40,331554,50,030008,Normal
24,5,5,4171,1801,5,ack,55,-----,6,13375,16091,883005,Switch1,client-5,41,789362,41,789362,41,799365,0,328,204851,18051,3,55,0,008440,41,769355,41,799365,1,030042,50,057349,Normal
24,9,9,142775,23,21,ack,55,-----,10,4445,16085,884675,Router,Switch1,14,664929,14,664929,14,674932,0,328,064183,18043,5,55,0,008437,0,14,654926,14,684936,1,030058,50,060098,Normal
24,2,2,44,13387,23,24,ack,1192,-----,23,489,8655,10316800,router,server1,13,478747,13,479731,13,489827,0,000924,962,684973,11,47520,1192,0,098389,0,032226,13,457571,13,489827,1,9,99048,UDP-Flood
5,2,24,47,27925,5,21,ack,1500,-----,26,1072,9108,13662000,client-5,Switch1,27,038857,27,038857,27,049457,0,00048,1016,522962,1524780,1500,0,130291,0,098759,27,038857,27,137616,1,9,959955,UDP-Flood
8,24,8,24510,23,24,top,1540,-----,9,756,16088,24775500,Router,server1,3,519381,3,519381,3,529704,0,328,11523,505294,1540,0,236278,0,3,499394,3,529704,50,031854,Normal
24,11,11,252941,24,23,ack,55,-----,12,7883,16103,885665,server1,Router,35,022235,35,022235,35,032238,0,328,522947,18068,8,55,0,008446,0,35,022235,35,032235,1,030019,50,046382,Normal
9,1,24,29,58826,9,21,ack,1000,-----,30,3481,6250,6250000,client-9,Switch1,52,848,52,848,52,85808,0,124,940828,124941,1000,0,059605,0,00144,52,848,52,87968,25,75,02368,Normal
24,8,8,161180,23,21,ack,55,-----,9,5022,16088,884640,Router,Switch1,16,402328,16,402328,16,412331,0,328,113525,18046,2,55,0,008438,0,16,392325,16,422334,1,030054,50,061864,Normal
24,14,14,205985,24,23,ack,55,-----,15,6451,16103,885665,server1,Router,20,606792,20,606792,20,616795,0,328,432001,18063,8,55,0,008446,0,20,606792,20,626802,1,030038,50,059974,Normal

```

Fig. 2 Sample of the Structure of the Data

```

1
@attribute SRC_ADD numeric 0
@attribute DES_ADD numeric 1
@attribute PKT_ID numeric 2
@attribute FROM_NODE numeric 3
@attribute TO_NODE numeric 4
@attribute PKT_TYPE {tcp,ack,cbp,ping} 5
@attribute PKT_SIZE numeric 6
@attribute FLAGS {-----A--} 7
@attribute FID numeric 8
@attribute SEQ_NUMBER numeric 9
@attribute NUMBER_OF_PKT numeric 10
@attribute NUMBER_OF_BYTE numeric 11
@attribute NODE_NAME_FROM {Switch1,Router,server1,router,client-4,client-2,Switch2,client-5,client-9,client-2,client-1,client-14,client-5,client-11,client-13,client-0,switch1,client-4,clienthttp,client-7,client-19,client-14,client-12,client-8,client-15,webserverlistin,client-18,client-1,switch2,client-6,client-10,client-7,webcache,client-10,client-17,client-16,client-17,client-18,client-12,client-8,client-0,client-16,client-13,client-11,client-5,client-3,client-9,client-19,http_client}
@attribute NODE_NAME_TO {Router,server1,Switch2,Switch1,client-1,client-5,client-7,switch1,client-11,client-15,client-13,client-3,client-9,client-6,router,client-4,client-14,switch2,client-8,clienthttp,webcache,client-10,client-12,webserverlistin,client-0,client-2,http_client,client-1,client-9,client-1,client-19,client-4,client-17,client-7,client-3,client-12,client-12,client-18,client-16,client-17,client-0,client-16,client-18,client-5,client-11,client-14,client-8,client-6,client-10,client-19,client-15}
@attribute PKT_IN numeric 14
@attribute PKT_OUT numeric 15
@attribute PKT_R numeric 16
@attribute PKT_DELAY_NODE numeric 17
@attribute PKT_RATE numeric 18
@attribute BYTE_RATE numeric 19
@attribute PKT_AVG_SIZE numeric 20
@attribute UTILIZATION numeric 21
@attribute PKT_DELAY numeric 22
@attribute PKT_SEND_TIME numeric 23
@attribute PKT_RESEVED_TIME numeric 24
@attribute FIRST_PKT_SENT numeric 25
@attribute LAST_PKT_RESEVED numeric 26
@attribute PKT_CLASS {Normal,UDP-Flood,Smurf,SIDDOS,HTTP-FLOOD} target variable(27).

```

Fig. 3 Features of the dataset



B. Data Pre-Processing

The data is present in the text format and requires pre-processing. In the file, each row is taken as one whole string and then converted to a string array using split function available in python. A separate list data structure is created for every feature. The value is then put to its respective list.

C. Conversion to Numeric Value

Some of the features had values which are of string datatype and had to be converted to a numeric value. The task here is to represent the each string value uniquely using a numeric value. The reason for this is that each of the string has a unique meaning associated with it. For example, consider the feature node_name_from which says from which type of node the packet is being sent. For this conversion, we use hashing. Hashing is the process of transforming string into a unique key and the value of the hash for a different string is completely different. The value returned by the hash function is also a string which can be then converted to a numeric value using ascii of the characters in the string. We have taken summation of the ascii value of the characters in the hash. This process makes sure that each unique string is assigned unique integer value for computation. However, for the values of the target variable we have assigned unique number as shown the figure 4.

```
for i in range(len(templist)):
    temptext = templist[i]
    hash_object = h.md5(temptext.encode())
    md5_hash = hash_object.hexdigest()
    value = 0
    for j in range(len(md5_hash)):
        value = value + ord(md5_hash[j])
    valuelist[templist[i]] = value
return valuelist
```

Fig. 3 Python Implementation for Conversion to Numeric Value

```
{'HTTP-FLOOD\n': 2, 'UDP-Flood\n': 3, 'SIDDOS\n': 4, 'Normal\n': 5, 'Smurf\n': 6}
```

Fig. 4 Numeric value assigned for the values of the target variable

D. Scaling of data

Scaling of data is important for getting the maximum accuracy for any classifier. Scaling can be defined as the process of fitting the data on a specific scale. Consider an example of a data which consists of currencies like rupee and dollar. We know that 1 US dollar = 80 rupees . So therefore, for consistent results, the whole data should be converted to either dollars or rupee. If scaling is not done on the data then the classifier considers a difference in price of 1 rupee as important as a difference of 1 US dollar. This would lead to wrong results and reduce the effectiveness of the classifier.

E. Ensemble Learning

In this paper we implement a ensemble model which is combination of four classifiers which are Stochastic Gradient Classifier , Multilayer Perceptron (also known as Artificial Neural Network) , Random Forest Classifier and Naïve Bayes Classifier. Brief working of each of the classifier has been given for completeness. Ensemble model is created to combine the results and improve the accuracy of the model .We combine the results of these models by using a mathematical function which is discussed below.

In the first step, we assign weights for each of the model. The weights signify how much the classifier shall contribute in giving the final prediction. Weights are the accuracy of the classifier when run separately on the dataset. Consider the Table I for weights assigned to the classifiers.



Table I

Serial No.	Classifier	Weights (Accuracy)
1	Stochastic Gradient Descent	98.50 %
2	Multilayer perceptron	98.62 %
3	Random Forest	97.23 %
4	Naïve Bayes	96.97 %

Figure 5 shows the 2 dimensional list containing the class probabilities of each data point for stochastic gradient descent. Class probabilities are nothing but the probability of the values of the target-variable given the values of the feature variable have already occurred. Now, in order to predict the final result we multiply the weight of the each classifier with the class probability. Once it is computed for all the classes (or values) of the target variable and the one which has the maximum value is chosen as the final prediction. Consider figure 6 which is the python implementation of the above idea. We design iteration on the values (class) of the target variable (type of distributed denial of service) and calculate according to the equation shown in figure 6. Once the calculation is completed for all the values of the target variable, then simply the one with maximum number is chosen as the predicted value.

```
[[0.00000826 0.01097373 0.00027594 0.98526049 0.00348158]
 [0.00000262 0.01066955 0.00022107 0.98515328 0.00395347]
 [0.00000122 0.01039839 0.00015766 0.98524586 0.00419688]
 ...
 [0.00000171 0.01053096 0.00019098 0.98517701 0.00409933]
 [0.00002233 0.01083333 0.00016394 0.98444659 0.0045338 ]
 [0.0000014 0.01058564 0.0002341 0.98541981 0.00375906]]
 0.9850214790307719
```

Fig. 5 two dimensional list of class probabilities of each data point for stochastic gradient descent

```
for i in range(len(X_test)):
    listofprobabilities.clear()
    dict_listofprobabilities.clear()
    j=0
    for key in dict_integer_target_variable:
        dict_listofprobabilities[dict_integer_target_variable[key]] = weightdecisiontree*y_probdecisiontree[i][j] + weightgnb*y_probgnb[i][j]+weightmultilayerperceptron*y_probmultilayerperceptron[i][j]+weightsgd*y_probsgd[i][j]
        listofprobabilities.append(weightdecisiontree*y_probdecisiontree[i][j] + weightgnb*y_probgnb[i][j]+weightmultilayerperceptron*y_probmultilayerperceptron[i][j]+weightsgd*y_probsgd[i][j])
        j=j+1
    listofprobabilities.sort()
    searchvalue = listofprobabilities[len(listofprobabilities)-1]
    counter = counter + 1
    for keys in dict_listofprobabilities:
        if float(dict_listofprobabilities[keys])==float(searchvalue):
            y_test_pred.append(keys)
print(counter)
```

Fig. 6 Python implementation of the ensemble model

For better understanding, we predict the value of the target variable for the first data point using our ensemble learning approach. The class probabilities of the first data point for naïve bayes figure 7, stochastic gradient descent figure 5, multilayer perceptron figure 8, random forest figure 9 is shown.

```
[[0. 0. 0. 0.99999998 0.00000002]]
```

Fig. 7 Naïve bayes class probabilities for first data point

```
[[0. 0.00866969 0.00013761 0.98697545 0.00421724]]
```

Fig. 8 Multilayer perceptron class probabilities for first data point



Fig. 9 Random forest class probabilities for first data point

Now, to predict the value of the target variable we multiply the weight and the class probability. The $f(x)$ represents the function of the ensemble learning and x takes all the values of the target variable. Therefore, the numeric value for the first class is $f(x_1) = 0.0 * 0.9697 + 0.0 * 0.9862 + 0.0 * 0.9723 + 0.00000826 * 0.985$ and $f(x_1) = 0$. $f(x_2) = 0.0 * 0.9697 + 0.00866969 * 0.9862 + 0.0 * 0.9723 + 0.010973 * 0.985$ and $f(x_2) = 0.0193591723$. Similarly, $f(x_3) = 0.0040$, $f(x_4) = 3.88593$, $f(x_5) = 0.00758839$. Clearly, $f(x_4)$ has the maximum value and x_4 is predicted as the output for the target variable. x_4 represents Normal packet. So, the first datapoint corresponds to a normal packet.

This approach takes into account that not all classifiers have same accuracy and classifiers which have better accuracy should have more contribution in the prediction. This method has increased the accuracy of each of the classifier and given a more robust, effective and more accurate model.

F. Naïve Bayes

Naive Bayes Classifier is a probability based classifier. The main mathematical concept behind the usage of this classifier is Bayes theorem. Bayes theorem is used to calculate the probability of an event A happening given that the event B has already occurred. Therefore, it can be said that the event B is the evidence and event A is the hypothesis. It is important to understand here that Naive Bayes assumes the features are independent. The presence of one particular feature does not affect the other. Since, broadly the classification problem can be further divided into either multiclass prediction or biclassification (yes or no). The mathematical formula differs for both the types. Consider the Figure 10 for biclassification problems and Figure 11 for multiclass problem [8].

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

Fig. 10. Bayes Theorem for bi-classification

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

Fig. 11 Bayes Theorem for multiclass classification

G. Stochastic Gradient Descent

Stochastic Gradient Descent Classifier has two main mathematical concepts. The first concept is the Gradient Descent Algorithm and the other part being Stochastic. Stochastic in plain term means random. Gradient Descent is an optimization algorithm. It is used to find local minima of a differential function. Consider an example of a parabola as shown in the figure 12. In order to find the local minima dx/dy should be equal to zero. Gradient Descent takes larger step size when the value is far from the local minima and small step sizes when the value is near the local minima. As mentioned earlier stochastic means random. However the dataset consists of multiple rows (or samples) and each row having multiple features and each contributing to the prediction of the target variable (type of Distributed Denial of Service attack). Stochastic Gradient Descent randomly takes one feature randomly instead of iterating over all features in the samples of the dataset. This reduces the computation enormously and makes it suitable as a classifier.

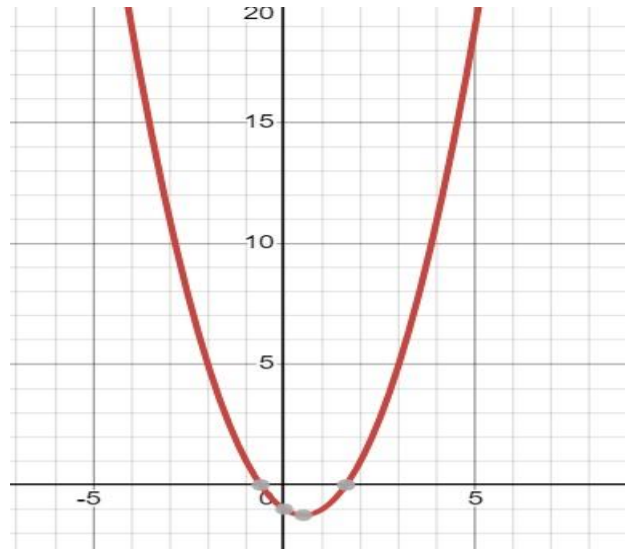


Fig. 12 Diagram of a parabola

H. Multilayer Perceptron

Multilayer Perceptron or sometimes also referred as the Neural Network or the Artificial Neural Network is based on the simple model of the brain. The building block of the Artificial Neural Network is the neurons. Each neuron takes weighted input signals and computes the output signal based on the activation function. There are mainly three layers in the multilayer perceptron – input layer, hidden layer, output layer. Usually, there is only one input layer, one output layer and many hidden layers. The activation function is the sigmoid function. The sigmoid activation function takes real numbers as input and converts it into value between 0 and 1.

I. Random Forest

Random forest Classifier applies decision tree algorithm on various subsets of the dataset and combines the results for prediction of the target variable. Decision tree constructs a flow-chart like tree structure where the internal node denotes the features, the branch denotes the values that the feature holds and the leaf node consists of the values held by the target variable.

IV. RESULTS AND DISCUSSION

The ensemble model is based on the mathematical equation as discussed in the implementation section. The ensemble model has an accuracy of 98.621 percent which is higher than all the individual classifier models. Though the difference in accuracy is in the range of 0.5 percent to 1.5 percent as compared to the famous classifiers used in this paper, however 1% increase in accuracy on a dataset containing 2 million dataset is predicting 20,000 more rows correctly which is a very good increase. Any machine learning model is judged based on the f1-score, precision and recall score.

Precision score quantifies the number of positive class predictions that actually belong to the positive class [9]. The figure 13 is the classification report. Recall quantifies the number of positive class predictions made out of all positive examples in the dataset. F-Measure provides a single score that balances both the concerns of precision and recall in one number. Confusion matrix as shown in the figure 14.



	precision	recall	f1-score	support
2	0.98	0.94	0.96	1352
3	1.00	0.90	0.95	66274
4	0.88	0.95	0.91	2130
5	0.99	1.00	0.99	639048
6	0.96	0.32	0.48	4217
accuracy			0.99	713021
macro avg	0.96	0.82	0.86	713021
weighted avg	0.99	0.99	0.99	713021

Fig. 13 Classification Report

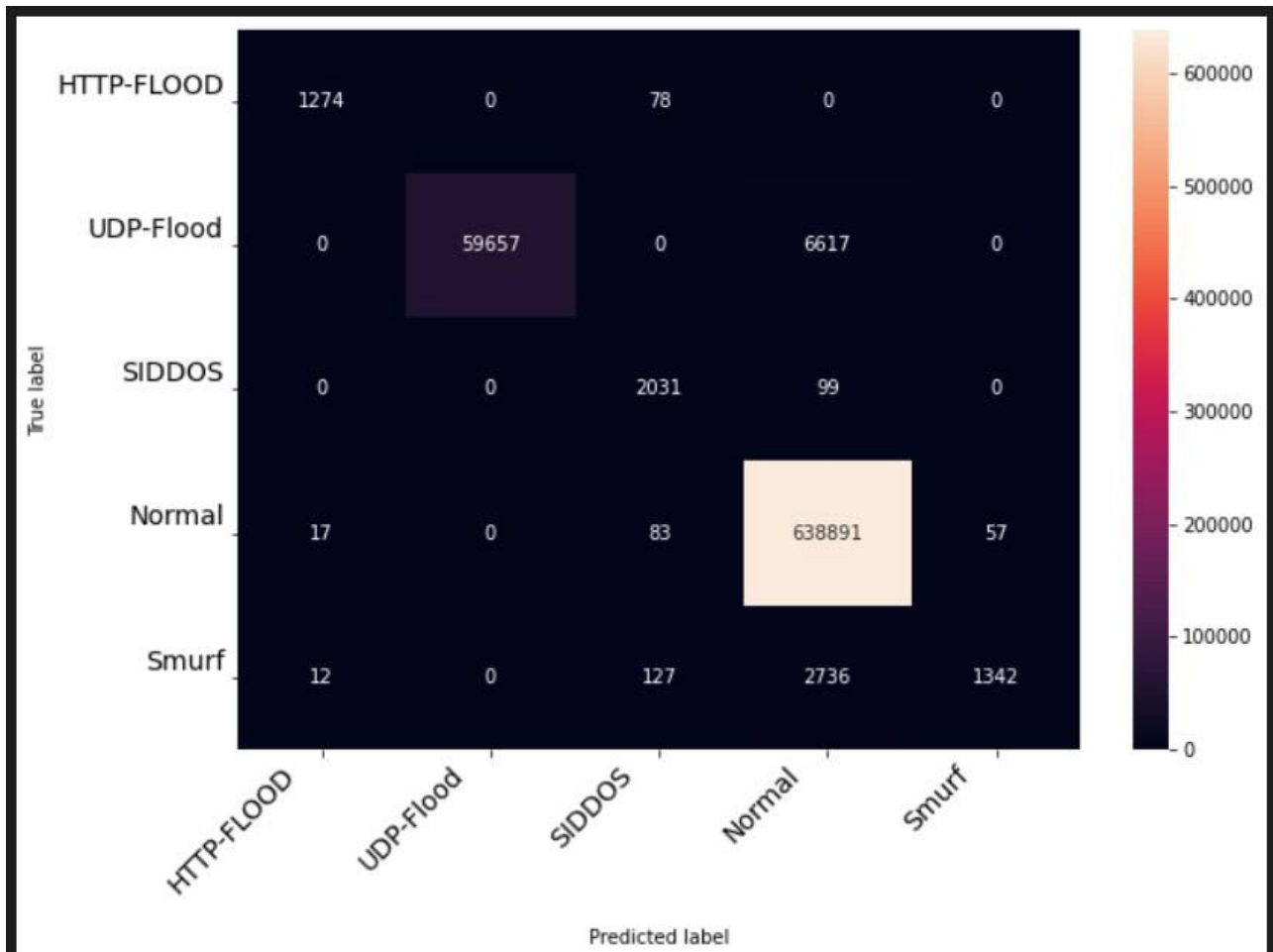


Fig. 14 Confusion matrix



V. CONCLUSION

Distributed Denial of Service (DDoS) attacks are a serious threat to the network security of various servers. With the ever expanding business activities on the internet, DDoS attack can cause billions of dollars in revenue for the organizations. They attacks have the capability of denying the services for the legitimate user by keeping the server busy by sending a flood of packets. In this paper, we have given a possible solution to help mitigate this problem. The novel ensemble model designed in this paper has an accuracy of 98.62%. The ensemble model has been developed using the four famous classifiers that is naïve bayes , stochastic gradient descent , multilayer perceptron and random forest. The data contains of total 2 million dataset and 28 features.

REFERENCES

- [1]. Pande, S., Khamparia, A., Gupta, D., Thanh, D.N.H. (2021), "DDoS Detection Using Machine Learning Technique. In: Khanna, A., Singh, A.K., Swaroop, A. (eds) Recent Studies on Computational Intelligence", Studies in Computational Intelligence, vol 921. Springer, Singapore.
- [2]. NG, B. A., & Selvakumar, S. (2019), "Deep radial intelligence with cumulative incarnation approach for detecting denial of service attacks", Neurocomputing.
- [3]. M. Wegmuller, J. P. von der Weid, P. Oberson, and N. Gisin, "High resolution fiber distributed measurements with coherent OFDR," in Proc. ECOC'00, 2000, paper 11.3.4, p. 109.
- [4]. Elejla, Omar E., Bahari Belaton, Mohammed Anbar, Basim Alabsi, and Ahmed K. Al-Ani, "Comparison of classification algorithms on ICMPv6-based DDoS attacks detection." In Computational Science and Technology, pp. 347-357, Springer, Singapore, 2019.
- [5]. S. Das, A. M. Mahfouz, D. Venugopal and S. Shiva, "DDoS Intrusion Detection Through Machine Learning Ensemble," 2019 IEEE 19th International Conference on Software Quality, Reliability and Security Companion (QRS-C), 2019, pp. 471-477.
- [6]. X. Yuan, C. Li and X. Li, "DeepDefense: Identifying DDoS Attack via Deep Learning," 2017 IEEE International Conference on Smart Computing (SMARTCOMP), 2017, pp. 1-8.
- [7]. Dataset used in this paper which contains 2 million rows and 28 features.[Online]. Available: <https://www.kaggle.com/datasets/jacobvs/ddos-attack-network-logs?resource=download>.
- [8]. Geeksforgeeks article describing naïve bayes classifier. [Online]. Available: <https://www.geeksforgeeks.org/naive-bayes-classifiers>.
- [9]. Definition of precision score, f1 score and recall score. [Online]. Available: <https://machinelearningmastery.com/precision-recall-and-f-measure-for-imbalanced-classification>.

Author Information

Nikhil Anand Mahendrakar, Sanjay A S, Manoj M
Siemens Healthineers, Bengaluru, India