



Cover Page



DOI: <http://ijmer.in.doi./2022/11.10.11>  
[www.ijmer.in](http://www.ijmer.in)

Digital Certificate of Publication: [www.ijmer.in/pdf/e-CertificateofPublication-IJMER.pdf](http://www.ijmer.in/pdf/e-CertificateofPublication-IJMER.pdf)

## DDOS ATTACK DETECTION USING SUPERVISED MACHINE LEARNING

Mohammad Edres Ansary and Preeti Sondhi

Universal Group of Institutions

Lalru, SAS Nagar, Mohali, Punjab, India

### Abstract

Websites can be common, either regularly visited by a large number of users or providing some useful information. It can then be accessed by a large number of users, often leading to an overload of servers, which can lead to an assault. Huge technological development is leading to many cases of hacking. Every day, numerous types of threats arise and their purpose is to make services inaccessible to the user. The DDoS attack can be determined using the amount of traffic they send per second to the host system, for example, small attacks could be measured by few megabits (Mbps), whereas it could be terabit per second in large attacks (Tbps). To execute the attacks efficiently, attackers use the botnet for large-volume attacks. If it occurs occasionally, the server will not respond to the legitimate user. The expert compares and records the incoming packet traffic with traffic signatures; network administrators secure internet devices from attack by using this tool.

**Keywords:** Attack, Ddos Attack, Machine Learning, Types of Ddos Attack, Types of Machine Learning Algorithms etc.

### Introduction

A Denial of Service (DoS) attack is an attempt by an attacker to render network resources by flooding the host of the service unresponsive to its legitimate users. A DoS attack originating from multiple sources is a distributed denial of service (DDoS) attack. In general, DoS attacks are initiated using an Internet connection from one computer or virtual machines, while DDoS attacks are initiated from several different compromised computers, virtual machines, to overwhelm victims' networks. DDoS is performed by simultaneously sending a large number of requests via botnets and compromised IoT devices to the target's exhaust computing resource (Bandwidth and Traffic). The compromised computers, also referred to as bot or zombie, run remotely under the supervision of one or more of the bot-masters and attack groups of bots (botnet) as in Figure 1. Bots can be either malicious users who are preparing for an attack or legitimate users who are infected.

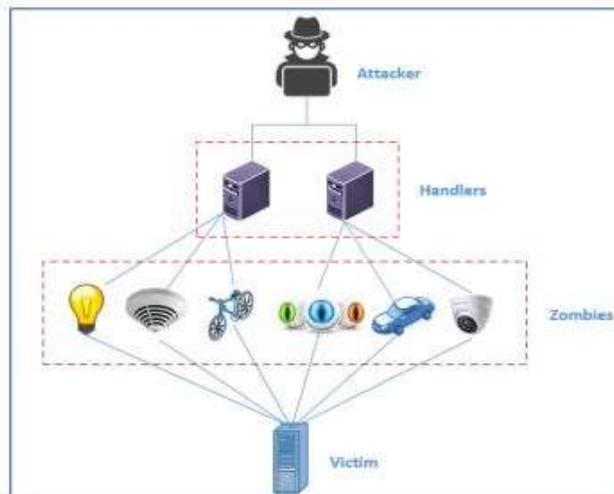


Fig. 1: DDOS attack network infrastructure

### Types of DDOS attacks

This malicious attack is done with several compromised services globally. There are three major types of DDoS attacks. They are,

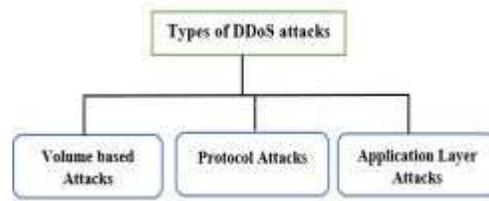


Fig. 2: DDoS attack types

**Volume based Attacks:** Also today, it is one of the most common attacks. This attack sends out massive quantities of data from the bandwidth of the available network. This leads to a congestion that leads to an acceleration that requests to send a large amount of traffic to the services of the victim that will affect the responses and result in the server being clogged. **Protocol Attacks:** Protocol attacks are also known as state exhaustion attacks that impact the victim's system's network layer. This interferes with the capability of table spaces, such as firewalls, and load balancers that forward requests to the target.

**Application Layer Attacks:** Interruption of the website or service of the victim with a large number of requests before the exhaustive stage of the request submitted is reached. The requests may include the downloading of large data or the database requesting the data. If the target is targeted with millions of requests at a time, the device process can slow down and can even be automatically shut down.

### Challenges with DDoS Attacks

DoS attacks are complex in nature, so the exact IP address from which the attack originated cannot be identified. It is possible to spoof its source and attack more than one system(node) as a source over the internet from various locations. Often, its source is compromised to launch attacks, too. Since it comes from multiple sources and needs human involvement, the type of attack is very complex. In 2017, when any device is affected by a DDoS attack, there is a lack of agile granular responses that can increase financial processing power issues.

### Machine Learning

Machine learning is defined as the ability of a system to learn and progress on the basis of learning without comprehensive programming. The methods of Machine Learning help to efficiently solve different problems in the IT field. It is now in vogue and helps to solve big network security issues. Using various learning strategies, it is narrowly categorised.

### Classification of Machine Learning Algorithms

Every task in Machine learning is broadly classified into several categories. They are

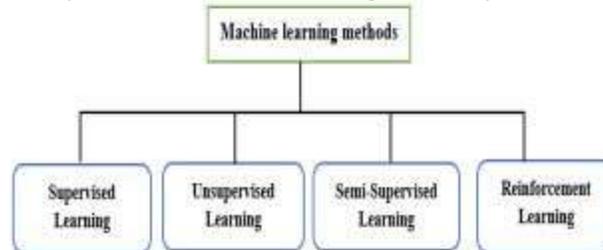


Fig. 3: Types of machine learning

**Supervised Learning Techniques:** The aim is to collect data and learn from it using various Machine learning based supervised techniques. It is a mathematical model consisting of set of inputs (labelled inputs) and desired outputs used for predictive modelling. Some commonly used algorithms are Nearest Neighbor, Decision tree, Support vector machines, Naïve Bayes, linear regression.

**Unsupervised Learning Techniques:** In this method, the data are unlabeled, the system trains itself and produces the output. This helps in picking the important data that is required for the analyses. When we require any information regarding the relationship of a data, we use unsupervised learning method. There are various clustering algorithms like association rules, k-means clustering.



Cover Page



**Semi-supervised learning Techniques:** It is a combination of supervised and unsupervised learning methods where we will use both labelled and unlabeled data. It produces a desired result having important parameters required for analyses.

**Reinforcement Learning Techniques:** Reinforcement training is based on trial and error method for a particular decision. It gains experiences from the previous trainings and gives accurate knowledge based on the response received.

### Techniques used

#### Logistic Regression

Logistic regression is one of the most popular Machine Learning algorithms, which comes under the Supervised Learning technique. It is used for predicting the categorical dependent variable using a given set of independent variables. Logistic regression predicts the output of a categorical dependent variable. Therefore, the outcome must be a categorical or discrete value. It can be either Yes or No, 0 or 1, true or False, etc. but instead of giving the exact value as 0 and 1, **it gives the probabilistic values which lie between 0 and 1.** Logistic Regression is much similar to the Linear Regression except that how they are used. Linear Regression is used for solving Regression problems, whereas **Logistic regression is used for solving the classification problems.** In Logistic regression, instead of fitting a regression line, we fit an "S" shaped logistic function, which predicts two maximum values (0 or 1). The curve from the logistic function indicates the likelihood of something such as whether the cells are cancerous or not, a mouse is obese or not based on its weight, etc. Logistic Regression is a significant machine learning algorithm because it has the ability to provide probabilities and classify new data using continuous and discrete datasets. Logistic Regression can be used to classify the observations using different types of data and can easily determine the most effective variables used for the classification.

#### Decision Tree

Decision Tree is a supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome.

In a Decision tree, there are two nodes, which are the Decision Node and Leaf Node. Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.

The decisions or the test are performed on the basis of features of the given dataset.

It is a graphical representation for getting all the possible solutions to a problem/decision based on given conditions.

It is called a decision tree because, similar to a tree, it starts with the root node, which expands on further branches and constructs a tree-like structure.

In order to build a tree, we use the CART algorithm, which stands for Classification and Regression Tree algorithm.

A decision tree simply asks a question, and based on the answer (Yes/No), it further split the tree into subtrees.

Below diagram explains the general structure of a decision tree:

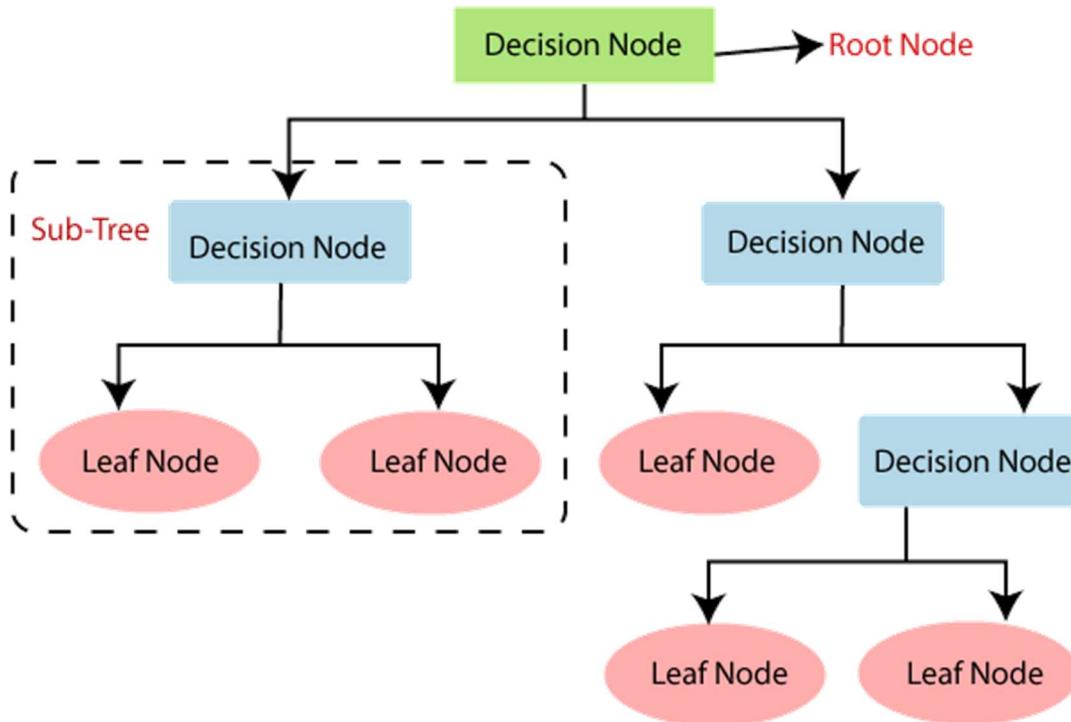


Fig. 4: Decision Tree

### Why use Decision Trees?

There are various algorithms in Machine learning, so choosing the best algorithm for the given dataset and problem is the main point to remember while creating a machine learning model. Below are the two reasons for using the Decision tree:

- Decision Trees usually mimic human thinking ability while making a decision, so it is easy to understand.
- The logic behind the decision tree can be easily understood because it shows a tree-like structure.

### Decision Tree Terminologies

**Root Node:** Root node is from where the decision tree starts. It represents the entire dataset, which further gets divided into two or more homogeneous sets.

**Leaf Node:** Leaf nodes are the final output node, and the tree cannot be segregated further after getting a leaf node.

**Splitting:** Splitting is the process of dividing the decision node/root node into sub-nodes according to the given conditions.

**Branch/Sub Tree:** A tree formed by splitting the tree.

**Pruning:** Pruning is the process of removing the unwanted branches from the tree.

**Parent/Child node:** The root node of the tree is called the parent node, and other nodes are called the child nodes.

### How does the Decision Tree algorithm Work?

In a decision tree, for predicting the class of the given dataset, the algorithm starts from the root node of the tree. This algorithm compares the values of root attribute with the record (real dataset) attribute and, based on the comparison, follows the branch and jumps to the next node.

For the next node, the algorithm again compares the attribute value with the other sub-nodes and move further. It continues the process until it reaches the leaf node of the tree. The complete process can be better understood using the below algorithm: 5

- **Step-1:** Begin the tree with the root node, says S, which contains the complete dataset.
- **Step-2:** Find the best attribute in the dataset using **Attribute Selection Measure (ASM)**.
- **Step-3:** Divide the S into subsets that contains possible values for the best attributes.
- **Step-4:** Generate the decision tree node, which contains the best attribute.

- **Step-5:** Recursively make new decision trees using the subsets of the dataset created in step -3. Continue this process until a stage is reached where you cannot further classify the nodes and called the final node as a leaf node.

### K Neighbours Classifier

- K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique.
- K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.
- K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suite category by using K- NN algorithm.
- K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for the Classification problems.
- K-NN is a **non-parametric algorithm**, which means it does not make any assumption on underlying data.
- It is also called a **lazy learner algorithm** because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset.
- KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much similar to the new data.

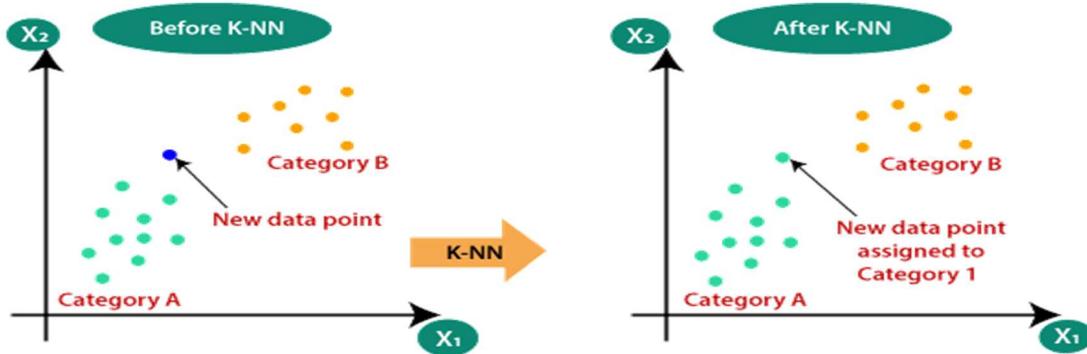
**Example:** Suppose, we have an image of a creature that looks similar to cat and dog, but we want to know either it is a cat or dog. So for this identification, we can use the KNN algorithm, as it works on a similarity measure. Our KNN model will find the similar features of the new data set to the cats and dogs images and based on the most similar features it will put it in either cat or dog category.



Fig. 5: KNN classifier

### Why do we need a K-NN Algorithm?

Suppose there are two categories, i.e., Category A and Category B, and we have a new data point  $x_1$ , so this data point will lie in which of these categories. To solve this type of problem, we need a K-NN algorithm. With the help of K-NN, we can easily identify the category or class of a particular dataset. Consider the below diagram:

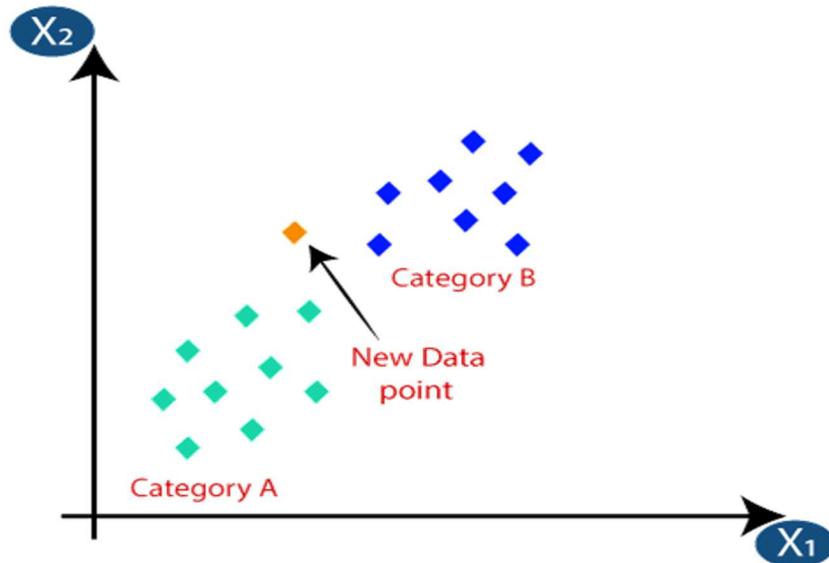


How does K-NN work?

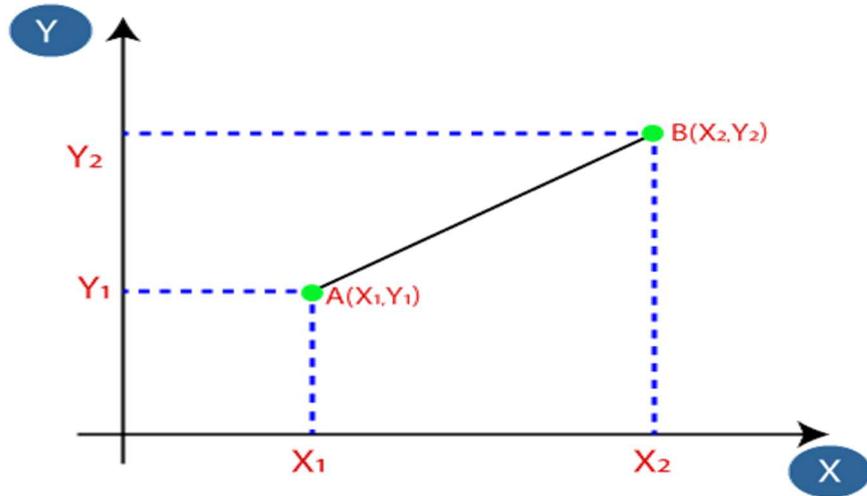
The K-NN working can be explained on the basis of the below algorithm:

- **Step-1:** Select the number K of the neighbors
- **Step-2:** Calculate the Euclidean distance of **K number of neighbors**
- **Step-3:** Take the K nearest neighbors as per the calculated Euclidean distance.
- **Step-4:** Among these k neighbors, count the number of the data points in each category.
- **Step-5:** Assign the new data points to that category for which the number of the neighbor is maximum.
- **Step-6:** Our model is ready.

Suppose we have a new data point and we need to put it in the required category. Consider the below image:

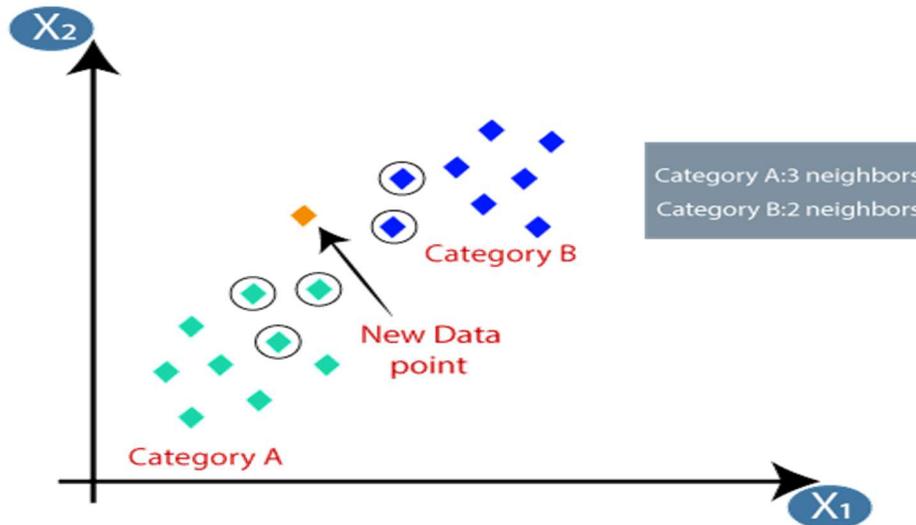


- Firstly, we will choose the number of neighbors, so we will choose the  $k=5$ .
- Next, we will calculate the **Euclidean distance** between the data points. The Euclidean distance is the distance between two points, which we have already studied in geometry. It can be calculated as:



$$\text{Euclidean Distance between } A_1 \text{ and } B_2 = \sqrt{(X_2 - X_1)^2 + (Y_2 - Y_1)^2}$$

- By calculating the Euclidean distance we got the nearest neighbors, as three nearest neighbors in category A and two nearest neighbors in category B. Consider the below image:



- As we can see the 3 nearest neighbours are from category A, hence this new data point must belong to category A.

**Advantages of KNN Algorithm**

- It is simple to implement.
- It is robust to the noisy training data
- It can be more effective if the training data is large.

**Disadvantages of KNN Algorithm**

- Always needs to determine the value of K which may be complex some time.
- The computation cost is high because of calculating the distance between the data points for all the training samples

**MLP Classifier**

Multi-Layer perceptron defines the most complex architecture of artificial neural networks. It is substantially formed from multiple layers of the perceptron. TensorFlow is a very popular deep learning framework released by, and this notebook will guide to

build a neural network with this library. If we want to understand what is a Multi-layer perceptron, we have to develop a multi-layer perceptron from scratch using Numpy. The pictorial representation of multi-layer perceptron learning is as shown below-

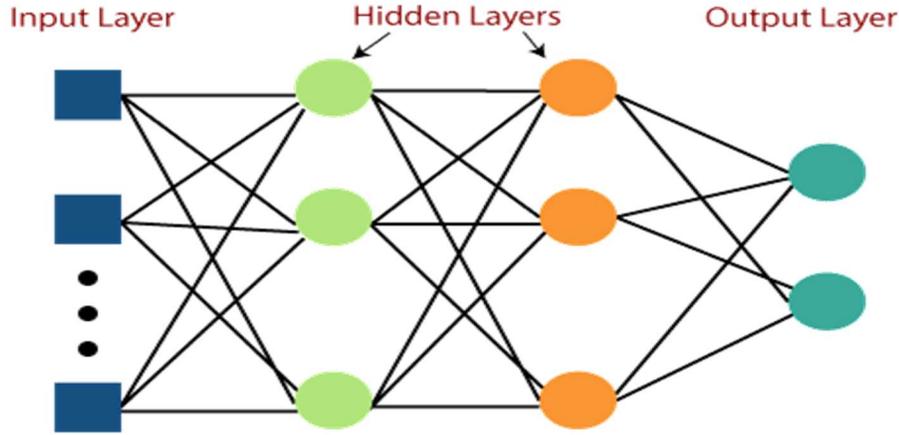


Fig. 6: multilayer perceptron

MLP networks are used for supervised learning format. A typical learning algorithm for MLP networks is also called back propagation's algorithm.

A multilayer perceptron (MLP) is a feed forward artificial neural network that generates a set of outputs from a set of inputs. An MLP is characterized by several layers of input nodes connected as a directed graph between the input nodes connected as a directed graph between the input and output layers. MLP uses backpropagation for training the network. MLP is a deep learning method.

### Proposed work

In all aspects of business and industry, including bank transfers, email, social media and university services, etc., network security has become of paramount importance. Web and network networks have been suffering from hacker attacks recently. New types of Distributed Denial of Service (DDoS) that operate on the application layer as well as the network layer are continually being created by hackers. In the above listed areas, the vulnerabilities allow hackers to refuse access to web services and slow down access to resources from the network. Machine learning is used to detect and identify network traffic based on certain features that are used to calculate and evaluate whether network traffic is regular or is a type of DDoS (average packet size, inter arrival time, packet size, packet rate, bit rate, etc.). For the most part, DDoS attacks have the same average packet size. Instead of the usual packet, the number of packets will increase in the attacked packet; the inter-arrival time will also be too small to allow attackers to quickly absorb resources. DDoS packets for network layer attacks often have a high bit rate. Attackers concentrate on any characteristics that allow them to absorb assets and make the service inaccessible to end users.

### Result and discussion

1)UDP Flood 2) TCP SYN Flood 3) ICMP Flood

1.) UDP FLOOD

Import libraries and read csv file

```
In [1]: import numpy as np
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt

In [3]: df = pd.read_csv(r"revised_kddcup_dataset.csv", index_col=0)

In [4]: df.head()

Out[4]:
```

	duration	protocol_type	service	flag	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	...	dst_host_srv_count	dst_host_same_srv_rate	dst_host_diff_srv_rate	dst_host_s...
0	0	udp	private	SF	105	146	0	0	0	0	...	254	1.0	0.01	
1	0	udp	private	SF	105	146	0	0	0	0	...	254	1.0	0.01	
2	0	udp	private	SF	105	146	0	0	0	0	...	254	1.0	0.01	



Cover Page



DOI: <http://ijmer.in.doi/2022/11.10.11>  
[www.ijmer.in](http://www.ijmer.in)

Digital Certificate of Publication: [www.ijmer.in/pdf/e-CertificateofPublication-IJMER.pdf](http://www.ijmer.in/pdf/e-CertificateofPublication-IJMER.pdf)

### Import, train and test models

```
In [21]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42, test_size=0.3)

In [22]: from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.neural_network import MLPClassifier
from sklearn.tree import DecisionTreeClassifier

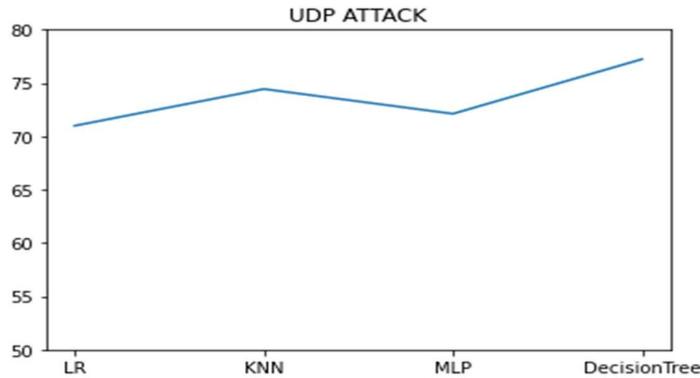
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score

In [26]: models = [LogisticRegression(), KNeighborsClassifier(n_neighbors=3),MLPClassifier(alpha=0.005),DecisionTreeClassifier()]
classifiers = ["LR", "KNN","MLP", "DecisionTree"]
scores = []

In [27]: for model in models:
model.fit(X_train,y_train)
y_pred = model.predict(X_test)
score = accuracy_score(y_test, y_pred)*100
scores.append(score)
print("Accuracy of the model is: ", score)
conf_matrix = confusion_matrix(y_test,y_pred)
report = classification_report(y_test,y_pred)
print("Confusion Matrix:\n",conf_matrix)
print("Report:\n",report)
print("\n*****")
```

### Models accuracy score

```
*****
In [28]: plt.plot(classifiers,scores)
plt.title("UDP ATTACK")
plt.ylim(50,80)
plt.show()
```



2. TCP SYN flood  
Model accuracy score



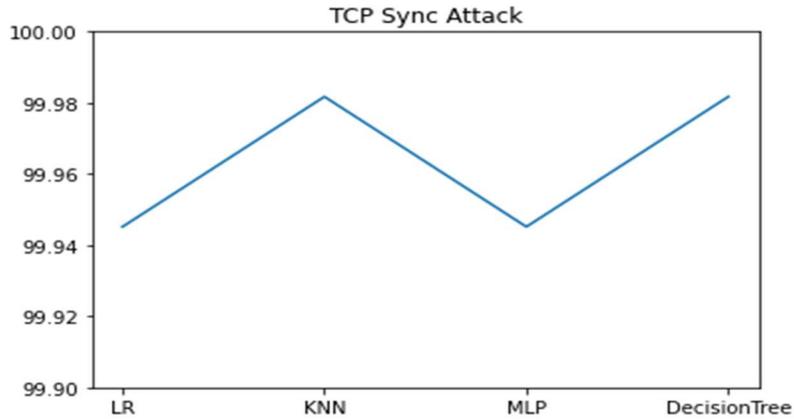
Cover Page



DOI: <http://ijmer.in.doi/2022/11.10.11>  
[www.ijmer.in](http://www.ijmer.in)

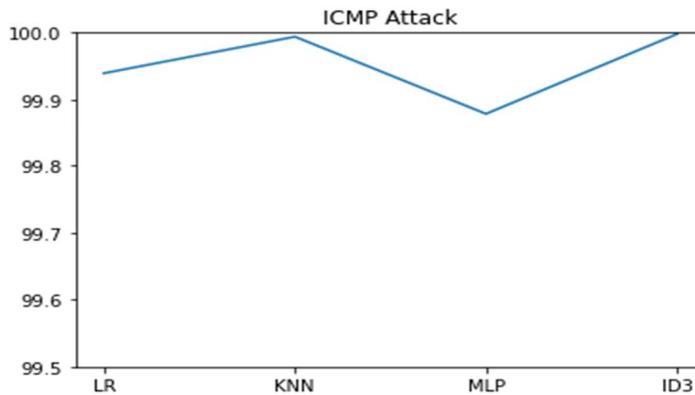
Digital Certificate of Publication: [www.ijmer.in/pdf/e-CertificateofPublication-IJMER.pdf](http://www.ijmer.in/pdf/e-CertificateofPublication-IJMER.pdf)

```
In [28]: plt.plot(classifiers,scores)
plt.title("TCP Sync Attack")
plt.ylim(99.9,100)
plt.show()
```



### 3. ICMP flood

```
In [37]: plt.plot(classifiers,scores)
plt.title("ICMP Attack")
plt.ylim(99.5,100)
plt.show()
```



### Conclusion

Attacks are the damage and disruption of a system's usual behaviour caused by the abuse of vulnerabilities by various methods and techniques. Attacks come with distinct motivations in various ways. An aggressive attack that tracks un-encrypted network traffic in order to find confidential information is one type of attack. The passive attack, which monitors weakly encrypted traffic to find authentication data, is another form of attack. Control attacks, physical attacks, distributed denial of service attacks, privacy attacks such as password base attacks, cyber espionage and eavesdropping are the most common attacks.



Cover Page



## References

1. N. Ravi and S. M. Shalinie, "Learning-Driven Detection and Mitigation of DDoS Attack in IoT via SDN-Cloud Architecture," in *IEEE Internet of Things Journal*, vol. 7, no. 4, pp. 3559-3570, April 2020, doi: 10.1109/JIOT.2020.2973176.
2. K. Huang, L. Yang, X. Yang, Y. Xiang and Y. Y. Tang, "A Low-Cost Distributed Denial-of-Service Attack Architecture," in *IEEE Access*, vol. 8, pp. 42111-42119, 2020, doi: 10.1109/ACCESS.2020.2977112.
3. S. Velliangiri, P. Karthikeyan & V. Vinoth Kumar (2020) Detection of distributed denial of service attack in cloud computing using the optimization-based deep networks, *Journal of Experimental & Theoretical Artificial Intelligence*, DOI: 10.1080/0952813X.2020.1744196
4. Asad, Mohammad & Khrais, Rami & Yateem, A.Rahman. (2020). DoS and DDoS Attack Detection Using Deep Learning and IDS. *International Arab Journal of Information Technology*. 17. 655-661. 10.34028/iajit/17/4A/10.
5. Dwivedi, Shubhra & Vardhan, Manu & Tripathi, Sarsij. (2020). Defense against distributed DoS attack detection by using intelligent evolutionary algorithm. *International Journal of Computers and Applications*. 1-11. 10.1080/1206212X.2020.1720951.
6. Tuan, N.N.; Hung, P.H.; Nghia, N.D.; Tho, N.V.; Phan, T.V.; Thanh, N.H. A DDoS Attack Mitigation Scheme in ISP Networks Using Machine Learning Based on SDN. *Electronics* **2020**, *9*, 413.
7. T. V. Phan and M. Park, "Efficient Distributed Denial-of-Service Attack Defense in SDN-Based Cloud," in *IEEE Access*, vol. 7, pp. 18701-18714, 2019, doi: 10.1109/ACCESS.2019.2896783.
8. Khalaf, Bashar & Mostafa, Salama & Mohammed, Mazin & Abdullaha, Wafaa & Mustapha, Aida. (2019). Comprehensive Review of Artificial Intelligence and Statistical Approaches in Distributed Denial of Service Attack and Defense Methods. *IEEE Access*. PP. 2169-3536. 10.1109/ACCESS.2019.2908998.
9. Sahoo, K.S., Panda, S.K., Sahoo, S. et al. Toward secure software-defined networks against distributed denial of service attack. *JSupercomput* 75, 4829–4874 (2019).
10. <https://doi.org/10.1007/s11227-019-02767-z>
11. Saritha et al. "Prediction of DDoS Attacks using Machine Learning and Deep Learning Algorithms." (2019).