

DISTRIBUTED DENIAL OF SERVICE ATTACKS CLASSIFICATION AND MITIGATION



KINZA TAMEEZ
01-249191-004

Master of Science in Data science

Supervisor: Dr. Kashif Naseer Qureshi

*Department of Computer Science
Bahria University, Islamabad*

MS-13
Thesis Completion Certificate

Student Name: Kinza Tameez Registration Number:
Program of Study: Masters of Science in Data Science

***Thesis Title: DISTRIBUTED DENIAL OF SERVICE ATTACKS
CLASSIFICATION
AND MITIGATION***

It is to certify that the above student's thesis has been completed to my satisfaction and, to my belief, its standard is appropriate for submission for evaluation. I have also conducted a plagiarism test of this thesis using HEC prescribed software and found a similarity index at 8% that is within the permissible set by the HEC. for MS/MPhil/Ph.D.

I have also found the thesis in a format recognized by the BU for MS/MPhil/Ph.D. thesis.



Principle Supervisor's Signature: _____

Principle Supervisor's Name: Dr. Kashif Naseer Qureshi
September 27, 2021

MS-14A
Author's Declaration

I, Kinza Tameez hereby state that my MS thesis titled "DISTRIBUTED DENIAL OF SERVICE ATTACKS CLASSIFICATION AND MITIGATION" is my work and has not been submitted previously by me for taking any degree from "Bahria University, Islamabad" or anywhere else in the country/world.

At any time if my statement is found to be incorrect even after my Graduate the university has the right to withdraw cancel my MS degree.

KINZA TAMEEZ
01-249191-004
September 27, 2021

MS-14B

Plagiarism Undertaking

I, Kinza Tameez solemnly declare that research work presented in the thesis titled **DISTRIBUTED DENIAL OF SERVICE ATTACKS CLASSIFICATION AND MITIGATION** is solely my research work with no significant contribution from any other person. Small contribution/help whenever taken has been duly acknowledged and that complete thesis has been written by me.

I understand the zero-tolerance policy of Bahria University and the Higher Education Commission of Pakistan towards plagiarism. Therefore, I as an author of the above-titled thesis declare that no portion of my thesis has been plagiarized and any material used is properly referred to / cited.

I undertake that if I am found guilty of any formal plagiarism in the above-titled thesis even after awarding of MS degree, the university reserves the right to withdraw/revoke my MS degree and HEC, and the university has the right to publish my name on HEC/University Website on which names of students who submitted plagiarized thesis are placed.

KINZA TAMEEZ
01-249191-004
September 27, 2021

ABSTRACT

Distributed Denial of Service (DDoS) attack introduces new threats and vulnerabilities to the network security, in addition to the ones that already existed. The DDoS attack type is one of the most violent attack type in current years, causing chaos on the entire network system. Most of the research is limited to the classification of DDoS attacks by Machine learning algorithms. Deep learning is not involved in the classification of DDoS attacks and the classes on which classification is applied are not enough. Another challenge is detecting and mitigating DDoS attacks which are not effective and lead to different errors and degraded the accuracy. The existing solutions have adopted old datasets whereas the new DDoS attacks have changed their patterns and types. The proposed model aims to classifying the different classes of DDoS attacks on the new dataset by using deep learning models. The proposed classification classifies the newly released dataset of CICDDOS 2019. The dataset contains complete and current DDoS attack types. The evaluation of the proposed model showed that Long Short Term Memory (LSTM) which is a modification of Recurrent Neural Network (RNN) for the classification of DDoS attacks gives the highest accuracy as compared to SVM (Support Vector Machine). The proposed model achieves better results in terms of accuracy and false positive rate.

ACKNOWLEDGMENTS

I am very grateful to Allah Almighty for enabling me to complete this project. I would like to thank my supervisor Dr. Kashif Naseer Qureshi, for his guidance, supervision, and encouragement. Without his continued support and interest, this thesis would not have been the same as presented here. I would like to avail this opportunity and express my sincere regards to all the faculty members of our department for their support. My sincere appreciation also extends to all my colleagues and others who have assisted on various occasions. Their views and tips are useful indeed. Unfortunately, it is not possible to list all of them in this limited space. I am grateful to all my family members.

KINZA TAMEEZ
Islamabad, Pakistan
September 27, 2021

Table of Contents

ABSTRACT	vii
CHAPTER 1	1
INTRODUCTION	1
1.1 Overview	1
1.2 Introduction	1
1.3 Intrusion Detection System in Cyber Security	2
1.4 Problem Background	3
1.5 Problem Statement	3
1.6 Research Questions	3
1.7 Research Objectives	4
1.8 Research Scope and Limitation	4
1.9 Significance of the Study	4
1.10 Thesis organization	4
CHAPTER 2	5
LITERATURE REVIEW	5
2.1 Overview	5
2.2 Intrusion Detection System	5
2.3 Types of Attacks	6
2.3.1 Types of DDoS Attacks	6
2.3.2 Drive-by Attack	11
2.3.3 Password Attack	11
2.3.4 SQL injection attack	12
2.3.5 Cross-site Scripting (XSS) Attack	13
2.3.6 Eavesdropping Attack	13
2.3.7 Birthday Attack	14
2.3.8 Malware Attack	14
2.4 Solutions	18
2.4.1 Naïve Bayes	18
2.4.2 Decision tree	18
2.4.3 Random Forest	19
2.4.4 Support Vector Machine (SVM)	20

2.4.5	K-Nearest Neighbor (KNN).....	20
2.4.6	Multilayer Perceptron	21
2.4.7	Ensemble classifier	21
2.5	Background Research.....	22
CHAPTER 3	29
METHODOLOGY	29
3.1	Overview	29
3.2	Research Methodology Framework	29
3.3	Confusion Matrix	31
3.3.1	Accuracy Method.....	32
3.3.2	Additional Metrics	32
CHAPTER 4	34
PROPOSED WORK DESIGN	34
4.1	Overview	34
4.2	Recurrent Neural Networks.....	35
4.3	LSTM.....	36
4.4	System Overview	36
4.4.1	Dataset in the Proposed Methodology.....	37
4.4.2	Training Dataset.....	39
4.4.3	Testing Dataset.....	39
4.4.4	Classification.....	40
CHAPTER 5	41
RESULTS	41
5.1	Experiment and Results.....	41
5.2	Environment and Tools	41
5.3	Optimizer.....	41
5.4	Training Model.....	41
5.5	Hyper-Parameter Tuning.....	42
5.6	Results	43
5.6.1	Result of Select K-Best.....	44
5.6.2	Accuracy Analysis	45
5.6.3	Training Accuracy	45
5.6.4	Testing Accuracy	45

CHAPTER 6.....	48
CONCLUSION.....	48
Future Work.....	48
References.....	49

List of Tables

Table 2.1: Comparison Table.....	26
Table 5.1: Performance metrics with different learning rates.....	43
Table 5.2: Results of balanced and unbalanced datasets	44
Table 5.3: Confusion Matrix	47

List of Figures

Figure 2.1: Architecture of DDoS Attacks	2
Figure 2.1: Types of DDoS Attacks	6
Figure 2.2: Decision Tree	19
Figure 2.3: SVM Classifier.....	20
Figure 3.1: CICCDDOS 2019 Dataset Distribution	31
Figure 4.1: Simple Neural Network	34
Figure 4.2: RNN Architecture	35
Figure 4.3: LSTM Architecture	36
Figure 4.4: System Overview	37
Figure 4.5: Feature Selection Method	39
Figure 5.1: LSTM Loss	42
Figure 5.2: Feature analysis using Select k Best Method.....	45
Figure 5.3: Accuracy of Training and Testing Dataset	46
Figure 5.4: ROC curve.....	46

CHAPTER 1

INTRODUCTION

1.1 Overview

This chapter presents the detail information of cybersecurity attacks, detection and prevention techniques. The research questions, objectives and problem background also discuss to achieve the research goal.

1.2 Introduction

After the evolution of the internet, the most important factor is network security [1]. Protecting computers, networks, electronic systems, and data from harmful attacks is called cybersecurity. It is also known as Information Technology Security. Cyber threats grow rapidly across the globe by breaching a large number of data each year. The term cybersecurity is used to state the security presented through online services to protect online digital data or information. With the increase of internet users, the security threats have increased and that can cause massive destruction. Cybersecurity is essential because it protects data from threats such as data theft or misuse, as well as a computer from different viruses.

According to a report in 2019, 7.9 billion records are exposed by data breaches and that is double in the year 2018 [2]. In cybersecurity, it is one of the most difficult problems to combat Distributed Denial of Services attacks (DDoS). In 2018 [2], GitHub is suffered from the largest and the major DDoS attack, which is around 1.35-terabit-per-second attack contrary to the site. This is the biggest DDoS attack on GitHub in history. DDoS attacks are currently the most prevalent and exponentially growing threats for any organization and are difficult to control [3]. Detection of DDoS attacks is compulsory to provide secure communication and services over the network. The DDoS attack is deliberately executed by the user to disturb and damage the services of the user. This type of attack is attempted on the network resources and servers by using some assigned computers and initiate overflowing of messages and connection requests and leads to the Denial of Service (DoS) to the user [4].

Nowadays DDoS attacks are dynamic and are discharged in diversity and make patterns, which it makes difficult for the static solutions to detect dynamic attacks [5]. Figure 1.1 shows the DDoS attacks architecture.

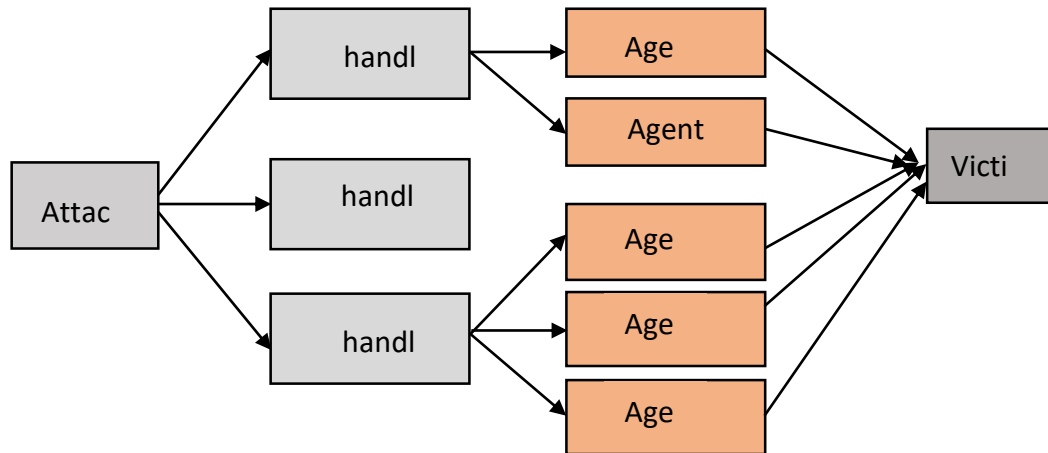


Figure 1.1: Architecture of DDoS Attacks

1.3 Intrusion Detection System in Cyber Security

An Intrusion Detection System (IDS) monitors network traffic for unusual behavior and sends out alerts when it is found. It's a software program that scans a network system for malicious activity, policy violations, and virus threats. Normally, any malicious activity or violation is recorded and reported to a system administrator. Detection of Intrusion plays an important role in cybersecurity. In many last years, numerous IDS are based on open and commercial signatures which have been proposed by the scientist. This gave rise to the detection of IDS which are based on anomaly detection [6]. Protection of the system from DDoS attacks is necessary as it corrupts the system. Detecting DDOS attacks at the application layer is difficult because of its Distributed Nature.

For availability of protection of services, safety and security of other communication services and resources, many different approaches and techniques are used to detect DDoS attacks. In [7], authors used the application of an artificial neural network for the detection of a DDoS attacks. DDoS are currently the most prevailing and rapidly growing threats for any organization and difficult to control [8]. The DDoS attack is deliberately executed by the user to disturb and damage the user services. This type of attack on network resources and servers is targeted by some assigned computers resulting in the overflowing of messages and connection requests that makes services unavailable for the user [9]. Nowadays DDoS attacks

are dynamic and are discharged in diversity and make patterns, which it makes difficult for the static solutions to detect dynamic attacks [10]. Another technique is proposed in [11], for the detection of DDoS attack flow monitoring services in a Software-Defined Network (SDN). DDoS attacks make online services complicated and sometimes disconnect and unavailable.

1.4 Problem Background

Problem background is the process of finding and investigating the problems in the research domain. Most of the research is limited to the classification of DDoS attacks by using machine learning algorithms. Deep learning does not involve in the classification of DDoS attacks and the classes on which classification is applied are not enough. The other challenge in the research is that detecting and mitigating DDoS attacks are not effective due to which errors occur and it does not give better accuracy and also the dataset used in previous research are old datasets. In the previous models, dataset collection has been always one of the biggest challenge because of unavailability of dataset. In the dataset of intrusion detection, privacy and legal issues are the important things for where the data is not easily available. Although, various researchers have defined their datasets but those are not feasible or over fit for the classification. Accuracy is one of another challenge in the previous researches. Because of the unstable and unclassified data, the model does not give better accuracy. While investigating the problem we came across to know that in the previous researches the classification is on the limited classes.

1.5 Problem Statement

The existing models are based on old datasets which are not feasible for IDS systems due to its unclassified type. The accuracy is another challenge in existing schemes which need attention due to rapidly increased DDoS attacks on digital data.

1.6 Research Questions

The research questions are as follows:

1. How to classifying the attacks?
2. How to recognize the different attacks classes?
3. How to improve the accuracy and classification feature to detect the DDoS attacks?

1.7 Research Objectives

The proposed research is based on different classes of DDoS attacks on the new dataset. The primary goal of the research is better classification of DDoS attacks by using deep learning models. The objectives of this research are as follows:

1. To design a monitoring mechanism for DDoS attacks.
2. To control and mitigate the attacks to protect the data, networks, and devices.
3. To improve the accuracy by adopting the better techniques and approaches.

1.8 Research Scope and Limitation

Mitigation and detection of different DDoS attacks are always deeply investigate because of their severity. The scope of this research is to classifying and mitigating DDoS attacks to prevent the computer devices and networks which are affecting by these attacks.

1.9 Significance of the Study

The main purpose of this study is to classifying the real DDoS attacks to mitigate them to avoid data loss. This classification enables detection, prevention, and mitigation of DDoS attacks before they attack data.

1.10 Thesis organization

The thesis is divided into 5 chapters. Chapter 2 covers the DDoS attacks types, and their classification techniques, background and literature work also presents in this chapter to find the problem background and research gap. Chapter 3 discusses the technical details of the proposed research methodology. Chapter 4 discusses the proposed work design and development. Chapter 5 discusses the experimental setup with results and discussion based on extracted results. Chapter 6 concludes the research and explain the future work.

CHAPTER 2

LITERATURE REVIEW

2.1 Overview

In this chapter, the discussion about the basic concept of Intrusion Detection Systems (IDS), its detection IDs, types, and attacks. The existing solutions, their achievements, and limitations are also discussing. The discussion and comparison section also presents the findings and conclusion of existing IDS solutions and motivation to design or proposed a new system.

2.2 Intrusion Detection System

In the present era of the internet and technology, identifying malicious activities in the network or a system is the greatest challenge. Cyber-attacks, targeting the information system and computer networks. On average, the system and network flaws or design errors made by human causes these cyber-attacks. To defend and protect the computer networks against these attacks, an Intrusion Detection System (IDS) has been proposed in [12]. An IDS is hardware or can be a software that inactively or actively identifies the intrusions which are caused by attackers and controls the network [13]. IDS is a type of alarm that monitors unusual activities [14]. It can be defined as computerization of intrusion detection process, a process of finding violation events of security policies [14], and knows that how to read and understand the contents of files from routers, firewalls, servers, and other network devices [15]. In addition, identifying the security incidents, IDS calls for specific properties which act as an inactive countermeasure. It monitors as a whole or a part of networks and also targets detection rate of attack and low false alarm rate [16].

2.3 Types of Attacks

A Denial-of-Service (DoS) attack overburdened the resources of the system and ignore the service requests and makes the system unavailable. Distributed Denial-of-Service (DDoS) and DoS attacks are almost the same but in DDoS attacks, they have launched more than one host computer or hosts can be many which are infected with malicious software that is monitored by an attacker. DoS attacks, unlike those that are intended to allow the attacker to gain or improve access, do not give direct benefits to the attacker. For some of them, only the service denial is enough to satisfy. If the attacker is attacking any business component, then his benefit would be substantial. Another use of a DDoS attack is when another type of attack is to be launched it makes the system offline. Like session hijacking. Some of the different types of DoS and DDoS attacks such as TCP SYN flood attack, smurf attack, ping-of-death attack, teardrop attack, and botnets. There are different types of DDoS attack as shows in Figure 2.1.

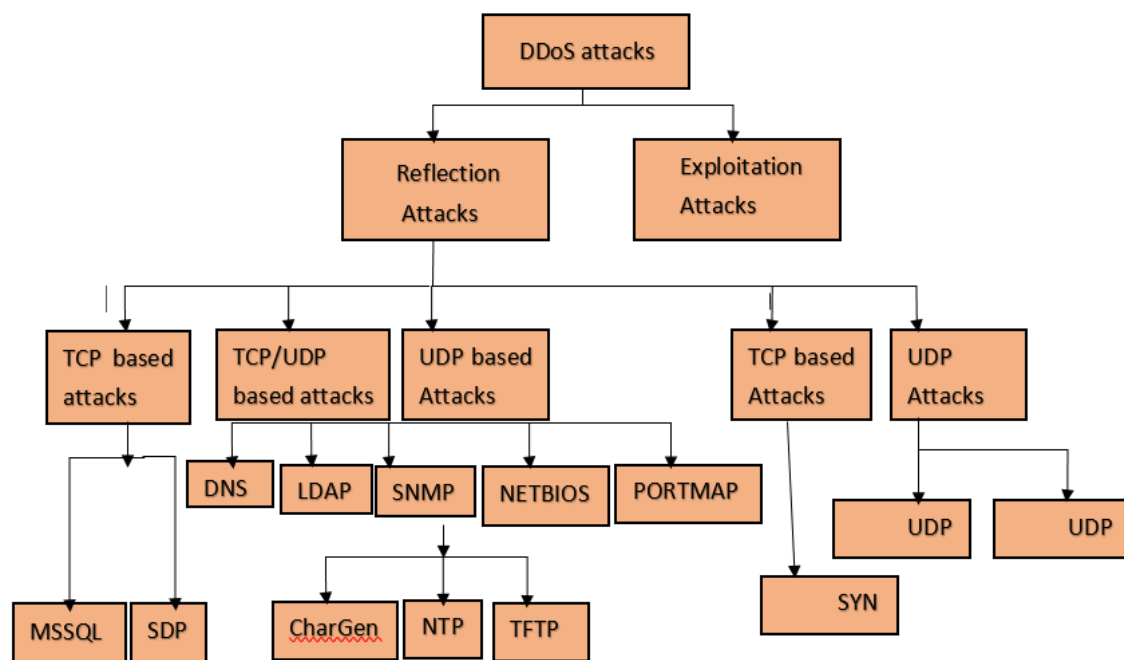


Figure 2.1: Types of DDoS Attacks

2.3.1 Types of DDoS Attacks

There are different types of DDoS attack as discuss in next sub sections.

1. TCP SYN Flood Attack

TCP SYN flood attack is a kind of DoS or DDoS attack. In this type of attack, an attacker attacks the server in such by using lot of connection requests on the server but the server does not respond to those requests and the system crash or become unable to use. It should also be kept in mind that the TCP SYN flood attack is implemented during the TCP session initialization handshake. Following are some of the steps that can be taken to overcome TCP SYN flood attacks.

- The network should have the ability to analyze the traffic that is coming from different parts of the network.
- There should be proper scalability to analyze and overcome attacks of different sizes ranging from low-end to high-end [16].

2. Teardrop attack

This type of attack attempts to make the computer resource unavailable by sending a lot of requests and data to the network or a server. For this type of attack, attackers use the fragmented packets and then send these packets to the network. Those fragmented packets then overlapped with each other because of the bugs during the reassemble phase of fragmentation, and then fails the server. Teardrop attacks are very vulnerable to many organizations because the organizations are still running legacy applications using the old and unpatched operating systems. There is one example that how F5's BIG-IP application delivers and overcomes the teardrop attacks. The F5 has a proper implementation that analyzes the packets that are coming into the network. If the incoming packets are rightly aligned, then the packets are kept otherwise packets dropped. If the users do not have the proper patches to prevent teardrop attacks, then SMBv2 should be disabled as well as block ports 139 and 445.

3. Smurf Attack

This type of attack is carried out by sending ICMP (Internet Control Message Protocol) echo requests (which is also known as ping) to the broadcast address of routers and involved other networks with some spoofed source addresses. Spoof means the address of the desired DoS attack. Those devices which are receiving the original ICMP echo request broadcast it to every other device that is connected by it and then each device gets that spoofed source address. Because of this, a high rate of requests is generated and the system become

unusable. Since 1999, most of the routers do not broadcast the ICMP request to the other connected devices and therefore system becomes safe from these attacks. For instance, if an IP network has 500 hosts then for each spoofed echo request 5000 responses will be produced. Following are some of the ways to prevent Smurf attacks.

- IP-directed broadcast on the router should be disabled.
- The operating system should also be reconfigured to prevent ICMP responses to IP broadcast requests.

4. Ping of Death Attack

This attack enables the attackers to destabilize and fail the targeted network or computer by sending a huge number of packets to the network using a simple ping command. In general, the normal IPv4 packet has a size of 65,535 bytes. Many traditional computer systems would crash if they receive such a huge size of packets. Since it is against the internet protocols, therefore the attackers transfer the packets in fragmented form. This attack is very easy because in this case, the attackers only have to know the IP address of the targeted systems. To prevent ping of death attacks, a firewall can be used that analyzes and check the size of fragmented IP packets.

5. Botnets

A botnet is generally a collection of different devices that are connected via the internet and the attackers can intervene in the system and might take full control of it using a botnet attack. This type of attack is generally implemented by cybercriminals and it results in data-stealing and theft. This type of attack has access to several networks or systems at one time and then using commands these attacks carry out malicious activities. These attacks also result in overwhelming the bandwidth of the targeted system. These types of attacks are difficult to analyze because the botnet varies in several geographical locations. Following are some of the countermeasures for botnet attacks.

- Use a firewall which is a basic cybersecurity tool
- System and software are always updated to the latest version
- Protection like Thor Foresight Enterprise that works very well at the DNS level

6. Man-in-the-Middle (MitM) Attack

When a hacker gets in between the communication of a client and the server, it's called a MitM attack. The following are some examples of man-in-the-middle attacks:

7. Session Hijacking

An attacker takes the control of a session between the client and a server. The computer that is being attacked changes its IP address to that of the trusted client and in the meantime, the server continues to communicate with the client. For instance, the attack could be explained as follows:

1. Client and server are connected
2. The computer of the attacker gains control of the client
3. It is then disconnecting the server from the client
4. The computer of the attacker recreates the client's sequence numbers and replaces the IP address of the client with its own
5. The attacker will continue talking to the server and it thinks that it is still in communication with the client

8. IP Spoofing

IP spoofing is used by an attacker to convince a system that is communicating with a known, and trusted entity and provide access to the system and to the attacker. It is done by sending a packet whose IP is known instead of its IP address so that the host could not know the scam and accept the packet.

9. Replay

Such an attack takes place when an attacker receives and keeps the old messages and tries circulating them around later on pretending to be one of the original recipients or point of contact in such a manner. The prevention of such an attack is possible with timestamps or nonce. This attack consists of unique attributes such as unethical and illegal means.

The volume of data, this attack contains, results in the upward trend of the cost associated with bandwidth which results in major delay or deletion of the transcripts of the original messages from the queue. The overall productivity of the system deteriorates as well.

If the frequency of replaying and deleting occurs, digital signature technology may prove to be obsolete and may result in a successful attack initiated by the attacker [17].

10. Phishing Attacks

In this type of attack, the email is sent to an individual who contains personal data. Attackers create the messages by themselves so that the victim does not know about the spam. It's very difficult to identify this type of attack because an individual could not know that the message is real or fake. In the "from" section the address is not real. That is why it is email spoofing. To reduce this type of attack and the risk some of the techniques are here: It also has different other types [18].

11. Spear Phishing

It is a type of email spoofing which sends the email to an individual containing personal and sensitive messages. Attackers spend their time researching the targets and creating personal and important messages. We concluded that detection of spear phishing is difficult and to defend is even more difficult. Email spoofing is one of the easiest ways for the hacker to conduct a spearfishing attack by falsifying an email in the "from" section. The host thinks that the email is from the known person or any company but in reality, it is from the hacker's side. Another technique that hackers use is website cloning. In website cloning, they copy the genuine website with the fake one so that one can open the link and enters the personal information which they can copy.

There are some techniques to reduce this type of attack:

12. Critical thinking

Do not ignore the sender's name and analyze the name and origin of the email after checking the details of the email then delete it or open it. Do not take it for granted by not seeing it.

13. Hovering over the links

Do not open the link directly. First, take the mouse on the link and see the details from which source it is coming.

14. Analyzing email headers

Email headers explain that how and from where an email was sent to you. The email's "Reply-to" and "Return-Path "parameters should not be from different domains.

15. Sandboxing

In sandboxing, by logging activities such as opening an attachment and clicking the links which are in the email, we can test the content of the email.

2.3.2 Drive-by Attack

The easiest method of spreading malware is Drive-by attacks. Hackers insert a malicious script into the code on the pages of vulnerable websites. When someone visits the site, malware is directly installed on the computer by the script, or it may refer the victim to a site operated by the hackers. When user visit the website and open an email or view a pop-up message, it results in the occurrence of drive-by attacks. A drive-by attack does not need a user to enable the attack actively. The user doesn't need to get infected by just opening or downloading a malicious email. A drive-by attack may take advantage of a security problem in an app, an operating system, or web browser because of failed or missing updates.

To avoid drive-by attacks users must keep their operating systems and browsers up to date, and avoid visiting websites that could have malicious code. Stick to the websites you're familiar with, but remember one thing that these too can be hacked. Don't load your device with too many useless programs and apps. The greater the number of plug-ins you have, the greater the number of flaws that can be exploited.

2.3.3 Password Attack

Passwords are the most widely used approach for the authentication of the user to an Information System. The password of a person can be easily accessed by peeking into a person's desk. Another technique sniffing the network connection to get the unencrypted password is obtained by gaining access to a password database or getting access to a password database. The final method should be applied randomly or systematically:

1. Brute-Force

In Brute force attack, passwords are being guessed in hope that one can find the right password. It's like a gamble to try passwords many times. Hackers try the password by connecting the person's name, company name, date of birth, etc.

2. Dictionary Attack

In this type of attack, different common passwords relating to that person to get access to a user's computer are used. One of the methods to get the password is to copy a password encrypted file then encrypt a dictionary of those common passwords which have been used, with the same encryption and then, compare the results. Lockout policy will protect you from to protect brute-force or dictionary attacks by implementing which locks the account after few attempts of failed passwords.

2.3.4 SQL injection attack

Structure Query Language (SQL) injection has become a major complaint with database-driven websites. It occurs when a malefactor uses the client-server input data to run a SQL query on the database. SQL commands are inserted into data-plane input to run predefined SQL commands (for example, instead of the login or password). SQL injection exploits can read sensitive data from databases and then modify data (insert, update, or delete) after that it performs database administration operations (like shutdown) then retrieve the contents of the given file, and in some cases, issue orders to the operating system. For example, when a user opens a web form from the website, it asks user information and sent it to the database to recover the information of the associated account using dynamic SQL.

SQL is vulnerable to this type of cybersecurity attack because it does not distinguish between the control and data planes. When a website uses dynamic SQL, SQL injections are most efficient. SQL injection is very common in the applications of PHP and ASP and because of these applications of J2EE and ASP.NET cannot easily abuse SQL injections due to the availability of design of the programmatic interfaces. Using the models of least privilege in your databases to defend yourself from SQL injection attacks. Stick to stored procedures and prepared statements (and make sure they don't contain any dynamic SQL) (parameterized queries).

Injection attacks must be avoided by the code that is not executed in the favor of the database. Additionally, at the application stage, verify the input data against the white list.

2.3.5 Cross-site Scripting (XSS) Attack

Third-party web tools are used in XSS attacks which run scripts in the web browser of the victim or scriptable program. A payload that contains malicious JavaScript is injected into the database of the website by the attacker. When a victim requests a page from the website, the website then sends the page to the browser of the victim along with the payload of the attacker that is embedded in the HTML body, which contains the malicious script. Like, it might send the cookie of the user to the server of the attacker, the attacker will then retrieve it and use it to hijack the victim's session. When XSS is used to operate the additional vulnerabilities, the most damaging effects arise. An intruder may use these flaws to steal cookies as well as log keystrokes, take screenshots, discover and collect network information, and remotely access and monitor the victim's computer.

Although XSS can be exploited in VBScript, ActiveX, and Flash, JavaScript is the most commonly exploited, owing to its widespread use on the internet. Developers should clean the data input by users in an HTTP request before reflecting it to prevent XSS attacks. Before echoing something back to the user, Like the query parameters' values during searches, confirm that all data is checked, filtered, or escaped. Convert all the special characters to their respective HTML or URL encoded equivalents.

2.3.6 Eavesdropping Attack

Interception of network traffic is used in eavesdropping attacks. Passwords, credit card numbers, and other sensitive information that a user might be transmitted over the network may be accessed by eavesdropping. Eavesdropping can be done in two ways: passively or actively:

a. Passive eavesdropping

By listening to the network's message transmission, a hacker may detect the content.

b. Active eavesdropping

By impersonating a friendly unit and sending queries to transmitters, a hacker actively obtains information. Probing, searching, or tampering are both terms for the same thing. Passive eavesdropping attacks are often more difficult to detect than active eavesdropping attacks since active attacks enable the attacker to first gain knowledge of friendly units through passive eavesdropping. Data encryption is the most effective eavesdropping deterrent.

2.3.7 Birthday Attack

Birthday attacks target hash algorithms, which are used to check the integrity of texts, applications, and digital signatures. A message digest (MD) of fixed length is generated by a hash function, irrespective of the length of the input message; this MD characterizes the message uniquely. When a hash function is used to process two random messages, the birthday attack states the likelihood of finding two random messages that produce the same MD. If the attacker measures the same MD for his message as the recipient, he will safely substitute the user's message with his, and even if the receiver compares MDs, he will not be able to detect the substitution.

2.3.8 Malware Attack

Malicious software is unwanted software that can be installed on the computer without your anyone's permission. It can attach itself to the real code and then spread; it can lurk in useful applications or replicate itself across the Internet. Here are some of the most common types of malware:

1. Macro viruses

A virus is a malicious program loaded onto a victim's computer without the user's knowledge to carry out malicious actions. A virus attaches itself to another program when loaded onto the computer and triggers only if that infected program is executed[19]. It replicates itself over the network and can paralyze the whole network. Viruses require user intervention to spread, whereas worms spread automatically. Due to these differences, infections are transmitted via email or Microsoft Word documents, which depend on the recipient to open the file or email to infect the system, which will be classified as a virus.

2. File infectors

Viruses that infect executable code, such as .exe files, are known as file infectors. When the code is loaded, the malware is installed. Another version of a file infector links to a file by making a virus file with the same name but a .exe extension. As a result, the virus code will run when the file is opened.

3. System or boot-record Infectors

A boot-record virus infects hard discs' master boot record. When the system boots up, it checks the boot sector for viruses and loads them into memory, where they can spread to other discs and machines.

4. Polymorphic Viruses

These viruses hide their presence through a series of encryption and decryption cycles. A decryption program first decrypts the encrypted virus and its associated mutation engine. The virus then infects a section of code. The virus encrypts the mutation engine and a copy of the virus with an algorithm corresponding to the new decryption routine, and the mutation engine implements a new decryption routine. The mutation engine and virus's encrypted package are attached to the new code, and the process is repeated. Because of the numerous changes to their source code, those viruses are hard to detect but still have a high level of entropy. This feature can be used to detect them by anti-virus software or free tools like Process Hacker.

5. Stealth Viruses

Stealth viruses cover the functions of the system to remain undetected. They accomplish this by hacking malware detection software, by causing it to account for an infected area as being uninfected. If there is an increase in the size of an infected file and changes the date and time of the file, are hidden by these viruses.

6. Trojans

Trojan horse is a computer virus. It is a type of computer software that is hidden just like normal software like utilities, games, and sometimes anti-virus programs. When it runs on a computer user will encounter a problem like disabling the background processing of the

system and deleting hard drive data or corrupting file association systems. Zeus, Dark comet, and Storm worm are some of the notorious Trojans.

7. Logic Bombs

It is a type of malevolent software that is attached to an application and is generated by some definite occurrence, as in a logical condition or a specific date and time.

8. Worms

Worms are amongst the most common form of malware. They work as a self-contained application and circulate over the network by replicating themselves [20]. Worms generally contain payloads that are designed to steal data or delete files. They reside in the memory without changing system files and keep on replicating themselves to make the system and network unresponsive[21]. Code Red, love letter worm, and Morris worm are famous worm attacks.

9. Droppers

Computers have programs names droppers to install viruses. One of the major advantages of a dropper is that it cannot be infected with malicious code. Therefore, it may not be detected by software that scans the viruses. It provides an update for the virus software that is already present in our system.

10. Ransomware

Ransomware is a type of malware program that infects, locks, or controls the system intending to solicit money from its owners. Ransomware and rapidly evolving variants are dangerous and harmful threats. It restricts users from gaining access to those files that are infected. The attacker detains the user and asks for money to get access again. Ransomware attacks are usually conducted through Trojans, and trick users into sending malware, ransomware, into their system.

11. Adware

Advertising-supported software is a type of malware designed to deliver advertisements to users. Often part of some free software or web pages. It may be associated with some Spyware that tracks activity and logs information about your preferences to

provide targeted advertising. Very common, especially for mobile applications, these applications are free when displaying advertising banners. You can then purchase the full version that removes this annoying banner. However, this cannot be considered adware because it is part of the initial agreement.

12. Spyware

Spyware, as the name implies, is commonly used to spy on users. This is a malicious program that aims to steal personal information from infected hosts without their knowledge and then send this information to the attacker. It is used to track user activity, steal passwords, bank account details, etc.

13. Bots

Bots, abbreviated as robots, are software that performs simple actions periodically. Some of them are used for purposes such as regular backup, auto auctions, and others. However, in the context of malware, it usually means periodically harmful actions. The best example is the DDoS attack described at the beginning of the chapter. Bots usually do not show visual interaction, so they can be hidden in the system without the user knowing.

14. Rootkits

Rootkits generally hide themselves or other harmful programs from victims[22]. Detection of rootkits is difficult because it uses root privileges, which enables corrupt system logs and intrusion detection software. Rootkits contain variations in tools, from programs that allow hackers to take your passwords to modules that can be used to easily steal credit cards or online banking information. Rootkits can also provide hackers the facility to hack or deactivate security software and track the keys you tap on keyboards, making it easy for criminals to steal your personal information.

15. Backdoors

It is an application in which an attacker is connected with a computer avoiding security mechanisms of the system which are hidden by the software in different secret ways. Backdoors provide an attacker with a remote shell (cmd.exe, bash, or special console) in the system to control the system remotely. Hackers are commonly used at the back of the house to hide their existence in the system after compromising the system. In addition, technical

support teams sometimes use backdoors to aid computer users. RAT, Remote Access Terminal is a popular backdoor tool [23].

2.4 Solutions

Machine Learning methods for classifications. In this section, we briefly describe the categories of machine learning algorithms used in this research. Malware detection tasks classify files into two categories: malware or benign. Therefore, in this study, we used five supervised learning techniques (classification model) for the task of detecting malware. In particular, the K-Nearest Neighbors (KNN), Decision tree, Support Vector Machines, Naïve Bayes, and Random Forests. An overview of the study techniques can be found in the section below.

2.4.1 Naïve Bayes

The Naive Bayesian classification comes from the Bayesian decision theory. This is a widely used classification method because of its manipulative capabilities and related possibilities based on user classification decisions and empirical performance. In the Bayesian classification, each class is represented by a summary of single possibilities. Assumptions are all class-independent independent attributes, hence the presence (or absence) of certain attributes of a class is not related to the presence (or absence) of other attributes. This means that Naive Bayesian assumes that all independent attributes contribute to the classification of possible classifications. Due to the precise model of possibilities, Naive Bayesian classification can be trained efficiently in a supervised learning environment. The authors provide an analysis of the average Bayesian classification case and provide experimental evidence for the usefulness of the Bayesian classification. The naive Bayesian classification has been successfully used in several complex applications, such as Bayesian spam filtering, which uses the Bayesian classification to identify spam emails [24].

2.4.2 Decision tree

Decision tables (also known as logic tables) or decision trees (also known as decision diagrams) describe situations related to a particular action or decision and the limitations of the associated action. In the decision table, rules are displayed as rows and conditions as columns with entries in each cell for each action which must be performed by people who do these acts. In comparison to other classifiers, the Decision Tree is the most straightforward

to comprehend and interpret. Figure 2.2 shows the decision tree. The basic steps in putting together a Decision Tree are as follows:

1. Determine the best attribute and place it on the top of the tree as a root node.
2. Divide the training dataset into subsets in such a way that every subset contains the only data which has the same attribute value.
3. Steps 1 and 2 should be repeated unless the leaf node of each branch of the tree could be found [25].

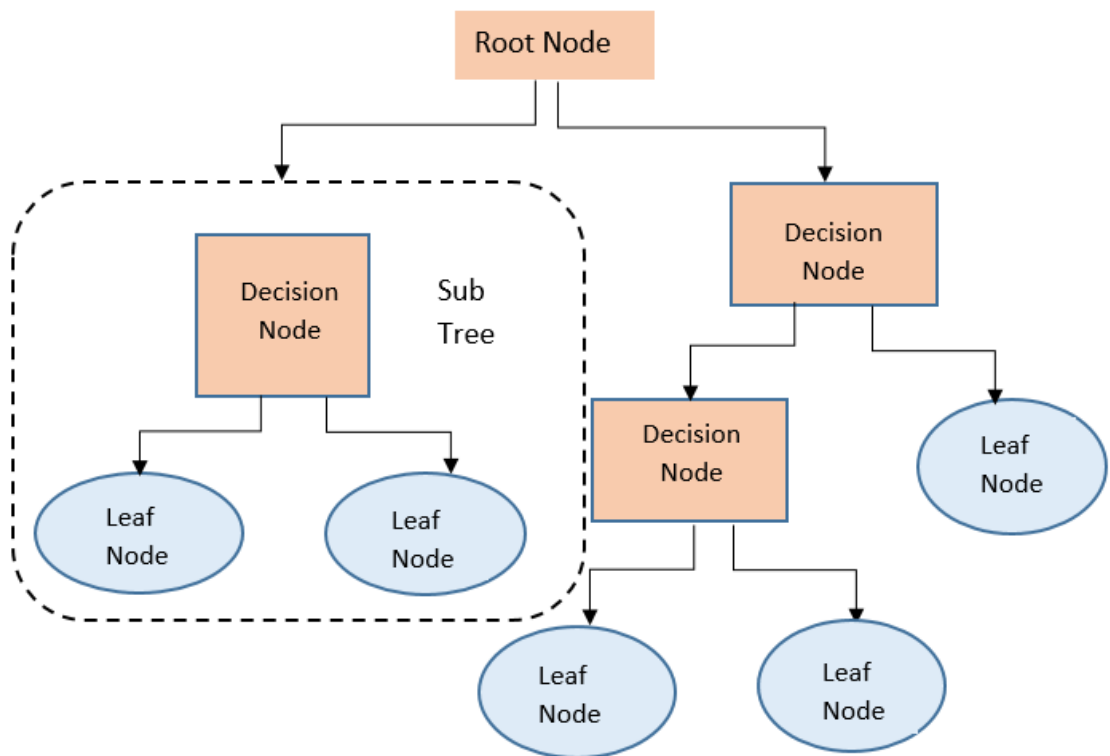


Figure 2.2: Decision Tree

2.4.3 Random Forest

A random forest is now a popular classifier in a variety of fields. A random forest is a mixture of trees that create a forest with several trees in it. In the random forest, the highest numeral of the tree means higher accuracy in results. It works step by step and takes decisions. Its use to identify diseases as well. The selection of the dataset randomization of the nodes while the construction of decision trees [26].

2.4.4 Support Vector Machine (SVM)

A support vector machine (SVM) is a machine learning model that is widely used in the field of Artificial Intelligence nowadays. In SVM high-dimensional hyperplane is constructed to separate the tags from two or more different classes to frame a model. Because the SVM deals with the high-dimensional data with a minimal preparing set of features, it is now popular for physiological data in clinical applications [7]. Figure 2.3 shows the SVM classifier.

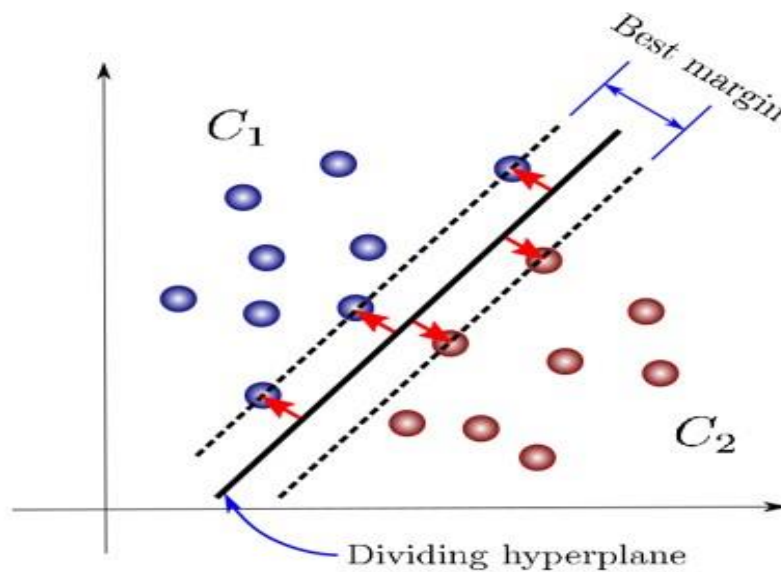


Figure 2.3: SVM Classifier

2.4.5 K-Nearest Neighbor (KNN)

The K-nearest neighbors (KNN) algorithm is a supervised machine learning (ML) algorithm that we can use in both classification and regression problems. However, in business, KNN is primarily used for classification and prediction problems. The following properties are used to define KNN.

- a) Lazy learning algorithm A lazy learning algorithm uses all of its data for training and does not have a specialized training phase during classification, KNN is a lazy learning algorithm.
- b) Non-parametric learning algorithm KNN is also a non-parametric learning algorithm because it makes no assumptions about the underlying data. The

KNN algorithm predicts the value data points that are new and based on 'feature similarity,' which means that the new data point is assigned based on its closeness with the points in the datasets.

Step 1: We need a dataset to implement any algorithm. As a result, during the first step of KNN, we must load both the training and test data.

Step 2: Next, we must select the value of K, i.e., the closest data points.

Step 3: Do follow each point that is in the test data: Calculate the distance between the test data and training data using by Manhattan, Euclidean, or Hamming distance. They are the common methods used for calculating distance.

- c) Arrange them in ascending order based on the value of the distance.
- d) It will then select the top row let say K from the sorted list.
- e) In the last step based on the most repeated rows of the class, it will assign it to the test point.

Step 4 – End

2.4.6 Multilayer Perceptron

The multilayer perceptron is a type of artificial neural network (ANN) technique that is commonly used in classification. It is a feedforward neural network that trains using backpropagation. It has a three-layer, input layer hidden layer, and output layer. The decision has not been reached about the comparative values of input variables individually, the number of inputs is to be by adjusting the weights all through the training phase, the training data distribution does not depend on the pre-assumptions in multilayer perceptron. Its uniqueness is that it has neurons in every layer [27].

2.4.7 Ensemble classifier

Ensemble classifiers, on the other hand, use autonomous algorithms to solve classification problems based on individual results provided by the primary algorithms. The boosted tree ensemble classifier, for example, employs a pre-set number of decision trees in such a way that the outcome of one tree is used to increase the number of contributing features in the next tree. As a result, the weighted average for final classification is calculated using a series of decision tree results. Independent decision trees are run in parallel in the bagged tree

ensemble classifier to provide results for the ensemble technique. The final classification is determined by a simple average or vote [28].

2.5 Background Research

In [1], the dataset of CAIDA 2007 and DARPA scenario-specific is used to evaluate algorithms. Supervised Machine Learning algorithms are used for the classification of a dataset. Random forest and Naïve Bayes algorithms are used for classification and the preprocessing of data ensemble algorithms is used. Random forest algorithm ranked first in accuracy and Naïve Bayes ranked second. The accuracy is detected by increasing recall and precision. In this paper, the proposed KNN algorithm is used with correlation analysis, thus the CKNN algorithm is used for the detection of DDoS attacks. The correlation approach increases the accuracy and decreases the overhead which is generally caused when the data is too big. The R-polling method is also used in this paper for reducing the training data [29].

In [30], authors detected DDoS attacks by using a novel method, for the feature reduction. A data Principle Component Analysis (PCA) method is proposed for the classification, and Naïve Bayes and KNN classifiers are also used. From the result, it is concluded that the average false positive rate is increased by 4.11% and the detection rate is 0.9%. In [31], data are collected through an online resource. Applications of the artificial neural network are used for the detection and classification of three types of DDoS attacks. The main target of this research is to develop a system using ANN for the detection and classification of DDoS attacks to increase the accuracy of those classes of DDoS traffic. The accuracy shown by the model is 82.1%. In another paper, the EPA-HTTP dataset is taken and the main objective is to detect DDoS attacks by using a multi-layer perceptron algorithm and a generic algorithm for learning data. Different parameters were selected instead of using the whole dataset This model shows an accuracy of 94% with 0.9962 sensitivity and 0.056 specificities [32].

The offline dataset is used in [33], for the detection of DDoS attacks by using passive monitoring. For training and testing, authors used different datasets for normal traffic and attack traffic and both datasets are obtained from different sources. In this paper, the Naïve Bayes algorithm method is used for classification. They used Naïve Bayes because it is simple and good for large datasets. It is based on Bayes Theorem. A deep learning approach is

proposed on the dataset of UNB ISCX Intrusion Detection Evaluation 2012 for the detection of DDoS attacks. The deep defense approach is based on CNN and RNN. RNN is used for feature extraction. They used four models of RNN and make a comparison between them and then they make a comparison between RNN and Random Forest in [34].

The dataset taken in [35], which is downloaded from the Wireshark tool. In this paper Packet Threshold Algorithm (PTA) together with SVM in his research for the detection of DDoS attacks. PTA SVM detects that whether the packet coming from the source is a normal packet or DDoS attack by setting the threshold for each type of packet. Four types of DDoS attacks (Smurf, TCP SYN flood, UDP flood, and Ping of Death) were detected. SVM is used for the classification of the trained dataset. After the detection of DDoS attacks, PTA-SVM was compared with other Machine Learning algorithms like Logistic Regression and K-means to get better accuracy. In [36], authors proposed a method for the detection of distributed denial of service attacks by using Naïve Bayes, SVM, and neural network classifier. The naïve bayed model calculates the probabilities and the instance which has the high probability was classified into their class. Then he uses the SVM model to classify and after that NN determines the labeled classes by calculating the weights. The dataset taken in this research was real-time data. Instead of training and testing the data separately, he created the topologies by using Mininet Emulator, and then the processed data had been trained. The Naïve Bayes algorithm gave 70 % accuracy which is slightly low because the data set was small. Then SVM gave 80% of accuracy and the highest accuracy was given by NN because it works best on complex features.

In [5], authors proposed a system for the detection of DDoS attacks by combining the three concepts: which is the execution of a distributed system by the classification algorithms and controlled by the fuzzy logic system. Naive Bayes, Random Forest, Decision Tree (Gini), and Decision Tree algorithms are used for the classification of packets. The system is examined by the performance of classification algorithms then the system is analyzed by using different sizes and numbers of datasets. Finally, the testing and the efficiency of the fuzzy logic system are examined. A fuzzy logic system is used because it selects the optimal classification algorithm at the right time. This research proposed a classification of DDoS attacks. He used Naïve Bayes and NN methods. The first step of his research was collecting a dataset from the Ahmad Dahlan University Research Laboratory network (LRis-UAD). The

next step was preprocessing and feature selection. They selected the features using the extraction method. The third step was the classification of packets by using NN and then by Naïve Bayes. Data was trained and tested by both the classifiers. The results showed that the accuracy of the Naïve Bayes classifier was higher than NN [37]. IoT is one of the most addressed areas for research nowadays.

In [8], authors proposed a classification technique to secure IoT for DDoS attacks. Due to the complexity of IoT DDoS attacks generates on the network layer. Naïve Bayes model with multi-agent based intrusion detection was used for the classification of data. Three phases were discussed in this research is the first phase data was processed for classification. In the second phase, the system collected the processed data and in the last step, the Naive Bayes model gives the accuracy on the test data and produces the results.

The data set taken in [38] is downloaded from CAIDA. (Center for Applied Internet Data Analysis) 2007 [38]. proposed an anomaly detection method in which these four steps are involved.

1. Preparing the data
2. Pre-processing of data
3. Processing of data
4. And post-processing of data for the detection of DDoS attacks.

In the preprocessing of data feature extraction is done on variance fractal dimension trajectory. For the processing of data, CNN is involved in the training and testing of data. The dataset is divided into training set 50% and testing set 50%. As it takes too much time to train the CNN on the training set, to overcome this time duration SVM is used for the post-processing of data. The DDoS attacks are detected with an accuracy of 87.3%. The dataset taken in [39] is downloaded from the Wireshark tool, proposed Packet Threshold Algorithm (PTA) together with SVM in his research for the detection of DDoS attacks. PCA-SVM detects that whether the incoming packet is a normal packet or DDoS attack by setting the threshold for each type of packet. SVM is used for the classification of the trained dataset. After the detection of DDoS attacks, PCA-SVM is compared with other Machine Learning algorithms like Logistic Regression and K-means but the PCA-SVM technique is better for the detection of DDoS attacks with 99.1% accuracy.

The dataset is taken in [32] and downloaded from CAIDA (Center for Applied Internet Data Analysis) proposed ensemble feature algorithms to select features to train the datasets for the classification. 16 features will consume time and resources so it is important to select the most appropriate and efficient features to deliver better accuracy. For the final classification of the dataset different machine learning algorithms are compared which include Naive Bayes, RBF, multilayer perceptron, and Random Forest. And plot the ROC curve for these algorithms. The ROC curve of multilayer perceptron is most efficient and is also reduces the computation time and gives better accuracy which is 98.3%. In [40], proposed classification algorithms for the detection of ICMPv6 based DDoS attacks. Five different machine learning algorithms are used and compared which are: Decision Tree, SVM, Naïve Bayes (KNN), and Neural Networks. The dataset used for classification is a flow-based dataset that is available publically. The results showed that KNN gives the best accuracy which is 85.7% among all algorithms in minimum time (0.12 seconds).

In [31], authors proposed a method to detect and classification of three classes of DDoS attacks traffic and one class of normal traffic by using an Artificial Neural Network. Dataset was collected from different online sources. The study of the results which are obtained by the simulation of the model proved that the extraction of parameters and by using an artificial neuron network can give high accuracy of 95.6% on the new dataset. This paper [41] focuses on the protection of an IoT infrastructure from DDoS attacks. For the detection of these intrusions, the Naive Bayes classification algorithm is applied along with the Multi-Agent System, which is responsible for the collection of the data which has to be analyzed. The data is further classified by the Naïve Bayesian classifier. Semi-supervised learning can be a way to obtain subsets of partially labeled or unlabeled datasets based on their similarities. Once the data has been recorded and the traffic flow observed, this data can be completely assigned a specific label based on the observations.

In [45], authors proposed a clustering approach to recognize disturbances in the flow of traffic and recognize normal traffic and DDoS attack traffic. The three strongest features are considered which can further be used at the target machine to recognize a DDoS attack. The author used K-means and Agglomerative clustering methods with feature extraction to obtain classes of the traffic of data. After the appropriate labeling, KNN, SVM, and Random

Forest were applied by the author to obtain models that could be used in the future for further classifications. Overall average accuracy of 94% was obtained between the three algorithms. In the end, the author also validated his method by applying different samples of DDoS attacks. In [42], authors used a Multi-Layer Perceptron along with a Genetic Algorithm to train the MLP neural network, for the detection of classification. The dataset of environmental protection agency-hypertext transfer was used and analyzed. This model gives an accuracy of 98.3%.

The dataset used in [2], from CICDDOS 2017. In this research different machine learning algorithms were used for the training and testing of data. Classification of port scanning and DDoS attacks was done in this research. ORT scanning is also a serious technique in cyber-attacks for scanning vulnerabilities of the target. The classifiers for this research were SVM, KNN, Ensemble Classifier, Decision tree. For feature extraction discrimination analysis was used. The results showed that the accuracy of SVM was better. Table 2.1 shows the comparison of discussed studies.

Table 2.1: Comparison Table

Authors	Approaches used	Scheme of the research	Published Year	Results
Robinson [1]	Naïve Bayes, KNN	Computer Networks and machine learning	2015	0.8% precision
Xiao [29]	CKNN (KNN with Correlation)	Computer Networks along with machine learning	2015	96% accuracy
Umarani, S Sharmila, D [30]	Naïve Bayes, KNN	Computer Networks along with machine learning	2015	False-positive rate with KNN 0.18% and with Naive Bayes 0.20%
Perakovic, Dragan Periša, Marko Cvitic, Ivan[31]	ANN	Computer networks with deep learning	2016	82.1%

Johnson Singh, Khundrakpam Thongam, Khelchandra De, Tanmay [32]	Multilayer Perceptron, Generic Algorithm	Cybersecurity with machine learning.	2016	98% accuracy
Fouladi, Ramin Fadaei Kayatas, Cemil Eren Anarim, Emin [33]	Naïve Bayes PCA, KNN	Computer Networks along with machine learning	2016	85%
Yuan, Xiaoyong Li, Chuanhuang Li, Xiaolin [34]	CNN, RNN	Computer Networks along with deep learning	2017	97% accuracy
Yusof, Mohd Azahari [35]	Naïve bayes	Computer Networks along with machine learning	2017	98% accuracy
Meti, Nisharani Narayan [36]	Naïve Bayes, SVM and Neural network classifier	Computer Networks along with machine learning and deep learning	2017	0.8% accuracy
Alsirhani, Amjad Sampalli, [5]	Naive Bayes, Random Forest, Decision Tree (Gini), and Decision Tree algorithms	Computer Networks along with machine learning	2018	85% accuracy
Yudhana, Anton Riadi, Imam Ridho [37]	Naïve Bayes, Neural Network	Computer Networks along with machine learning and deep learning	2018	84% accuracy
Mehmood, Amjad Mukherjee, Mithun Ahmed [38]	Naïve Bayes	Computer Networks along with machine learning	2018	60% accuracy

Ghanbari, Maryam, Witold Kinsner [39]	CNN, SVM	Computer Networks along with machine learning	2019	87.3% accuracy
Johnson Singh, Khundrakpam Thongam [32]	Naive Bayes, RBF, multilayer perceptron, and Random Forest	Computer Networks along with machine learning	2019	98% accuracy
Elejla, Omar E Belaton [40],	Decision Tree, SVM, Naive Bayes (KNN) and Neural Networks	Cloud computing with machine learning	2019	85.7% for KNN and 73% accuracy for SVM
Perakovic, Dragan Periša, [31]	Artificial Neural Network	Computer Networks along with deep learning	2020	95.6% accuracy
Mehmood, Amjad Mukherjee [41]	Naïve Bayes classification algorithm and Multi Agent System for data collection	Computer Networks along with machine learning	2020	85.6% accuracy
Aamir, Muhammad Zaidi [45]	KNN, SVM and Random Forest Agglomerative clustering methods for feature extraction	Computer Networks along with machine learning	2020	95%, 92% and 96.6% accuracy
Koay, Abigail Chen, Aaron Welch, Ian Seah, Winston KG [42]	Multi-Layer Perceptron Along with Genetic Algorithm	Computer Networks along with machine learning	2021	90.04% accuracy
Weiss, Jamie [42]	SVM, KNN, Ensemble Classifier, Decision tree	Computer Networks along with machine learning	2021	78% accuracy

CHAPTER 3

METHODOLOGY

3.1 Overview

In this chapter, the proposed methodology of the research discusses. Moreover, this chapter presents the three phases of research methodology framework.

3.2 Research Methodology Framework

Phase I: Problem Investigation

In the previous researches, dataset collection is always one of the biggest challenge because in the dataset of intrusion detection privacy and legal issues are the important things for which data is not easily available. Researchers make their data and it causes overfitting in classification. Accuracy is one of another challenge in the previous researches. Because of the data, the model does not give better accuracy in the previous researches. While investigating the problem we came across to know that in the previous researches the classification is in the limited classes. This study tries to overcome these challenge by using 13 classes.

Phase II: Design

In the design phase, the overview of the proposed methodology presents to design the proposed model.

A. System Overview

The dataset is used to classify with the help of a domain expert. After that, we use the machine learning technique to train proposed model. In addition to the hand-crafted features, we intend to investigate feature extraction based on select K best features. In the former case, an additional classifier is required to classify the DDoS attacks class.

B. Dataset

The accessibility of datasets is one of the most important challenge for intrusion detection approaches using Machine learning/deep learning. The datasets are not easily

available in the intrusion detection domain because of privacy and legal considerations. The network traffic includes highly confidential information, and its availability can expose the privacy of customers and businesses, as well as personal correspondence. To prevent any sensitive issues, many researchers generate their data to fill the preceding gap. However, in these conditions, the majority of the datasets created are incomplete, and the row samples used to cover the application behaviors are insufficient [34]. The most used and popular datasets which are publically open for research are CICIDS2017, KDDCUP99 [45], Kyoto 2006+, NSL-KDD [53], and ISCX2012 [54].

We used CICDDoS 2019 dataset in this research which is a newly released dataset [55]. This dataset is having a large number of DDoS attacks which could be completed by application-layer protocols that are using TCP/UDP. The attacks in the dataset are categorized as either exploitation-based or reflection-based attacks. For training and testing purposes, the dataset is obtained on two different days. The training package, which is recorded on January 12th, 2019, includes 12 different types of DDoS attacks, each in its PCAP format. The types of attacks included for training are UDP, SNMP, NetBIOS, LDAP, TFTP, NTP, SYN, WebDDoS, MSSQL, UDP-Lag, DNS, and SSDP DDoS based attacks. On 11th March 2019 testing data was formed, containing 7 DDoS attacks which are SYN, MSSQL, UDP-Lag, LDAP, UDP, PortScan, and NetBIOS. Dataset has been categorized in fig. The dataset has more than 8 features which were extracted by using CICFlowMeter tools [43]. The CICDDoS2019 dataset is available on the Canadian Institute for Cybersecurity website in PCAP file and based on flow format. Figure 3.1 shows the CICDDoS 2019 Dataset Distribution.

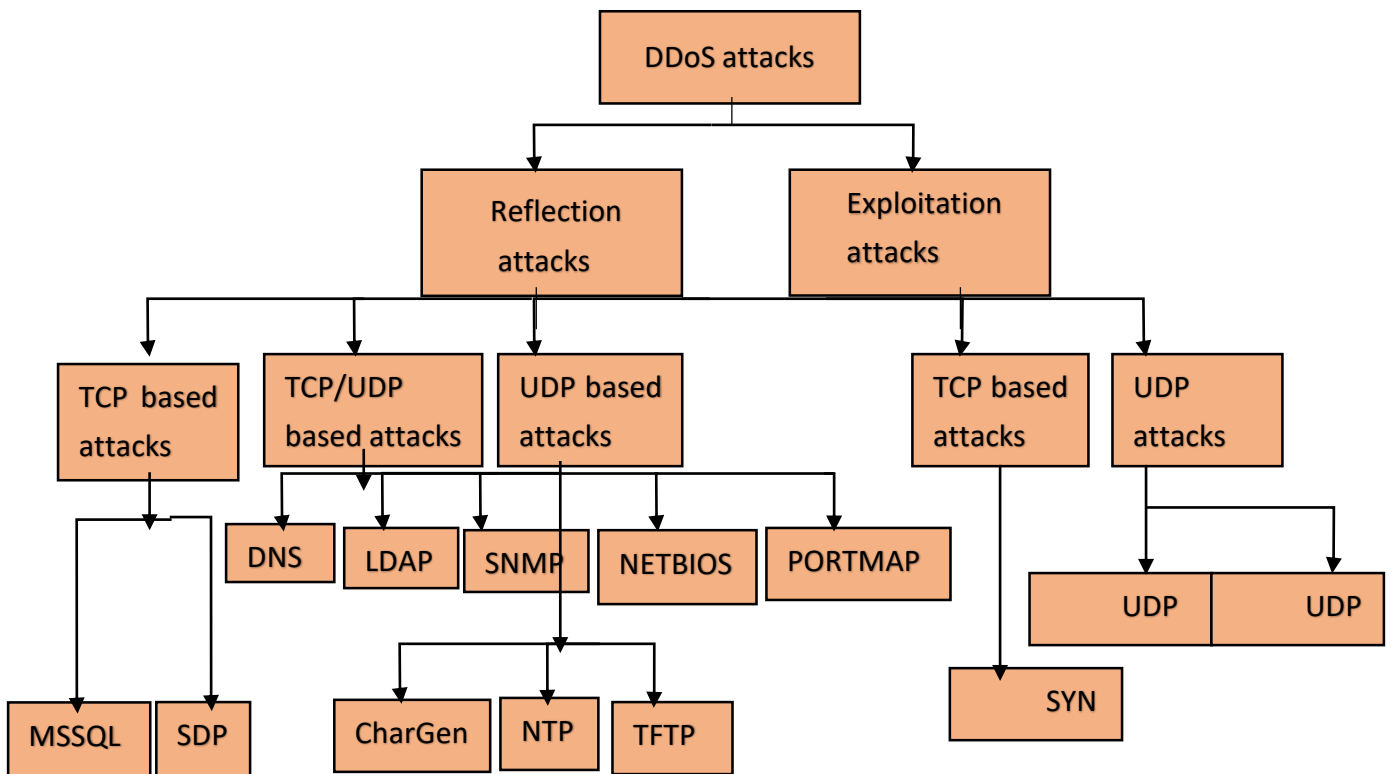


Figure 3.1: CICCDDOS 2019 Dataset Distribution

Phase III: Evaluation

This section presents the evaluation metrics, used when evaluating the detection and classification. Because theory about these metrics has already been thoroughly described in the previous report, this section briefly describes them.

3.3 Confusion Matrix

The confusion matrix provides an overview of the predictions made by the classifier. This includes the following four metrics:

- a) True Positive (TP) is when the condition is correctly predicted as positive.
- b) True Negative (TN) is when the condition is correctly predicted as negative.
- c) False Positive (FP) is when the condition is incorrectly predicted as positive.
- d) False Negative (FN) is when the condition is incorrectly predicted as negative.

In addition to using accuracy as a measure for performance, confusion matrix is another simple but effective method for describing the performance of an algorithm. It further

extends the functionality of the accuracy parameter in such a manner that it also keeps records of the incorrectly classified data along with the correct data. A confusion matrix is based on four possible values depending on the statistics of the algorithm when it is run on a test data instance:

a. True Positive

Data instance belongs to a specific class and is correctly classified by the algorithm that data belongs to the same class. In the case of signature, an attack is correctly classified as an attack.

b. False Negative

The algorithm detects that the data does not pertain to a specific class, however it belongs to that class. In other words, an attack that is wrongly classified as legitimate traffic.

c. False Positive

The data instance does not belong to a specific class, but it is incorrectly identified by the algorithm as belonging to that class. In this case, legitimate traffic is wrongly classified as an attack. This case is most likely occurred while using honeypots. The aim is to minimize this rate to lower false-positive alerts.

d. True negative

The algorithm identifies that data does not belong to a specific class; however, the data belongs to another class.

3.3.1 Accuracy Method

For calculating the accuracy prediction errors, we going to use this accuracy equation.

$$\text{Accuracy} = (T P + T N) / (T P + T N + F P + F N)$$

3.3.2 Additional Metrics

a. Precision

Precision is used to measure the positive patterns that are correctly predicted from the total predicted patterns in a positive class.

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

b. **Recall:** Recall is used to measure the fraction of positive patterns that are classified correctly.

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

c. False Positive Rate (FPR) = $\text{FP} / (\text{FP} + \text{TN})$

d. Detection Rate = $\text{TP} / (\text{TP} + \text{FN})$

CHAPTER 4

PROPOSED WORK DESIGN

4.1 Overview

In this chapter, technical details of the proposed model presents. Over the years, the advancement in processing computing power has to lead the harnessing of the powers of Neural Networks. A network consists of 3 basic layers input layer, hidden layer, and output layer as shown in Figure 4.1.

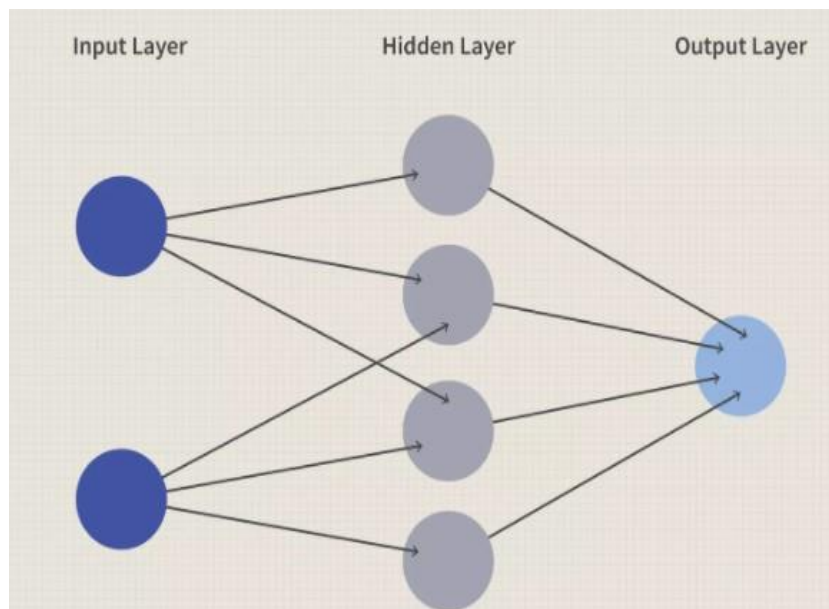


Figure 4.1: Simple Neural Network

4.2 Recurrent Neural Networks

Recurrent Neural Networks are derived from feedforward neural networks. RNNs can overcome the traditional feed-forward neural networks problem which results in creating more powerful models with high accuracy of classification [46]. They use their internal state which is also called memory to process variable length of the sequence of inputs. The output value of the input sequence depends on past computed value one of input or output or both can be a sequence [47]. Figure 4.2 shows the RNN architecture.

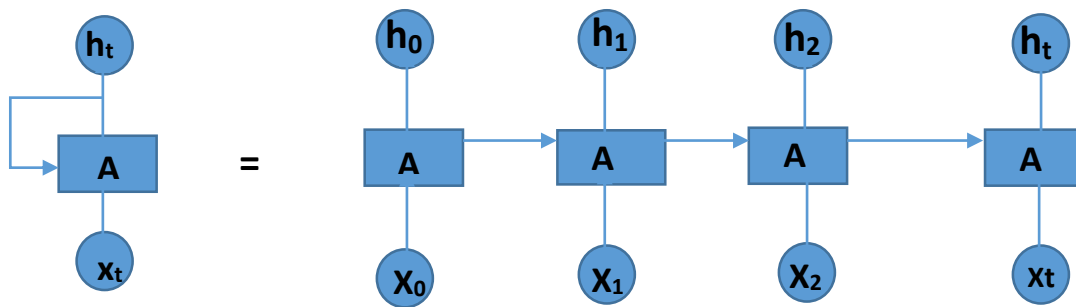


Figure 4.2: RNN Architecture

In RNN, the weight matrix always repeats its multiplication during the training and it creates gradient vanishing or explosion problems. RNN has different variants like the long short-term memory (LSTM) and gated recurrent unit (GRU) [48]. They overcome the problem of gradient vanishing. The hyperbolic tangent and sigmoid functions are used as the activation functions in LSTM and GRU, resulting in gradient vanishing over layers. As a result, building and training a deep LSTM or GRU-based RNN network is very difficult. Before that, in CNN Relu is used as an activation function that can stack into a deep network (for example more than 20 layers by using the simple convolutional layers and more than 100 layers with residual connections) and still can train efficiently [45]. RNN can be applied to different kinds of applications for example translation, Captioning [43], Sentimental Classification [49], and Speech recognition [50].

4.3 LSTM

Long Short Term Memory (LSTM) is presented to avoid the problem of backpropagation errors from exploding. LSTM introduced forget gates to prevent the problem of long-term dependency. The forget gate monitors how information is used in the cell states. Because of their ability to capture long-term dependencies, LSTMs can outperform conventional RNNs in capturing nonlinear dynamics in time series sensory data and learning effective representations of system conditions. Given that LSTMs have been implemented in a variety of applications successfully, which includes speech recognition, image captioning, recognition of handwriting, genomic analysis, and natural language processing since they can identify long-range correlations and nonlinear dynamics in time series data [51]. Figure 4.3 shows the LSTM architecture.

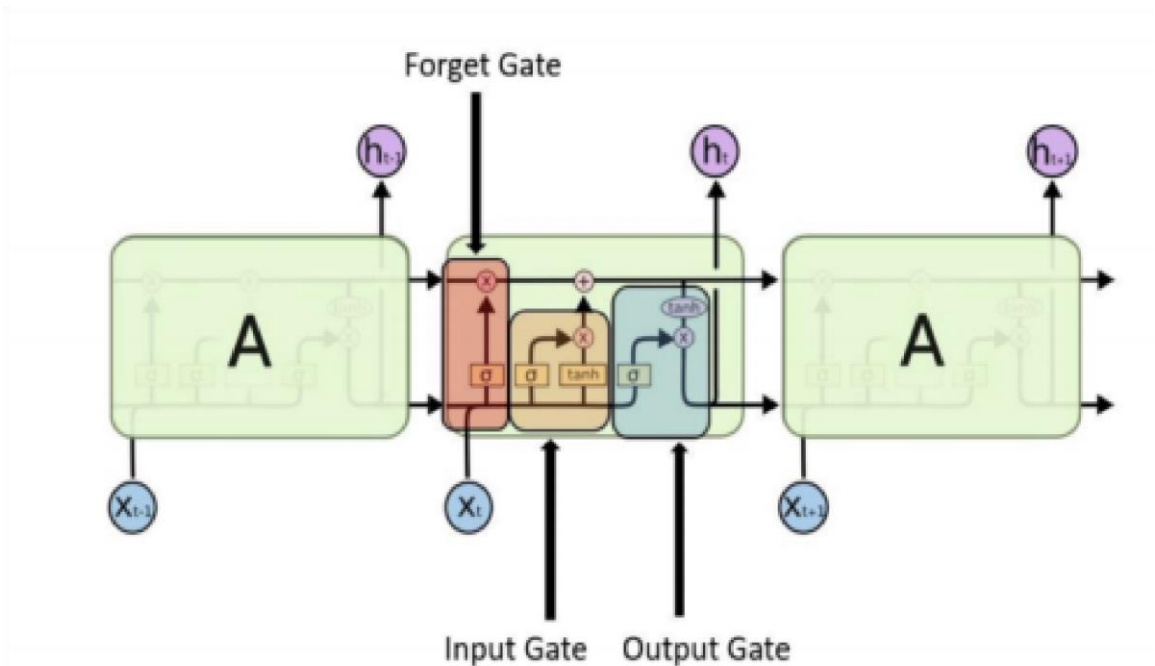


Figure 4.3: LSTM Architecture

As discussed earlier that LSTM uses hyperbolic tangent(\tanh) and sigmoid functions as gate functions [45]. In LSTM when multiple layers are combined in the deep layer, it overcomes the problem of gradient vanishing [52].

4.4 System Overview

An overview of the intended classifying system is presented in Figure 4.4.

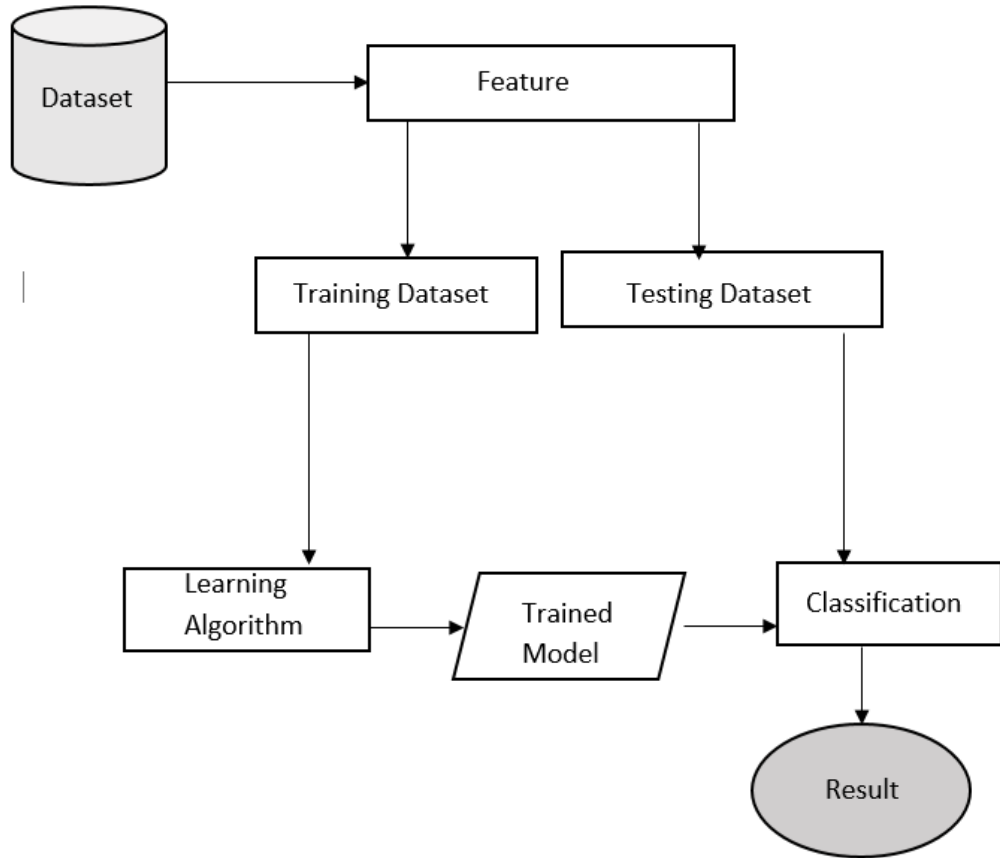


Figure 4.4: System Overview

4.4.1 Dataset in the Proposed Methodology

The first step of this study is Dataset. The dataset taken in this research is from CICDDOS 2019. It is publicly available for free. The dataset consists of 12 classes. This dataset is having a large number of DDoS attacks which could be completed by application-layer protocols that are using TCP/UDP. The attacks in the dataset are categorized as either exploitation-based or reflection-based attacks. For training and testing purposes, the dataset was obtained on two different days. The training package, which was recorded on January 12th, 2019, includes 12 different types of DDoS attacks, each in its PCAP format. The types of attacks included for training are UDP, SNMP, NetBIOS, LDAP, TFTP, NTP, SYN, Web DDoS, MSSQL, UDP-Lag, DNS, and SSDP DDoS based attacks.

A. Feature Extraction

Dataset was in the raw form and contains too many redundant values and zeros. Then Feature extraction process is executing. Feature Selection was done by using different

learning algorithms. The filter method is used. One of the best approaches of filter feature selection method is the uni-variate method is used to identify the important features of the dataset. SVM, LSTM, Decision Tree, and Naive Bayes methods are used in this research to select k best features to form the given dataset.

B. Feature Preprocessing and Selection

Feature preprocessing and feature selection is the technique that chooses the best features automatically from the dataset. In feature selection process achieves the subsection of the best-related features of the dataset deprived of reusing them. It is used for improving the information to achieve the best accuracy. Feature selection uses in different fields, such as 22 data mining and ML applications. The benefits of making feature selection are removing over-fitting, improving accuracy, and reducing training time [43]. The feature selection method increases the significant data from present features and accomplishes the most noteworthy correctness of classifiers [44]. Feature selection is used in medical field for the best problem-solving systems. Feature selection has worth for referencing approaches. We performed different feature techniques to check suitable accuracy in various states in the classification algorithm. Following methods are used for selecting features.

1. Remove the columns which are greater than the given threshold and have which have missing fractions.
2. Remove the features which have only one unique value.
3. Remove the collinear features which are greater than the given value as identified by the correlation coefficient.
4. Remove the features which have 0.0 importance from a gradient boosting machine.
5. Remove features that are not a part of the specified cumulative feature importance from the gradient boosting machine.

C. Filter Method

The filter method depends on the information's overall uniqueness to be assessed and pick highlight subset, excluding any mining calculation. Filter techniques utilize the specific review rule, which incorporates separation, data, reliance, and consistency. Filter techniques use the vital standards of positioning procedure and utilization's the rank requesting technique for variable determination. The purpose behind utilizing the positioning strategy is

straightforwardness, produce significant and pertinent highlights. The classification method will sift through unessential highlights before the characterization measure begins. The data preprocessing step filter method is used for feature selection, independent of any learning algorithm. It gives statistical scores that determine the correlation in the output variable.

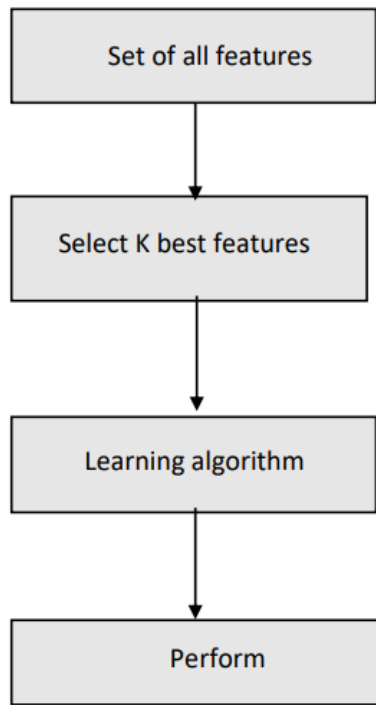


Figure 4.5: Feature Selection Method

One of the best approach for filter feature selection method is the univariate method that is used to identify the important feature of the dataset. Each feature is selected according to specified criteria and based on higher scores and ranks. In the univariate feature selection method select the best k features through `SelectKBest ()` class [45].

4.4.2 Training Dataset

Dataset is split into training and testing in the ratio of 30% and 70% for testing and training. Learning algorithms are applied to the training dataset to train.

4.4.3 Testing Dataset

We split the dataset into 30% for testing and 70% for training. For testing, we used LSTM and SVM classifiers on the dataset.

4.4.4 Classification

Classification is done using LSTM and then it is compared with SVM.

CHAPTER 5

RESULTS

5.1 Experiment and Results

This section presents the facts of different experiments with the results that we performed. We present feature selection results with traditional and deep learning classifiers, then data processing, and a trained model for the results. We divided the dataset into 70% of the trained dataset and 30% in the testing dataset, examining the performance variation.

5.2 Environment and Tools

Google Colab Environment is utilized to prepare and test the specified model. It is like a Jupiter notebook within the cloud environment and it handles all setup configurations. With the assistance of Google Colab, one can compose and execute through the browser. Google Colab permits free GPU for a few hours a day. For data preprocessing, I have used Spyder, which is a free source.

5.3 Optimizer

Adam optimizer is used with the default initial learning rate and weight decay.
optimizer = Adam (lr=0.001, decay=0.9)

5.4 Training Model

The training details of the model are listed below • number of DDoS: 125376

- number of benign: 84330
- epochs are =500
- batch size = 250

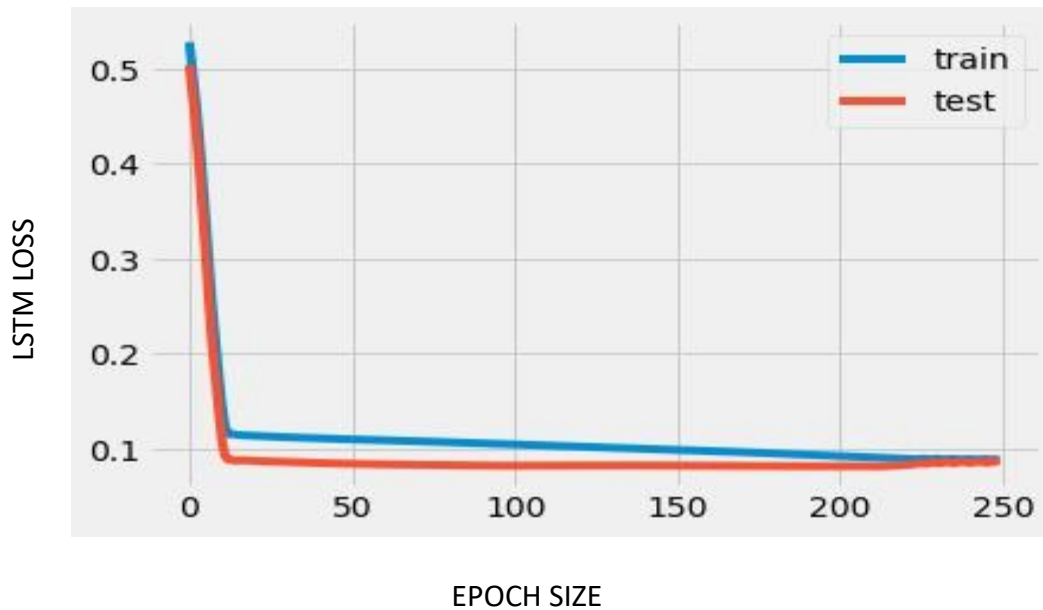


Figure 5.1: LSTM Loss

LSTM has shown in Figure 5.1, the best results and its training and testing loss is also almost the same after 200 epochs.

5.5 Hyper-Parameter Tuning

Hyperparameter values control the behavior of the trained model directly. The selection of the best hyperparameter values depends on human knowledge and practice. Different experiments are conducted to select the best values of hyperparameters. To test the model performance, we used different values of learning rates like 0.1, 0.001, etc. The performance of the LSTM model is shown in the figure below which shows the accuracy with multiple learning rates. Based on the results, we concluded that when we reduce the value of the learning rate, the model takes a long time. However, it gives better results when we used the small learning rate value. The model gave the best accuracy of 99% when using 0.0001 learning rate. We also changed the activation functions with learning rates and also hidden layers.

The data is divided into two sections.

1. Training
2. Testing

The training set is used to build the model by adjusting the weights on NN. The training set is used for the parameters of the experiment like hidden layers in the model finely. The test set is used to evaluate the accuracy of the model. In my thesis, I used the train-test split technique for the evaluation of the model. I didn't use k-fold cross validation though it is a widely used classification evaluation method. Cross-validation is not very suitable for time series data. Table 5.1 shows the performance metrics with different learning rates.

Table 5.1: Performance metrics with different learning rates

Learning Rates	Performance Metrics			
	Recall	Precision	F-Score	Accuracy
0.1	0.97	0.96	0.97	97%
0.01	0.99	0.98	0.98	98%
0.001	0.99	0.98	0.98	98%
0.0001	0.99	0.99	0.99	99%

5.6 Results

The CICDDos Attack dataset consists of 20414839 instances and 84 features variables and 7 different subtypes of DDoS attacks. These all sub-type of datasets have unbalanced instances. We have divided into two formats of datasets (Balance and Unbalance) for experiments. In balance datasets, we have an equal number of instances. A minimum number of instances in each class have 52187 while unbalancing dataset number of instances in each class. The maximum number of instances in MSSQL class is 5787453 and the minimum number of instances in UDP Lag has 52187 in unbalancing dataset. The results were obtained by the traditional and deep learning classifiers (SVM and LSTM) on balance and unbalance datasets. We have experimented 10 times and noted average accuracy in the given table. We observed LSTM showed high accuracy than the SVM classifier. It is common nature of deep learning model performs high accuracy on a large number of samples. Table 5.2 shows the results of balanced and unbalanced datasets.

Table 5.2: Results of balanced and unbalanced datasets

Classes	Number of Samples	
MSSQL	5787453	
LDAP	1915122	
BENIGN	56965	
NetBIOS	3657497	
Portmap	186960	
SYN	4891500	
UDP	3867155	
UDP Lag	52187	
Total	20414839	
Dataset	SVM	LSTM
Balance	88.3%	92.5%
Unbalance	95.8%	98.3%

5.6.1 Result of Select K-Best

Analysis of important features of classification models using Select k-Best. It gives different accuracy result in a different range of features. It shows the highest accuracy on features shown in Figure 5.2 below:

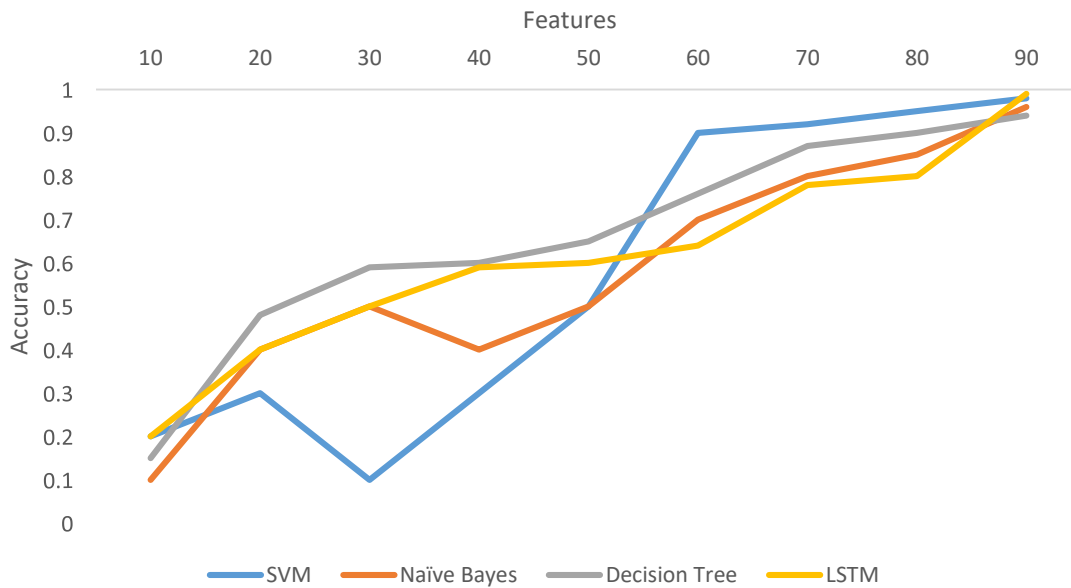


Figure 5.2: Feature analysis using Select k Best Method

5.6.2 Accuracy Analysis

Our system performances depend on the accuracy which we obtained through the analysis. We used the five classification models with different feature ranges. Accuracy provides the ratio between total instances examined and total correctly identified instances in a dataset. Accuracy is compared against every algorithm in each data set. The process of measuring accuracy is further segregated into three types named training accuracy, test accuracy, and cross-validation accuracy.

5.6.3 Training Accuracy

Training accuracy measures the classifying precision of an algorithm when it is applied to data instances belonging to a certain data set. In our research, training accuracy is measured as an average taken from the multiple runs of an algorithm on a data set. These multiple runs of the algorithm provide a more reliable result training accuracy.

5.6.4 Testing Accuracy

Testing accuracy means measuring the accuracy of an algorithm when it is applied to unseen data instances. Simply put, testing accuracy measures how the algorithm will classify data instances previously not seen by the algorithm. In testing accuracy, the average value is

taken based on the result of multiple measurements. Figure 5.3 shows the accuracy of training and testing dataset.

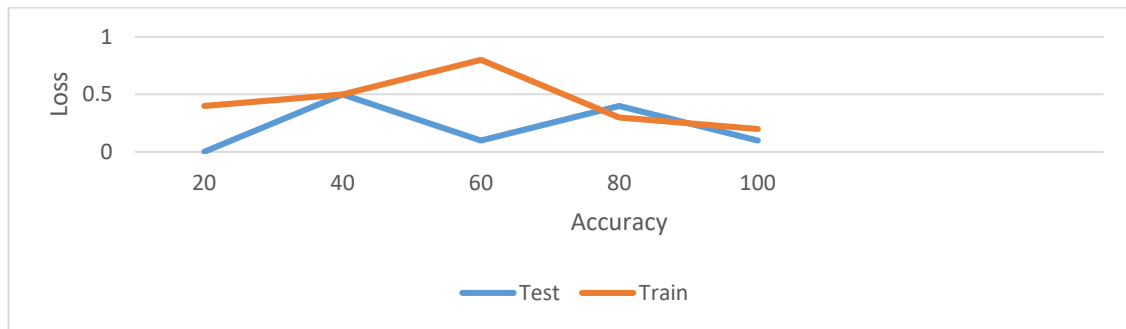


Figure 5.3: Accuracy of Training and Testing Dataset

The graph shows the accuracy and loss of the data. And it shows that LSTM shows better accuracy and loss. Figure 5.4 shows the ROC curve.

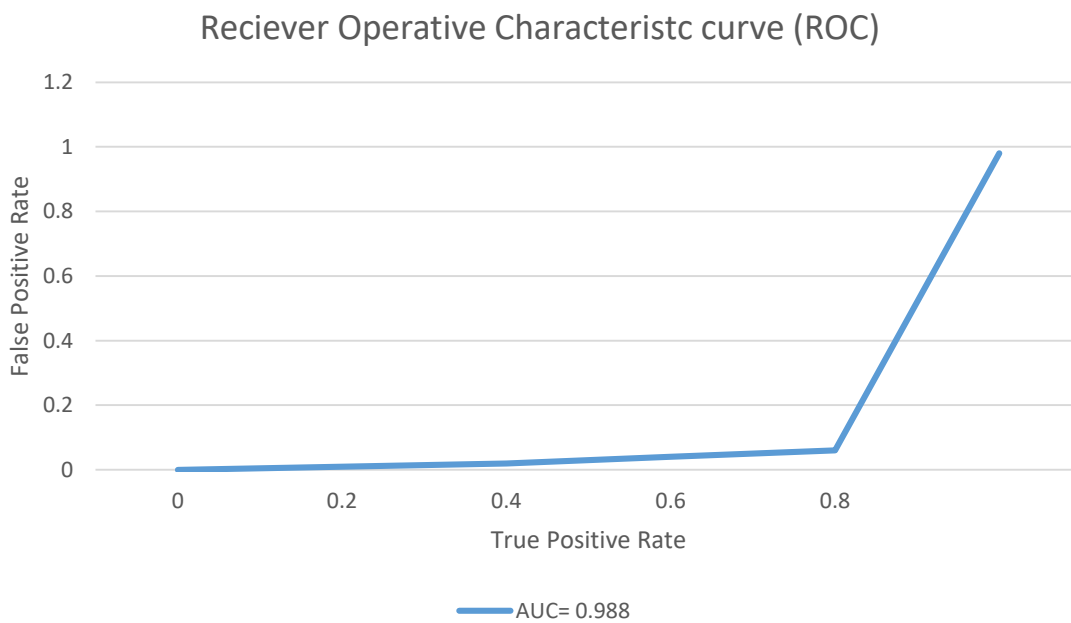


Figure 5.4: ROC curve

We are using the ROC curve to evaluate the accuracy of our model. It specifies the relation between true and false classes. In ROC the area which is under the ROC curve measures the separation between true and false-positive rates. Our model gives the 99.8% AUC which means our model is separating 99.8% of classes into negative and positive.

Furthermore, we used a confusion matrix for the evaluation of our model. It sums up the problems in true and false predictions. Table 5.3 shows the confusion matrix.

Table 5.3: Confusion Matrix

Attack	0.99	0.01
Normal	0.11	0.99
	Attack	Normal

CHAPTER 6

CONCLUSION

Recent deep learning techniques enable the system to be end-to-end trainable, this technology advancement opens the doors for researchers to use deep learning techniques in computer vision and make the computer vision systems robust and end-to-end learnable. Network virtualization introduces new threats and vulnerabilities to the conventional network, in addition to the ones that already existed. The DDoS attack type is one of the most violent attack types in current years, wreaking havoc on the entire network system. This study proposed the classification techniques using the deep learning model to classify the newly released dataset of CICDDOS 2019. The dataset contains complete and current DDoS attack types. The evaluation of proposed model showed that Long Short Term Memory (LSTM) which is a modification of Recurrent Neural Network is used for the classification of DDoS attacks which gives the highest accuracy as compared to SVM.

Future Work

In the future, we are planning to test our proposed system on other different datasets. In this research, we divided the dataset into balanced and unbalanced datasets classify them. On the other hand, each class of DDoS attacks must be classified independently in the future. We plan to spread our work to a multi-class classification model in the future.

References

1. Robinson, R.R. and C. Thomas. *Ranking of machine learning algorithms based on the performance in classifying DDoS attacks*. in *2015 IEEE Recent Advances in Intelligent Computational Systems (RAICS)*. 2015. IEEE.
2. Weiss, J., *DDoS Detection Using Deep Neural Networks on Packet Flows*. 2019.
3. Haider, S., et al., *A Deep CNN Ensemble Framework for Efficient DDoS Attack Detection in Software Defined Networks*. IEEE Access, 2020. **8**: p. 53972-53983.
4. Wani, A.R., et al. *Analysis and Detection of DDoS Attacks on Cloud Computing Environment using Machine Learning Techniques*. in *2019 Amity International Conference on Artificial Intelligence (AICAI)*. 2019. IEEE.
5. Alsirhani, A., S. Sampalli, and P. Bodorik, *DDoS Detection System: Using a Set of Classification Algorithms Controlled by Fuzzy Logic System in Apache Spark*. IEEE Transactions on Network and Service Management, 2019. **16**(3): p. 936-949.
6. Tama, B.A., M. Comuzzi, and K.-H. Rhee, *TSE-IDS: A two-stage classifier ensemble for intelligent anomaly-based intrusion detection system*. IEEE Access, 2019. **7**: p. 94497-94507.
7. Roopak, M., G.Y. Tian, and J. Chambers. *Deep learning models for cyber security in IoT networks*. in *2019 IEEE 9th annual computing and communication workshop and conference (CCWC)*. 2019. IEEE.
8. Mehmood, A., et al., *NBC-MAIDS: Naïve Bayesian classification technique in multi-agent system-enriched IDS for securing IoT against DDoS attacks*. The Journal of Supercomputing, 2018. **74**(10): p. 5156-5170.
9. Singh, K.J. and T. De, *Efficient classification of DDoS attacks using an ensemble feature selection algorithm*. Journal of Intelligent Systems, 2020. **29**(1): p. 71-83.
10. Jia, Y., et al., *Flowguard: an intelligent edge defense mechanism against IoT DDoS attacks*. IEEE Internet of Things Journal, 2020. **7**(10): p. 9552-9562.
11. Tidjon, L.N., M. Frappier, and A. Mammari, *Intrusion detection systems: A cross-domain overview*. IEEE Communications Surveys & Tutorials, 2019. **21**(4): p. 3639-3681.
12. Kumar, D.A. and S. Venugopalan, *INTRUSION DETECTION SYSTEMS: A REVIEW*. International Journal of Advanced Research in Computer Science, 2017. **8**(8).

13. Benkhelifa, E., T. Welsh, and W. Hamouda, *A critical review of practices and challenges in intrusion detection systems for IoT: Toward universal and resilient systems*. IEEE Communications Surveys & Tutorials, 2018. **20**(4): p. 3496-3509.
14. Kim, K. and M.E. Aminanto. *Deep learning in intrusion detection perspective: Overview and further challenges*. in *2017 International Workshop on Big Data and Information Security (IWBIS)*. 2017. IEEE.
15. Karatas, G., O. Demir, and O.K. Sahingoz. *Deep learning in intrusion detection systems*. in *2018 International Congress on Big Data, Deep Learning and Fighting Cyber Terrorism (IBIGDELFT)*. 2018. IEEE.
16. Bawany, N.Z., J.A. Shamsi, and K. Salah, *DDoS attack detection and mitigation using SDN: methods, practices, and solutions*. Arabian Journal for Science and Engineering, 2017. **42**(2): p. 425-441.
17. Sun, Y., et al., *Attacks and countermeasures in the internet of vehicles*. Annals of Telecommunications, 2017. **72**(5-6): p. 283-295.
18. Borkar, A., A. Donode, and A. Kumari. *A survey on Intrusion Detection System (IDS) and Internal Intrusion Detection and protection system (IIDPS)*. in *2017 International conference on inventive computing and informatics (ICICI)*. 2017. IEEE.
19. Sihwail, R., et al., *Malware detection approach based on artifacts in memory image and dynamic analysis*. Applied Sciences, 2019. **9**(18): p. 3680.
20. Spafford, E.H. *The internet worm incident*. in *European Software Engineering Conference*. 1989. Springer.
21. Yang, L.-X., et al., *A novel computer virus propagation model and its dynamics*. International Journal of Computer Mathematics, 2012. **89**(17): p. 2307-2314.
22. Qamar, A., A. Karim, and V. Chang, *Mobile malware attacks: Review, taxonomy & future directions*. Future Generation Computer Systems, 2019. **97**: p. 887-909.
23. Sugunan, K., T.G. Kumar, and K. Dhanya, *Static and dynamic analysis for android malware detection*, in *Advances in Big Data and Cloud Computing*. 2018, Springer. p. 147-155.
24. Rukmawan, S., et al. *Cerebral Infarction Classification Using the K-Nearest Neighbor and Naive Bayes Classifier*. in *Journal of Physics: Conference Series*. 2021. IOP Publishing.
25. Wang, F., et al., *A linear multivariate binary decision tree classifier based on K-means splitting*. Pattern Recognition, 2020. **107**: p. 107521.

26. Chen, L., et al., *Detection of DNS DDOS attacks with random forest algorithm on spark*. Procedia computer science, 2018. **134**: p. 310-315.
27. Sahi, A., et al., *An efficient DDoS TCP flood attack detection and prevention system in a cloud environment*. IEEE Access, 2017. **5**: p. 6036-6048.
28. Aamir, M., et al., *Machine learning classification of port scanning and DDoS attacks: a comparative analysis*. Mehran University Research Journal Of Engineering & Technology, 2021. **40**(1): p. 215-229.
29. Xiao, P., et al., *Detecting DDoS attacks against data center with correlation analysis*. Computer Communications, 2015. **67**: p. 66-74.
30. Umarani, S. and D. Sharmila, *Predicting application layer DDoS attacks using machine learning algorithms*. International Journal of Computer and Systems Engineering, 2015. **8**(10): p. 1912-1917.
31. Peraković, D., et al. *Artificial neuron network implementation in detection and classification of DDoS traffic*. in *2016 24th Telecommunications Forum (TELFOR)*. 2016. IEEE.
32. Johnson Singh, K., K. Thongam, and T. De, *Entropy-based application layer DDoS attack detection using artificial neural networks*. Entropy, 2016. **18**(10): p. 350.
33. Fouladi, R.F., C.E. Kayatas, and E. Anarim. *Frequency based DDoS attack detection approach using naive Bayes classification*. in *2016 39th International Conference on Telecommunications and Signal Processing (TSP)*. 2016. IEEE.
34. Yuan, X., C. Li, and X. Li. *DeepDefense: identifying DDoS attack via deep learning*. in *2017 IEEE International Conference on Smart Computing (SMARTCOMP)*. 2017. IEEE.
35. Yusof, M.A.M., F.H.M. Ali, and M.Y. Darus. *Detection and defense algorithms of different types of DDoS attacks using machine learning*. in *International Conference on Computational Science and Technology*. 2017. Springer.
36. Meti, N., D. Narayan, and V. Baligar. *Detection of distributed denial of service attacks using machine learning algorithms in software defined networks*. in *2017 international conference on advances in computing, communications and informatics (ICACCI)*. 2017. IEEE.
37. Yudhana, A., I. Riadi, and F. Ridho, *DDoS classification using neural network and naïve bayes methods for network forensics*. Int. J. Adv. Comput. Sci. Appl, 2018. **9**(11): p. 177-183.

38. Ghanbari, M. and W. Kinsner, *Detecting DDoS Attacks Using Polyscale Analysis and Deep Learning*. International Journal of Cognitive Informatics and Natural Intelligence (IJCINI), 2020. **14**(1): p. 17-34.
39. Raihan-Al-Masud, M. and M.R.H. Mondal, *Data-driven diagnosis of spinal abnormalities using feature selection and machine learning algorithms*. Plos one, 2020. **15**(2): p. e0228422.
40. Elejla, O.E., et al., *Comparison of Classification Algorithms on ICMPv6-Based DDoS Attacks Detection*, in *Computational Science and Technology*. 2019, Springer. p. 347-357.
41. Aamir, M. and S.M.A. Zaidi, *Clustering based semi-supervised machine learning for DDoS attack classification*. Journal of King Saud University-Computer and Information Sciences, 2019.
42. Koay, A., et al. *A new multi classifier system using entropy-based features in DDoS attack detection*. in *2018 International Conference on Information Networking (ICOIN)*. 2018. IEEE.
43. Sharafaldin, I., et al. *Developing realistic distributed denial of service (DDoS) attack dataset and taxonomy*. in *2019 International Carnahan Conference on Security Technology (ICCST)*. 2019. IEEE.
44. Ghogh, B., et al., *Feature selection and feature extraction in pattern analysis: A literature review*. arXiv preprint arXiv:1905.02845, 2019.
45. Divekar, A., et al. *Benchmarking datasets for anomaly-based network intrusion detection: KDD CUP 99 alternatives*. in *2018 IEEE 3rd International Conference on Computing, Communication and Security (ICCCS)*. 2018. IEEE.
46. Li, S., et al. *Independently recurrent neural network (indrnn): Building a longer and deeper rnn*. in *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.
47. Elsayed, M.S., et al. *Ddosnet: A deep-learning model for detecting network attacks*. in *2020 IEEE 21st International Symposium on "A World of Wireless, Mobile and Multimedia Networks"(WoWMoM)*. 2020. IEEE.
48. Wiest, L., *Recurrent Neural Networks-Combination of RNN and CNN*. Published on, 2017. **7**.

49. Clotet, X., J. Moyano, and G. León, *A real-time anomaly-based IDS for cyber-attack detection at the industrial process level of critical infrastructures*. International Journal of Critical Infrastructure Protection, 2018. **23**: p. 11-20.
50. Draper-Gil, G., et al. *Characterization of encrypted and vpn traffic using time-related*. in *Proceedings of the 2nd international conference on information systems security and privacy (ICISSP)*. 2016.
51. Zhao, R., et al., *Learning to monitor machine health with convolutional bi-directional LSTM networks*. Sensors, 2017. **17**(2): p. 273.
52. Sherstinsky, A., *Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network*. Physica D: Nonlinear Phenomena, 2020. **404**: p. 132306.
53. Shiravi, A., et al., *Toward developing a systematic approach to generate benchmark datasets for intrusion detection*. computers & security, 2012. **31**(3): p. 357-374.
54. Tavallaee, M., et al. *A detailed analysis of the KDD CUP 99 data set*. in *2009 IEEE symposium on computational intelligence for security and defense applications*. 2009. IEEE.
55. Rajagopal, S., P.P. Kundapur, and K. Hareesha, *Towards effective network intrusion detection: from concept to creation on Azure cloud*. IEEE Access, 2021. **9**: p. 19723-19742.