

# **Anomaly Detection in Networks Using Machine Learning**

**Umar Zaib**

A report submitted for the Course of Machine Learning

Supervisor: Dr. Arshad  
School of Computing Sciences  
University of PAF-IAST

January 2024

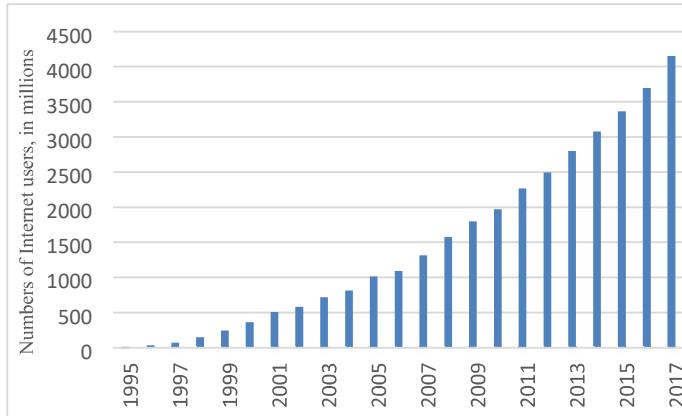
## **Abstract**

Everyday millions of people and hundreds of thousands of institutions communicate with each other over the Internet. In the past two decades, while the number of people using the Internet has increased very fast. Parallel to these developments, the number of attacks made on the Internet is increasing day by day. Although signature-based methods are used to prevent these attacks, they are abortive against zero-day attacks. On the other hand, the Anomaly-based approach is an alternative solution to the network attacks and can detect zero-day attacks as well. In this study, it is aimed to detect network anomaly using machine learning methods. In this context, the CICIDS2017 has been used as dataset because of its up-to datedness, and wide attack diversity. On this dataset, feature selection was made by using the Random Forest Regressor algorithm. Seven different machine learning algorithms have been used in the application step and achieved high performance. Machine learning algorithms and success rates are as follows: Naive Bayes 86%, QDA 86%, Random Forest 94%, ID3 95%, AdaBoost 94%, MLP 83%, and K Nearest Neighbors 97%

# 1 Introduction

## 1.1 Motivation

Everyday millions of people and hundreds of thousands of institutions communicate with each other over the Internet. In the past two decades, while the number of people using the Internet has increased very fast, today this number has exceeded 4 billion and this increase is continuing rapidly [1, 2].



**Figure 1.** Increase in The Number of the Internet Users According to Years (1995-2017) [1, 2].

Parallel to these developments, the number of attacks made on the Internet is increasing day by day. Against these attacks, there are two basic methods used to detect the attacks to ensure information security; identification based on signature, and detection based on anomaly.

In this project, it was aimed to develop a system that detects network anomaly quickly and effectively by means of machine learning methods.

## 1.2 Goals and Objectives

**1.2.1 Goals** The goals that are aimed to achieve at the end of this project are as follows:

- Examination of machine learning algorithms that can be used to detect network anomalies.
- To detect network attacks in a fast and effective way by studying network anomaly with machine learning methods.
- To determine the success level of the study by comparing the results obtained in it with the studies previously conducted in this area.
- Contributing to the literature by obtaining close results from previous studies in the detection of network anomalies.

### 1.2.2 Objectives

The objectives that are aimed to achieve at the end of this study are as follows:

- To examine the previous work done in the field by doing extensive field research.
- Selecting the appropriate dataset by performing comprehensive research on the alternatives to the dataset.
- Choosing suitable algorithms by conducting extensive research on machine learning algorithms.

- Deciding on right algorithms by performing exhaustive research on machine learning methods.
- Selecting the appropriate software platform.
- Choosing the suitable hardware/equipment platform.
- Deciding on the right evaluation criteria.
- Choosing the benchmark studies to be compared during the evaluation phase.

## 2 Background and Related Work

### 2.1 Datasets

In the detection of network anomaly by machine learning methods, there is a need for a large amount of harmful and harmless network traffic for training and testing steps. However, it is not possible for real network traffic to be used publicly because of privacy issues. To meet this

#### CICIDS 2017

Finally, dataset to be addressed is the CICIDS 2017 (Intrusion Detection Evaluation Dataset) created by the Canadian Institute for Cybersecurity at the University of New Brunswick. This dataset consists of a 5-day (3rd July- 7th July 2017) data stream on a network created by computers using up-to-date operating systems such as Windows Vista / 7 / 8.1 / 10, Mac, Ubuntu 12/16 and Kali. Details of the dataset can be seen from Table 1[4, 7].

| Flow Recording Day<br>(Working Hours) | pcap File size | Duration  | CSV File Size | Attack Name  | Flow Count |
|---------------------------------------|----------------|-----------|---------------|--|------------|
| Monday                                | 10 GB          | All Day   | 257 MB        | No Attack  | 529918     |
| Tuesday                               | 10 GB          | All Day   | 166 MB        | FTP-Patator, SSH-Patator   | 445909     |
| Wednesday                             | 12 GB          | All Day   | 272 MB        | DoS Hulk, DoS GoldenEye, DoS slowloris, DoS Slowhttptest, Heartbleed | 692703     |
| Thursday                              | 7.7GB          | Morning   | 87.7 MB       | Web Attacks (Brute Force, XSS, Sql Injection)                        | 170366     |
|                                       |                | Afternoon | 103 MB        | Infiltration   | 288602     |
| Friday                                | 8.2GB          | Morning   | 71.8 MB       | Bot  | 192033     |
|                                       |                | Afternoon | 92.7 MB       | DDoS   | 225745     |
|                                       |                | Afternoon | 97.1 MB       | PortScan   | 286467     |

Table 1. Details of CICIDS2017 Dataset

### 2.3 Attacks

In this section, the types of attacks that the data set contains are examined in detail.

All data in the dataset is tagged with 15 labels. One (Benign) of these tags represents normal network movements while the other 14 represent attacks. The benign record formed using Mail services, SSH, FTP, HTTP, and HTTPS protocols represent a non-harmful / normal data stream on the network, created by simulating real user data.

### 2.3.1 DoS HULK

HULK (HTTP (Hypertext Transfer Protocol) Unbearable Load King) can be used to perform both DoS and DDoS attacks. Because this attack influences the server load first-hand, it can disable the server in a very short time, like a minute or so. HULK has the capability to conceal the real user agent and use a different template for each attack request. During the attack, it creates TCP-SYN (Transmission Control Protocol - Synchronization) floods and multiple HTTP-GET flood requests. To better understand this attack, TCP-SYN flood and HTTP-GET flood methods can be examined.

#### TCP-SYN Flood

The TCP-SYN flood attack is intended to exploit the 3-way handshake method to consume server resources and render the server out of service. The 3-way handshake method involves the following steps:

- The client that wants to connect to the server sends a SYN packet to the server to start the communication.
- The server receiving the SYN packet sends the corresponding SYN-ACK packet.
- Finally, the client initiates communication by sending an ACK packet to the server.

In the TCP-SYN flood attack, the first two steps occur; the attacker (client) sends the SYN packet to the victim (server), and the server responds with the SYN-ACK packet in response. However, the attacker does not send the ACK packet that should be sent, and the connection waits as incomplete. The server reserves a resource by keeping these requests in a log queue. When the number of requests increases, the server becomes inaccessible.

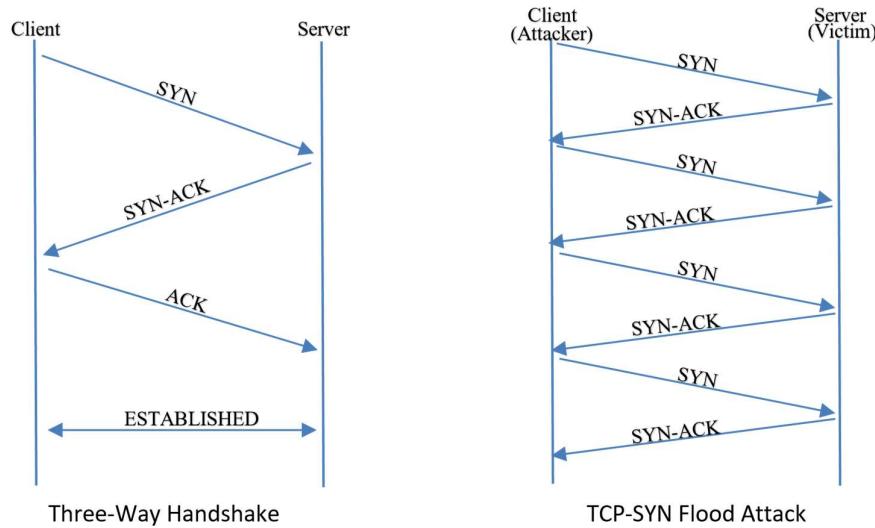
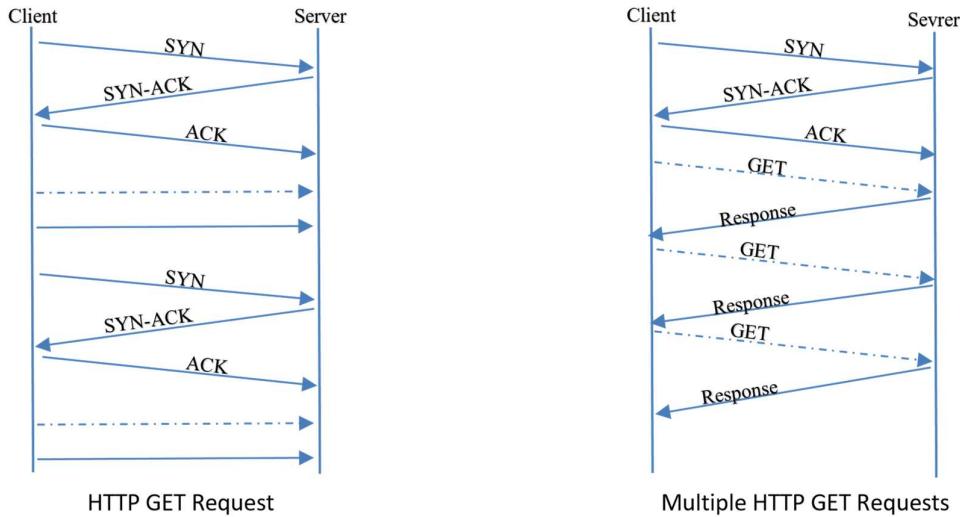


Figure 4. A comparison of TCP-SYN Flood Attack with a successful Three-Way Handshake.

#### HTTP-GET Flood

HTTP GET is the command used by the client to request a file from a web server. The purpose of an HTTP GET flood attack is to reduce the number of HTTP connections that the web server can keep open by consuming server resources, rendering the server unresponsive to legitimate HTTP GET requests. In this type of attack, the attacker sends too many HTTP GET requests to

the target server. The attacker can send an HTTP GET request on each TCP connection or can send multiple HTTP GET requests for each TCP connection using the HTTP feature. Moreover, these requests do not require high network traffic.



**Figure 5.** The comparison of singular and multiple HTTP GET requests in HTTP GET Flood attack [17]

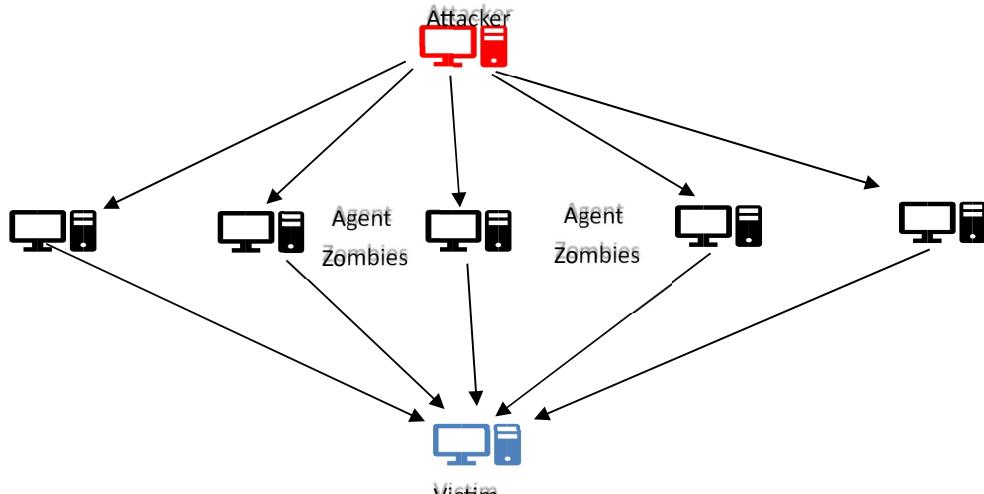
### 2.3.2 PortScan

The PortScan attack in the CICIDS2017 dataset is made by the program named Nmap and is located on Friday afternoon section. Ports are abstract connection points named with numbers between 0 and 65535 and are used to determine what type of service you want to receive from the reached computer. Port Scan is a popular attack type used to gather information about how attackers can get into the system. With this attack, the attacker can learn much critical information about the connected devices, such as operating system, running services, and port status. In a PortScan attack, the attacker sends a message to each port. The type of response received from the victim gives information about whether the port is used or not. By using this information, the attacker learns about the used system and detects its weaknesses. Port Scanning types can be listed as follows:

### DDoS

The DDoS attack, that sent HTTP, TCP and UDP requests. is located in Friday-afternoon section. A DDoS (Distributed Denial-of-Service) attack is another type of attack designed to prevent a legitimate user from accessing a website or web service, such as a DoS attack. Unlike the DoS attack, the DDoS attack is not made from a central computer but is made from many computers to a single source.

Hundreds of proxy computers can be used to attack a computer. These computers used were infected by malicious software so that they could be used remotely during the DDoS attack. These devices are called "zombies", and the network created by these devices is called "botnet".



**Figure 6.** Representation of a DDoS Attack

### FTP-Patator

FTP-Patator attack is made using the Patator, a multithreaded tool written in the Python program. This attack takes place on Tuesday-morning-records within the CICIDS2017 dataset. FTP (File Transfer Protocol) is a network protocol that provides file transfer between client and server in a network. To transfer a file using FTP, a valid username and password are required on the network to be sent. The FTP-Patator attack is a brute-force attack aimed at seizing this username and password.

**Brute-force** is a cryptographic attack designed to obtain passwords. In this attack, the attacker tries to log on to the system as a legitimate user by trying possible usernames and passwords on the system. Software that automatically generates passwords is often used for this process.

A large number of unsuccessful entry attempts over a short period of time is a characteristic feature for brute force attacks. In this context, it is possible to see dense packet flow during brute-force attack. In addition, failed entries do not contain high sizes files, so bandwidth consumption and bytes count are low.

### SSH-Patator

SSH-Patator attack is made using the Patator, a multithreaded tool written in the Python program. This attack takes place on Tuesday-afternoon-records within the CICIDS2017 dataset.

SSH (Secure Shell) is a cryptographic protocol that allows secure operation of various network services over a network in an unsecured environment (e.g. The Internet). The most common use of this protocol is to log on to a remote system.

In this attack, the intent of the attacker is to gain remote access to a system and gain full control over it. This attack consists of three steps:

**Scanning Step:** This is the first phase of the SSH-Patator attack. The goal in this phase is to learn about the system to be attacked (See 2.2.2 PortScan). The attacker tries to find the host

running SSH by performing a scan on a specific port number for IP blocks on a network or subnets.

**Brute-Force Step:** At this stage, the attacker tries to log in to the system he has discovered during the scanning phase with a large number of username and password combinations.

**Die-off Step:** The attacker gains the legitimate user authority by successfully logging in.

While during the first step of this attack, a large number of incomplete TCP packets with SYN flags are observed, in the second step (brute-force), small sized completed TCP packets are seen. The number of packets in the stream is high, but the packet sizes are low.

### **Botnet**

Botnet attack is made using Ares, an attack tool written in the Python program. This attack takes place on Friday-morning-records within the CICIDS2017 dataset.

Botnet is a network created by malware-infected computers. These computers that make up botnets are called bots or zombies. Because of the malicious software, these computers perform various harmful activities without the knowledge of their owners. These computers can be remotely controlled by bot-masters and used to make SPAM (unsolicited emails) or DDoS attacks.

The detection of the bots' network behavior is only possible in the active state (during the attack). When the botnet is passive, there is no network activity and cannot be detected. In the active state, the number of bytes per packet and the number of packets with PSH flag make it possible to detect this attack.

### **Web Attack**

Web Attack takes place on Thursday-afternoon -records in the CICIDS2017 dataset. In this group of attacks, three different web-based attacks were launched against a vulnerable PHP (Hypertext Pre-processor - a script programming language)/ MySQL (an open-source database administration framework) web application identified as the victim. These attacks are:

**XSS:** XSS (Cross-Site Scripting), a popular web-based attack type, is implemented by injecting code into a web page. When the victim wants to view this web page, this malicious code fragment can lead to undesirable results such as data theft, session seizing, special code execution, and fraud.

**SQL Injection:** SQL (Structured Query Language – a database management system) Injection attack is one of the most popular and most dangerous attacks on the web. A typical dynamic web page keeps the user's various information in the database for later use and occasionally makes various queries in the database to reach this data.

These attacks are usually done by exploiting a security vulnerability in the database layer of a web application. The attacker can access the database using the exploits of the connection between the web application and the database. Then he could be able to steal, delete, or modify the data by sending malicious commands to the database server.

## **Heartbleed**

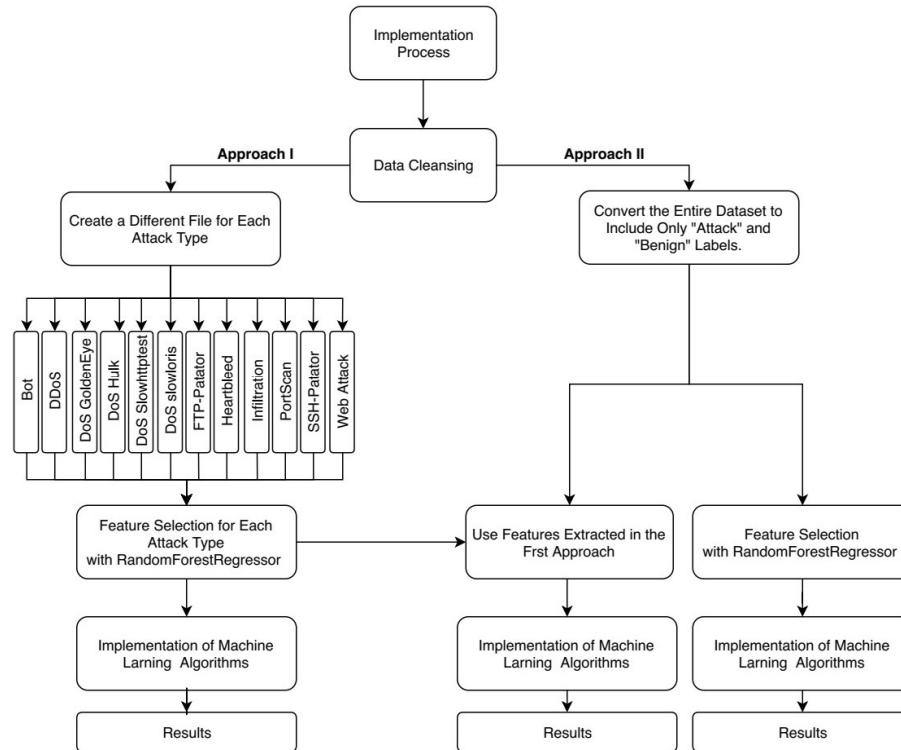
This attack is made using the Heartleech, a Heartbleed exploit tool written in C programming language. This attack takes place on Wednesday - afternoon -records within the CICIDS2017 dataset. The normal operation of heartbeat is as follows: The client transmits to the server a request consisting of random payload. the server replies to this request with a packet containing the same payload.

During the attack, a specially created heartbeat request is sent to the server. This demand is empty but points out that it carries very high amounts of data. At this stage, due to a vulnerability in OpenSSL, the server sends a piece of memory in the specified length in response to this message. These parts contain various confidential information such as personal information, usernames and passwords that should not be sent. The size of these packages can be up to 16 KB and the attacker can repeat this operation without any limit.

During this attack, some anomalies on the network flow are observed. The length of the incoming messages is less than the heartbeat minimum message size of 20 bytes. Outgoing message lengths are in the kilobyte level, and for a heartbeat response, this size is too big. The difference in size between outgoing and incoming messages is another way to understand the attack. In normal heartbeat messages, the lengths of the outgoing and incoming messages are the same. In case of attack, incoming message lengths are too small, outgoing message lengths are too large.

### 3.2 Implementation

In this section, various pre-processing and actual applications are performed to detect anomaly by machine learning techniques. For this purpose, the data cleansing process is performed in the first step and the dataset is cleaned from mistakes and defects. Then, the data set is divided into two parts, training, and test. After these operations, the properties to be used by the algorithms are decided at the step of feature selection. Finally, the section ends with the implementation of machine learning algorithms. Figure 12 illustrates the implementation process in detail.



**Figure 12.** The Implementation Process

### Data Cleansing

It may be necessary to make some changes to the dataset before using it in practice, making it more efficient. For this purpose, in this section, some defects of the CICIDS2017 dataset are corrected, and some data are edited.

The dataset file contains 3119345 stream records. The distribution of these stream records can be seen from Table 2. When these records are examined, the 288602 record is incorrect / incomplete. The first step in the pre-processing process will be to delete these unnecessary records.

| <b>Label Name</b>          | <b>Number</b> |
|----------------------------|---------------|
| Benign                     | 2359289       |
| Faulty                     | 288602        |
| DoS Hulk                   | 231073        |
| PortScan                   | 158930        |
| DDoS                       | 41835         |
| DoS GoldenEye              | 10293         |
| FTP-Patator                | 7938          |
| SSH-Patator                | 5897          |
| DoS slowloris              | 5796          |
| DoS Slowhttptest           | 5499          |
| Bot                        | 1966          |
| Web Attack – Brute Force   | 1507          |
| Web Attack – XSS           | 652           |
| Infiltration               | 36            |
| Web Attack – SQL Injection | 21            |
| Heartbleed                 | 11            |

**Table 2.** Distribution of stream records in the CICIDS2017 dataset.

Another error about the dataset is in the columns that make up the features. The dataset file consists of 86 columns that define the flow properties such as Flow ID, Source IP, Source Port etc. However, the Fwd Header Length feature (which defines the forward direction data flow for total bytes used) was written two times (41st and 62nd columns). This error is corrected by deleting the repeating column (column 62).

However, although the “Label” tag is a categorical feature, no changes have been made on it. The reason is that during the processing, the original categories are needed in order to classify the attack types in different forms and to try different approaches.

### 3.2.2 Creation of Training and Test Data

During the machine learning process, data is needed so that learning can take place. The data sets used are the result of this need. In addition to the data required for training, test data is needed to evaluate the performance of the algorithm and to see how well it works. The algorithm acquires a skill on the training data and applies it to the test data. The result of the test data is the performance of the machine learning algorithm [40].

However, the CICIDS2017 dataset used does not contain dedicated training and test data, but it contains a single unbundled dataset. Therefore, the data should be divided into training and test data parts. In the application phase, a Sklearn command, train\_test\_split[59] is used. This command divides the data into 2 parts at the sizes specified by the user. Generally preferred partitioning is 20% test, 80% training data [40] and this ratio is also preferred in this application. The train\_test\_split command makes the selection random when creating data groups. This process is known as cross-validation. In order to ensure that the results obtained

during the application are solid, the creation of the training and the test data have been performed 10 times in succession. The results obtained are the arithmetic mean of the repeated operations.

### 3.2.3 Feature Selection

In this section, the features in the dataset are evaluated to determine which features are important to define which attack.

#### 3.2.3.1 Feature Selection According to Attack Types

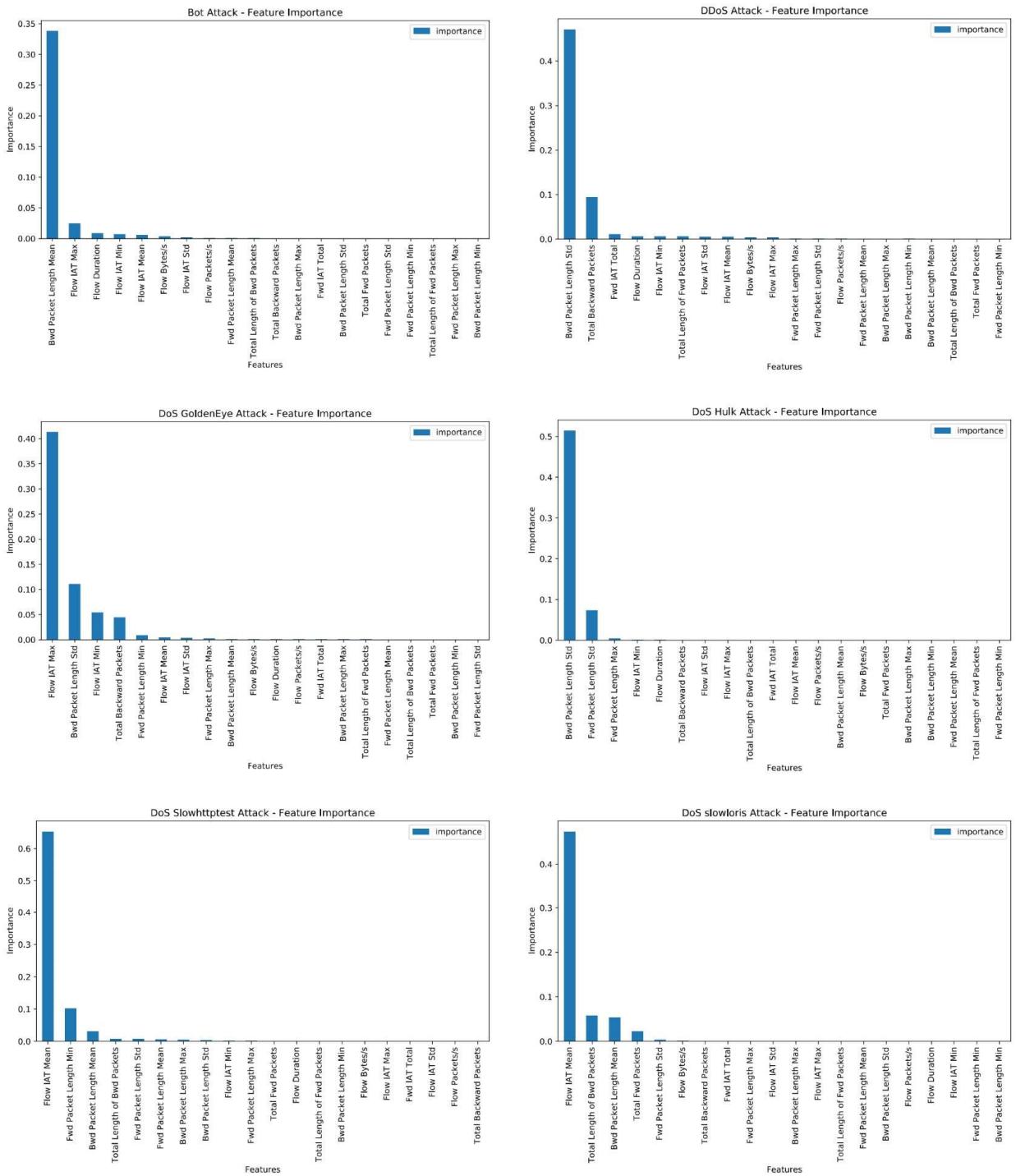
To do this calculation, a special file is created for every kind of attack by isolating the attack from other attacks. This file contains the entire stream identified as the attack and the data stream identified as randomly selected "Benign" (30% Attack, 70% Benign).

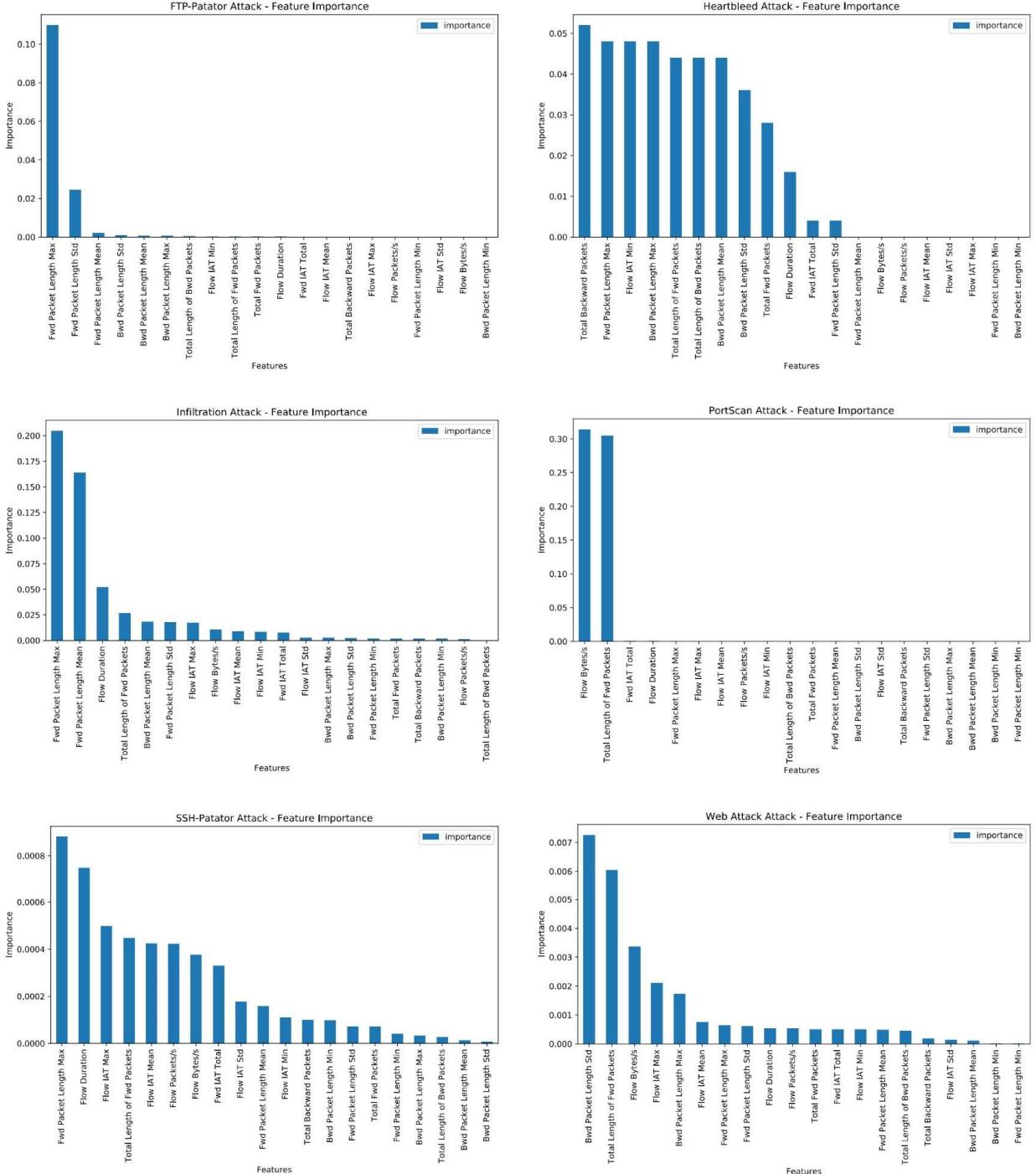
The Random Forest Regressor class of Sklearn is used when importance weights of features are calculated. This algorithm creates a decision-forest. In this decision forest, each feature is given a weight of importance as to how useful they are in the construction of the decision-tree. When the process is finished, these importance weights of features are compared and sorted. The sum of the importance weights of all the properties gives the total importance weight of the decision tree. The comparison of the score of any feature to the score of the whole tree gives information about the importance of that feature in the decision tree.

The distribution of features with the most significance value can be seen below:

|   |   |   |
|---|---|---|
| <b>Bot attack importance list:</b>  | <b>Dos Slowhttptest attack importance list:</b>   | <b>Infiltration attack importance list:</b>   |
| <b>Features</b>   | <b>Features</b>   | <b>Features</b>   |
| 1 Bwd Packet Length Mean<br>2 Flow IAT Max<br>3 Flow Duration<br>4 Flow IAT Min<br>5 Flow IAT Mean<br>6 Flow Bytes/s<br>7 Flow Std<br>8 Flow Packets/s<br>9 Fwd Packet Length Mean<br>10 Total Length of Bwd Packets<br>11 Total Backward Packets<br>12 Bwd Packet Length Max<br>13 Fwd IAT Total<br>14 Bwd Packet Length Std<br>15 Total Fwd Packets<br>16 Fwd Packet Length Std<br>17 Fwd Packet Length Min<br>18 Total Length of Fwd Packets<br>19 Fwd Packet Length Max<br>20 Fwd Packet Length Min     | 1 Flow IAT Mean<br>2 Fwd Packet Length Min<br>3 Bwd Packet Length Mean<br>4 Total Length of Bwd Packets<br>5 Fwd Packet Length Std<br>6 Fwd Packet Length Mean<br>7 Bwd Packet Length Max<br>8 Bwd Packet Length Std<br>9 Fwd IAT Min<br>10 Fwd Packet Length Max<br>11 Total Fwd Packets<br>12 Flow Duration<br>13 Total Length of Fwd Packets<br>14 Bwd Packet Length Min<br>15 Flow Bytes/s<br>16 Flow IAT Max<br>17 Fwd IAT Total<br>18 Flow IAT Std<br>19 Flow Packets/s<br>20 Total Backward Packets        | 1 Fwd Packet Length Max<br>2 Fwd Packet Length Mean<br>3 Flow Duration<br>4 Total Length of Fwd Packets<br>5 Bwd Packet Length Mean<br>6 Fwd Packet Length Std<br>7 Flow IAT Max<br>8 Flow Bytes/s<br>9 Flow IAT Mean<br>10 Flow IAT Min<br>11 Fwd IAT Total<br>12 Flow IAT Std<br>13 Bwd Packet Length Max<br>14 Bwd Packet Length Std<br>15 Fwd Packet Length Min<br>16 Total Fwd Packets<br>17 Total Backward Packets<br>18 Bwd Packet Length Min<br>19 Flow Packets/s<br>20 Total Length of Bwd Packets |
| <b>DoS attack importance list:</b>  | <b>Dos slowloris attack importance list:</b>  | <b>PortScan attack importance list:</b>   |
| <b>Features</b>   | <b>Features</b>   | <b>Features</b>   |
| 1 Bwd Packet Length Std<br>2 Total Backward Packets<br>3 Fwd IAT Total<br>4 Flow Duration<br>5 Flow IAT Min<br>6 Total Length of Fwd Packets<br>7 Flow IAT Std<br>8 Flow IAT Mean<br>9 Flow Bytes/s<br>10 Flow IAT Max<br>11 Fwd Packet Length Max<br>12 Fwd Packet Length Std<br>13 Flow Packets/s<br>14 Fwd Packet Length Mean<br>15 Bwd Packet Length Max<br>16 Bwd Packet Length Min<br>17 Bwd Packet Length Mean<br>18 Total Length of Bwd Packets<br>19 Total Fwd Packets<br>20 Fwd Packet Length Min | 1 Flow IAT Mean<br>2 Total Length of Bwd Packets<br>3 Bwd Packet Length Mean<br>4 Total Fwd Packets<br>5 Fwd Packet Length Std<br>6 Flow Bytes/s<br>7 Total Backward Packets<br>8 Fwd IAT Total<br>9 Fwd Packet Length Max<br>10 Flow IAT Std<br>11 Bwd Packet Length Max<br>12 Flow IAT Max<br>13 Total Length of Fwd Packets<br>14 Fwd Packet Length Mean<br>15 Bwd Packet Length Std<br>16 Flow Packets/s<br>17 Flow Duration<br>18 Flow IAT Min<br>19 Fwd Packet Length Min<br>20 Bwd Packet Length Min       | 1 Flow Bytes/s<br>2 Total Length of Fwd Packets<br>3 Fwd IAT Total<br>4 Flow Duration<br>5 Fwd Packet Length Max<br>6 Flow IAT Max<br>7 Flow IAT Mean<br>8 Flow Packets/s<br>9 Flow IAT Min<br>10 Total Length of Bwd Packets<br>11 Total Fwd Packets<br>12 Fwd Packet Length Mean<br>13 Bwd Packet Length Std<br>14 Flow IAT Std<br>15 Total Backward Packets<br>16 Fwd Packet Length Std<br>17 Bwd Packet Length Max<br>18 Bwd Packet Length Mean<br>19 Bwd Packet Length Min<br>20 Fwd Packet Length Min |
| <b>Dos GoldenEye attack importance list:</b>  | <b>FTP-Patator attack importance list:</b>  | <b>SSH-Patator attack importance list:</b>  |
| <b>Features</b>   | <b>Features</b>   | <b>Features</b>   |
| 1 Flow IAT Max<br>2 Bwd Packet Length Std<br>3 Flow IAT Min<br>4 Total Backward Packets<br>5 Fwd Packet Length Min<br>6 Flow IAT Mean<br>7 Flow IAT Std<br>8 Fwd Packet Length Max<br>9 Bwd Packet Length Mean<br>10 Flow Bytes/s<br>11 Flow Duration<br>12 Flow Packets/s<br>13 Fwd IAT Total<br>14 Bwd Packet Length Max<br>15 Total Length of Fwd Packets<br>16 Fwd Packet Length Mean<br>17 Total Length of Bwd Packets<br>18 Total Fwd Packets<br>19 Bwd Packet Length Min<br>20 Fwd Packet Length Std | 1 Fwd Packet Length Max<br>2 Fwd Packet Length Std<br>3 Fwd Packet Length Mean<br>4 Total Length of Bwd Packets<br>5 Bwd Packet Length Mean<br>6 Bwd Packet Length Max<br>7 Total Length of Bwd Packets<br>8 Flow IAT Min<br>9 Total Length of Fwd Packets<br>10 Total Fwd Packets<br>11 Flow Duration<br>12 Fwd IAT Total<br>13 Flow IAT Mean<br>14 Total Backward Packets<br>15 Flow IAT Max<br>16 Flow Packets/s<br>17 Fwd Packet Length Min<br>18 Flow IAT Std<br>19 Flow Bytes/s<br>20 Bwd Packet Length Min | 1 Fwd Packet Length Max<br>2 Flow Duration<br>3 Flow IAT Max<br>4 Total Length of Fwd Packets<br>5 Flow IAT Mean<br>6 Flow Packets/s<br>7 Flow Bytes/s<br>8 Fwd IAT Total<br>9 Flow IAT Std<br>10 Fwd Packet Length Mean<br>11 Flow IAT Min<br>12 Total Backward Packets<br>13 Bwd Packet Length Min<br>14 Fwd Packet Length Std<br>15 Total Fwd Packets<br>16 Fwd Packet Length Min<br>17 Bwd Packet Length Max<br>18 Total Length of Bwd Packets<br>19 Bwd Packet Length Mean<br>20 Bwd Packet Length Std |
| <b>Dos Hulk attack importance list:</b>   | <b>Heartbleed attack importance list:</b>   | <b>Web Attack attack importance list:</b>   |
| <b>Features</b>   | <b>Features</b>   | <b>Features</b>   |
| 1 Bwd Packet Length Std<br>2 Fwd Packet Length Std<br>3 Fwd Packet Length Max<br>4 Flow IAT Min<br>5 Flow Duration<br>6 Total Backward Packets<br>7 Flow IAT Std<br>8 Flow IAT Max<br>9 Total Length of Bwd Packets<br>10 Fwd IAT Total<br>11 Flow IAT Mean<br>12 Flow Packets/s<br>13 Bwd Packet Length Mean<br>14 Flow Bytes/s<br>15 Total Fwd Packets<br>16 Bwd Packet Length Max<br>17 Bwd Packet Length Min<br>18 Fwd Packet Length Mean<br>19 Total Length of Fwd Packets<br>20 Fwd Packet Length Min | 1 Total Backward Packets<br>2 Fwd Packet Length Max<br>3 Flow IAT Min<br>4 Bwd Packet Length Max<br>5 Total Length of Fwd Packets<br>6 Total Length of Bwd Packets<br>7 Bwd Packet Length Mean<br>8 Bwd Packet Length Std<br>9 Total Fwd Packets<br>10 Flow Duration<br>11 Fwd IAT Total<br>12 Fwd Packet Length Std<br>13 Fwd Packet Length Mean<br>14 Flow Bytes/s<br>15 Flow Packets/s<br>16 Flow IAT Mean<br>17 Flow IAT Std<br>18 Flow IAT Max<br>19 Fwd Packet Length Min<br>20 Bwd Packet Length Min       | 1 Bwd Packet Length Std<br>2 Total Length of Fwd Packets<br>3 Flow Bytes/s<br>4 Flow IAT Max<br>5 Bwd Packet Length Max<br>6 Flow IAT Mean<br>7 Fwd Packet Length Max<br>8 Fwd Packet Length Std<br>9 Flow Duration<br>10 Flow Packets/s<br>11 Total Fwd Packets<br>12 Fwd IAT Total<br>13 Flow IAT Min<br>14 Fwd Packet Length Mean<br>15 Total Length of Bwd Packets<br>16 Total Backward Packets<br>17 Flow IAT Std<br>18 Bwd Packet Length Mean<br>19 Bwd Packet Length Min<br>20 Fwd Packet Length Min |

The distribution of features with the most significance can be seen below:





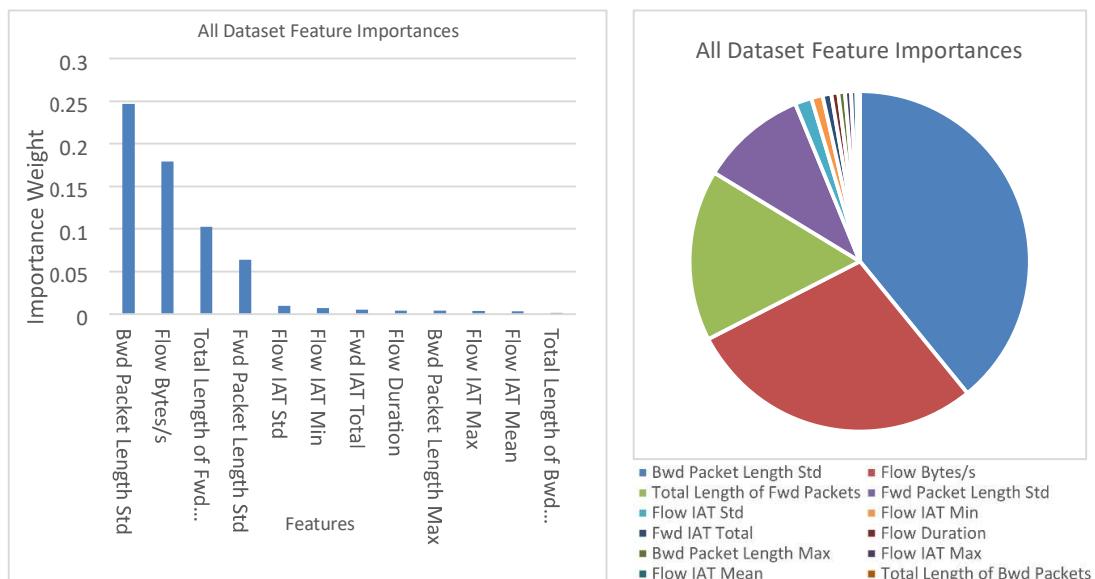
### Feature Selection According to Attack or Benign

Another approach in feature selection is to apply the Random Forest Regressor operation to the whole dataset by collecting all attack types under a single label; “attack”. So, the data in this file contains only the attack and benign tags. As a result of this operation, the feature list obtained is shown in Table 4 and graphics of feature in Figure 16.

| Feature Name          | Importance Weight | Feature Name  | Importance Weight |
|-----------------------|-------------------|---------------|-------------------|
| Bwd Packet Length Std | 0.246627          | Flow IAT Mean | 0.003266          |

|                             |          |  |                             |          |
|-----------------------------|----------|--|-----------------------------|----------|
| Flow Bytes/s                | 0.178777 |  | Total Length of Bwd Packets | 0.001305 |
| Total Length of Fwd Packets | 0.102417 |  | Fwd Packet Length Min       | 0.000670 |
| Fwd Packet Length Std       | 0.063889 |  | Bwd Packet Length Mean      | 0.000582 |
| Flow IAT Std                | 0.009898 |  | Flow Packets/s              | 0.000541 |
| Flow IAT Min                | 0.006946 |  | Fwd Packet Length Mean      | 0.000526 |
| Fwd IAT Total               | 0.005121 |  | Total Backward Packets      | 0.000169 |
| Flow Duration               | 0.004150 |  | Total Fwd Packets           | 0.000138 |
| Bwd Packet Length Max       | 0.004007 |  | Fwd Packet Length Max       | 0.000125 |
| Flow IAT Max                | 0.003579 |  | Bwd Packet Length Min       | 0.000084 |

**Table 4.** According Attack and Benign Labels Feature Importance Weight List



**Figure 16.** Graphs of Feature Importance Weights According to Attack and Benign Labels

## Implementation of Machine Learning Algorithms

### Naïve Bayes

Naïve Bayes is a machine learning algorithm that is simplified with the addition of the independence condition on Bayes' theorem.

Bayes' theorem is expressed by the following equation:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \quad 2.1$$

$P(A | B)$ : The probability of occurrence of A in the case of B occurrences.

$P(B | A)$ : The probability of occurrence of B in the case of A occurrences.  $P(A)$  and  $P(B)$ : Prior probabilities of A and B.

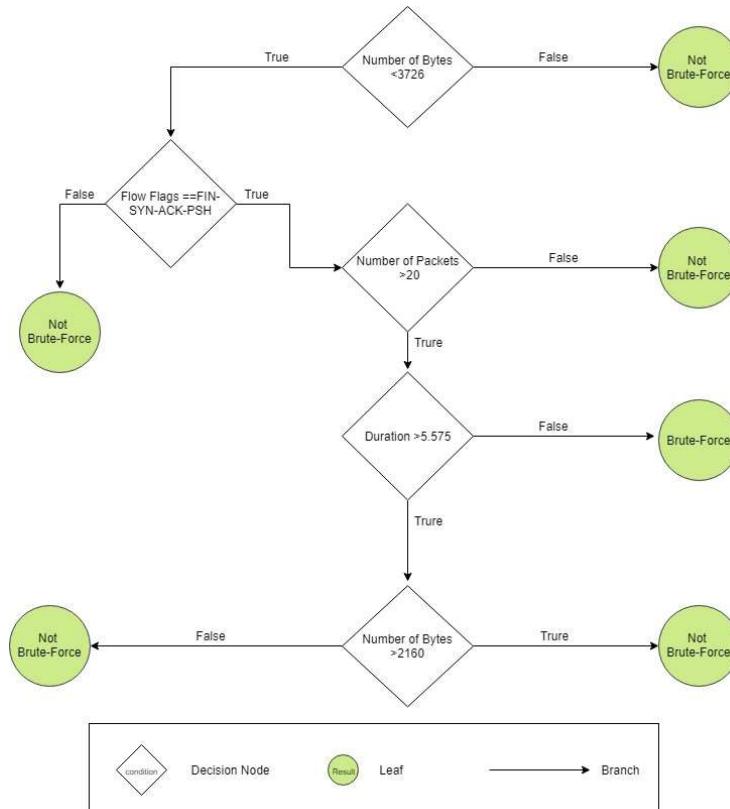
Naive Bayes is a network of probabilities consisting of a parent node representing the unobserved state and multiple child nodes representing the observed states.

In this statistical network, it is assumed that the child nodes are independent of each other. Although this hypothesis is the basis of Naive Bayes' theory, the assumption that sub nodes are independent of one another is often not correct. This is the underlying reason for the Naive Bayes method to achieve lower accuracy when compared to other machine learning methods. Despite this disadvantage, it is one of the most preferred machine learning methods because the training period is very short, and the computational cost is very low.

### Decision Trees

Decision trees are one of the popular classifiers used in machine learning methods. In this approach, the rules used are straightforward and understandable.

Each decision tree consists of nodes (root-node and sub-nodes), branches and leaves. Within each node, there is a decision statement. According to the result of this decision, the algorithm chooses one of the two branches in the next step (the number of branches may be more than two in some sub-algorithms). This selected branch takes the algorithm to the next node. This process ends with the last element, the leaf.



**Figure 7.** Detection of brute-force attack using decision trees.

Decision trees method applies the divide-and-conquer strategy. It makes very large and meaningless data smaller and group them into a meaningful one. It is therefore widely used for the classification of large and complex data.

Unfortunately, these complex trees lead to overfitting<sup>1</sup> and prevent getting fair results. Another disadvantage is that a slight misclassification of the data at the beginning of decomposition can lead to different branches and misleading results. To cope with this, usually, the data to be used must go through pre-processing.

In the implementation phase, a decision tree algorithm, ID3(Iterative Dichotomiser 3), is used. This algorithm is suitable for situations where the training set contains many features. It also stands out with its remarkable features such as giving a reasonable value without doing too much computation and connecting more than two branches to decision nodes.

### **Random Forest**

Random forest is a machine learning approach that uses decision trees. In this method, a "forest" is created by assembling a large number of different decision tree structures which are formed in different ways.

Random forest generates N decision trees from training set data. During this process, it randomly resamples the training set for each tree. Thus, n decision trees are obtained, each of which is different from the other. Finally, voting is performed by selecting new estimates from estimates made by N trees. The value with the highest rating is determined as the final value.

The random forest has many advantages. These can be listed as follows:

- This algorithm can work well with very large and complicated datasets.
- The overfitting problem frequently encountered by decision trees is very rare in this algorithm.
- It can be applied to a lot of kind of machine learning problem.
- It is also good to deal with missing values in the data set. It replaces these lost values with their own created values.
- In addition, it calculates and uses the importance level of the variables when making the classification. Thanks to this, it is also used for **feature selection** in the machine learning.

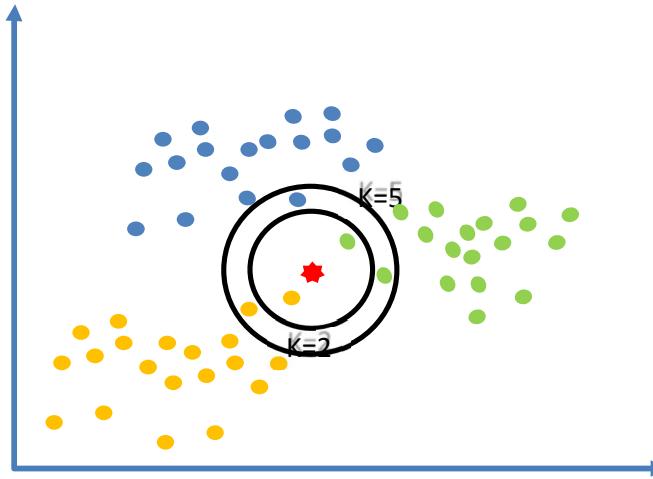
On the other hand, the structure of Random Forest is quite complicated because it is made up of many decision trees. Another disadvantage is that It is pretty difficult to understand its functioning.

### **K Nearest Neighbor**

KNN (K Nearest Neighbor), which is a sample-based method, is one of the most used machine learning algorithms with its simple and fast structure. This algorithm depends on the assumption that the examples in a dataset will exist close to the examples with similar properties in another dataset.

In this context, KNN identifies the class of new data that is not classifiable by using training data of known class type. This determination is made by observing the nearest neighbors of the new sample, for which no classifications are specified.

In a plane with N properties, the number of neighbors to be looked at for an unclassified sample is specified by the number K. For the unknown sample, the distances to the neighbors are calculated and the smallest K numbers are chosen from these distance values. The most repeated property within the K values is assigned as the unknown instance property. In Figure 5, a visual is created for the values 2 and 5 of K in a 3-dimensional plane ( $N = 3$ ).



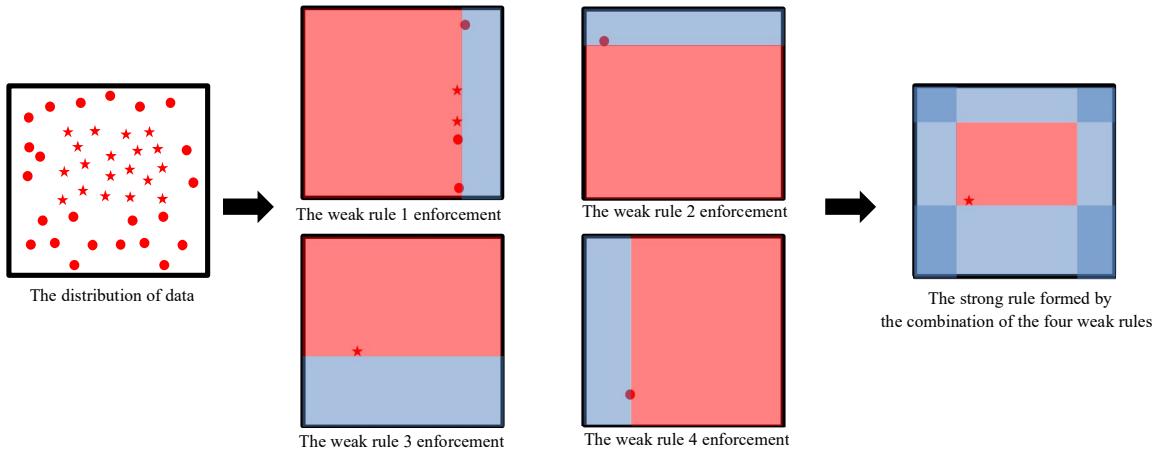
**Figure 8.** Operation of KNN algorithm for  $K = 2$  and  $K = 5$  values.

KNN, which provides good performance over multidimensional data and is a fast algorithm during the training phase, is relatively slow in the estimation stage [43].

### **AdaBoost**

AdaBoost (Adaptive Boosting), a boosting method, is a machine learning algorithm developed to improve classification performance. The basic working principle of Boosting algorithms can be explained as follows: The data are first divided into groups with rough draft rules. Whenever the algorithm is run, new rules are added to this rough draft rules. In this way, many weak and low performance rules called "basic rules" are obtained.

Once the algorithm has been working many times, these weak rules are combined into a single rule that is much stronger and more successful. During this process, the Algorithm assigns a weighting coefficient to each weak rule, giving the highest coefficient value to the lowest error rate. These weight values come into play when final rules are selected. The final rule is created by giving priority to the high scored weak rules.



**Figure 9.** Demonstration of the operation of the AdaBoost algorithm.

AdaBoost algorithm has many advantages. These can be listed as follows:

- There is no need for variable transformation to use it in this algorithm.
- It can perform operations on too many weak rules.
- The overfitting problem very rare in AdaBoost.
- It is also good to deal with missing values in the data set.

On the other hand, it can be mentioned as disadvantages that it is weak against noise and extreme values, and that the predictive value does not have the highest values when compared to other algorithms.

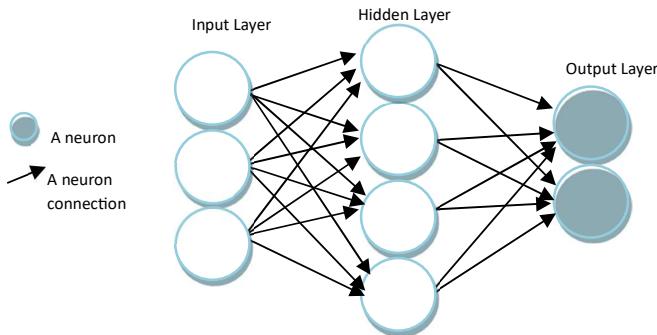
### MLP

MLP (Multi-Layer Perceptron) is a genre of artificial neural networks. Artificial neural networks (ANN) is a machine learning method that takes inspiration from the way the human brain works. The intention of this method is to imitate the properties of the human brain, such as learning, decision making, and deriving new information. While the human brain is made up of interconnected cells called neurons, artificial neural networks are made up of interconnected hierarchical artificial cells.

MLP consists of three stages. These stages are the Input layer, the hidden layer, and the output layer. The input layer is the stage of the MLP that is responsible for receiving data. No information processing is performed on this layer. Only the received information is transmitted to the next layer, the hidden layer. Each neuron in this step bonds with all the neurons in the hidden layer.

In the hidden layer, the data sent from the input layer is processed and transmitted to the next layer, the output layer. In MLP, the number of hidden layers or the number of artificial neurons in this layer may change. For example, the MLP used in this study has a 3-step hidden layer and 13 neurons in each layer. The change in the layer count and neuron number directly affects MLP performance and yield. By increasing the number of layers and neurons, MLP can deal with more complex problems. However, the increase in the number of these elements will also lead to an increase in the duration of the algorithm.

In the output layer, which is the last layer, each cell is tied to all the cells in the hidden layer, and the results of the processed data in the hidden layer are served at this stage.



**Figure 10.** Demonstration of a three-layer MLP

The advantages MLP provides are as follows:

- It is good at coping with complicated problems.
- It can work with missing data.
- It can generalize after the learning process. Thus, it serves a wider area than the other machine learning algorithms.

But from the other side:

- It is difficult to build the network structure.
- The user should decide the appropriate network structure.
- Overfitting problems may be encountered.
- It is difficult to interpret and understand.

### QDA

QDA (Quadratic Discriminant Analysis) is a discriminant analysis method. Discriminant Analysis is a statistical technique for assigning a measured data to one group among many groups. When this assignment is made, the observed data must be assigned to the group to which it belongs. If it is assigned a group that does not belong to it, an error occurs.

This mistake is known as "the error rate". The purpose of discriminant analysis is to perform the assignment process with a minimum error rate. Assuming that the data group has normal distribution and the variances are equal, the analysis is called Linear Discriminant Analysis. However, if the assumption of equality of groups' variances is not correct, a Quadratic Discriminant Analysis is obtained.

In order to be able to apply the Quadratic Discriminant Analysis, the number of samples observed must be greater than the number of groups.

Two different approaches have been used to apply machine learning algorithms to the dataset. In the first method, the files created in the [Feature Selection](#) section and the attributes obtained in the same section are used. These files contain 30% attack and 70% benign data and each of

them named by the type of attack it contains. The seven machine learning methods are applied to each file 10 times, resulting in a separate outcome for each attack type. With this method, it is aimed to observe the effectiveness and performance of different machine learning methods on different attack types.

In the second approach, the entire data set is used as a single file. All attacks contained in this file are collected under a single common name; "attack". So, the data in this file contains only the attack and benign tags. The set of features to be used consists of combining the 4 features with the highest importance-weight achieved for each attack in approach 1 under a single roof. Thus, 4 features are obtained from each of the 12 attack types, resulting in a pool of features consisting of 48 attributes. After the repetitions are removed, the number of features is 18. The list of these features can be seen in Table 5

|                        |                        |                             |
|------------------------|------------------------|-----------------------------|
| Bwd Packet Length Max  | Flow IAT Mean          | Fwd Packet Length Min       |
| Bwd Packet Length Mean | Flow IAT Min           | Fwd Packet Length Std       |
| Bwd Packet Length Std  | Flow IAT Std           | Total Backward Packets      |
| Flow Bytes/s           | Fwd IAT Total          | Total Fwd Packets           |
| Flow Duration          | Fwd Packet Length Max  | Total Length of Bwd Packets |
| Flow IAT Max           | Fwd Packet Length Mean | Total Length of Fwd Packets |

Table 5. The feature list created for all attack types.

## Results and Discussion

In this section, the results of the project done in the implementation section are presented. In this context, in the assessment carried out, the evaluation criteria are presented via the data of the F-measure.

The performance evaluation procedures are repeated 10 times for each machine learning algorithm. The numbers given in the tables are the arithmetic mean of these 10 processes. Box and whisker graphs are created to illustrate the consistency of the results and the change between them.

### 4.1 Approach 1 - Using 12 Attack Types

Seven different machine learning methods are applied to 12 different attack types and the obtained results are presented in Table 7.

Among the results presented in Table 7, the biggest and smallest achievements are emphasized as follows: The greatest scores are bold, the smallest scores are underlined and italics. In the results of the algorithms, if there is an equality in F-measure, the following values are examined in order to eliminate equality, respectively: accuracy, precision, recall, and time (In Table 7, only F-Measure values are given.).

| Attack Names  | F-Measures  |      |      |             |             |             |             |
|---------------|-------------|------|------|-------------|-------------|-------------|-------------|
|               | NB          | RF   | KNN  | ID3         | AB          | MLP         | QDA         |
| Bot           | <u>0.54</u> | 0.96 | 0.95 | 0.96        | <b>0.97</b> | 0.64        | 0.68        |
| DDoS          | 0.77        | 0.96 | 0.92 | <b>0.96</b> | 0.96        | 0.76        | <u>0.34</u> |
| DoS GoldenEye | 0.81        | 0.99 | 0.98 | <b>0.99</b> | 0.99        | <u>0.64</u> | 0.71        |

|                  |             |      |      |             |             |             |      |
|------------------|-------------|------|------|-------------|-------------|-------------|------|
| DoS Hulk         | <u>0.23</u> | 0.93 | 0.96 | <b>0.96</b> | 0.96        | 0.95        | 0.36 |
| DoS Slowhttptest | <u>0.35</u> | 0.98 | 0.99 | 0.98        | <b>0.99</b> | 0.78        | 0.38 |
| DoS slowloris    | <u>0.37</u> | 0.95 | 0.95 | <b>0.96</b> | 0.95        | 0.74        | 0.46 |
| FTP-Patator      | <b>1.00</b> | 1.00 | 1.00 | 1.00        | 1.00        | 1.00        | 1.00 |
| Heartbleed       | <b>1.00</b> | 0.99 | 1.00 | 0.95        | 0.93        | <u>0.66</u> | 1.00 |
| Infiltration     | 0.78        | 0.92 | 0.88 | 0.89        | <b>0.92</b> | <u>0.52</u> | 0.83 |
| PortScan         | <u>0.39</u> | 1.00 | 1.00 | <b>1.00</b> | 1.00        | 0.61        | 0.85 |
| SSH-Patator      | <u>0.33</u> | 0.96 | 0.95 | <b>0.96</b> | 0.96        | 0.83        | 0.41 |
| Web Attack       | 0.74        | 0.97 | 0.93 | <b>0.97</b> | 0.97        | <u>0.60</u> | 0.84 |

Table 7. Distribution of results according to type of attack and machine learning algorithm.

When looking at the results, it is noticed that Random Forest, KNN, ID3 and Adaboost algorithms, have achieved over 90% success in detecting almost all attack types. Among these four algorithms, ID3, which is the most successful, has completed 7 of 12 tasks with the highest score. In fact, in the 6 of these 7 tasks (DDoS, DoS GoldenEye, DoS Hulk, PortScan, SSHPatator and Web Attack) ID3 shares the highest score with at least one algorithm. However, low processing time puts it ahead of other algorithms.

Another point of interest is that all machine learning algorithms achieve a full score in the FTPPatator attack. This could be due to the fact that the features used to describe the FTP-Patator attack may become a characteristic that can be easily distinguished between normal and attack data, trapped in a very narrow range relative to normal data (See Figure 17).

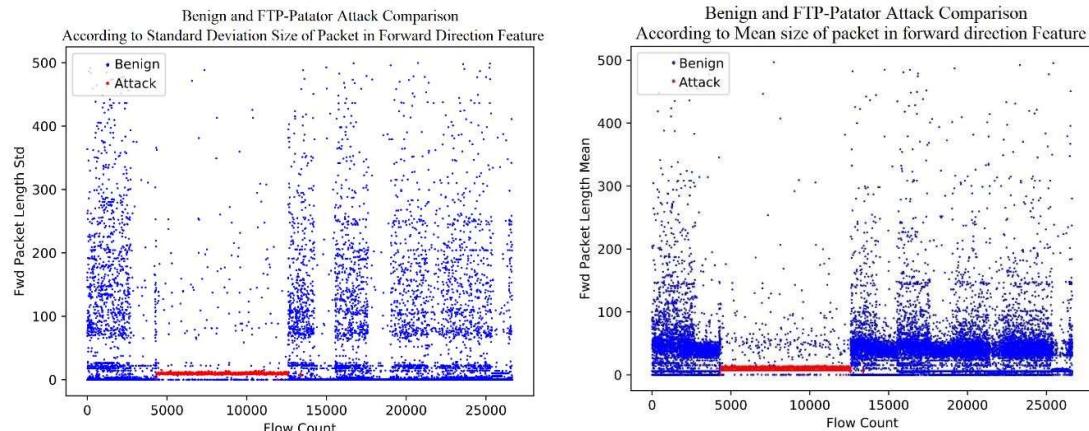
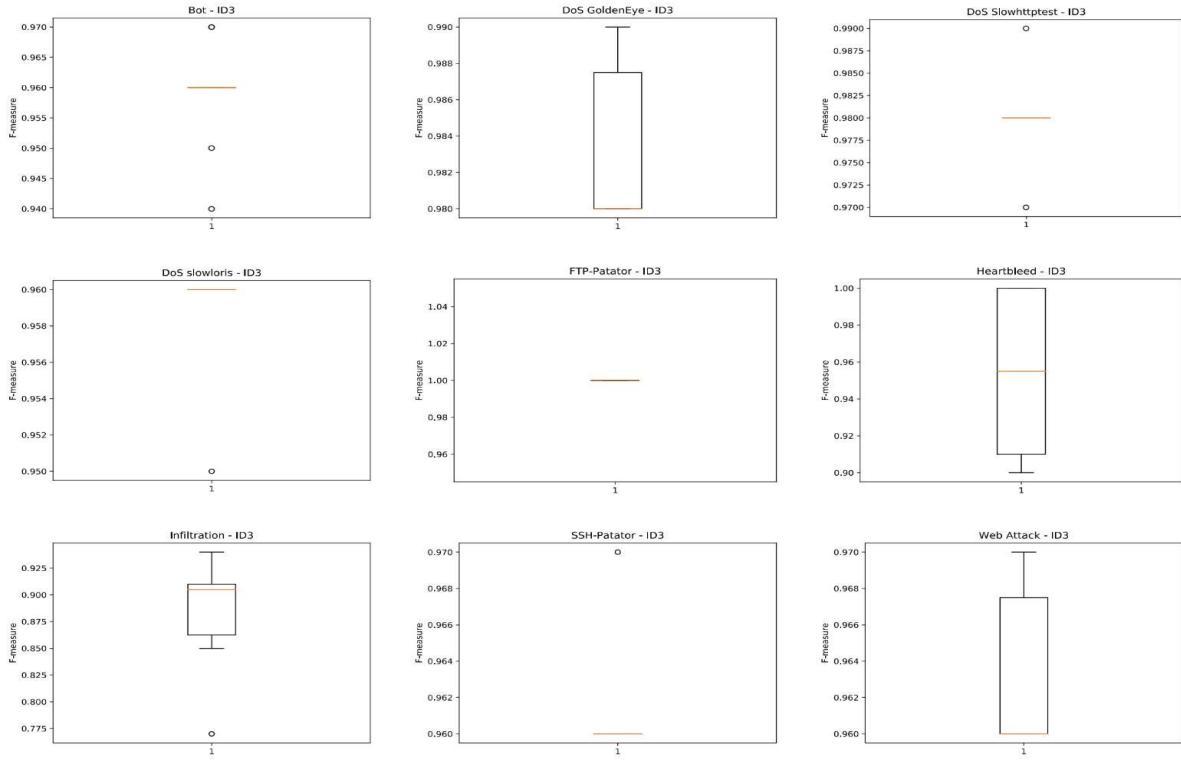


Figure 17. "Fwd Packet Length Std" and "Fwd Packet Length Mean" features in FTP-Patator and Benign flow.



**Figure 19.** Box and whisker graphics containing the results of applying the ID3 algorithm to various attack types.

#### 4.2 Approach 2 - Using Two Groups: Attack and Benign

In this section, the entire data set is used as a single dataset file. All attacks contained in this file are collected under a single common name, "attack". Seven different machine learning methods are applied to this dataset. In this approach, the 18 features obtained in the [Feature Selection According to Attack or Benign](#) section are used.

##### Using Features Selection for All Dataset.

An implementation that uses an alternative feature selection method can be seen in Table 9. In this method, the 18 features obtained in the [Feature Selection According to Attack or Benign](#) section are used. The changing parts are highlighted in red. When Table 8 and Table 9 are compared, there is no significant change in the algorithms of Random Forest, ID3, AdaBoost, and MLP based on F-measure. However, In Naive Bayes and QDA an increase of 2 and 11 points were observed respectively.

| Machine Learning Algorithms | Evaluation Criteria |           |        |          |          |
|-----------------------------|---------------------|-----------|--------|----------|----------|
|                             | F-Measure           | Precision | Recall | Accuracy | Time     |
| Naive Bayes                 | 0.81                | 0.8       | 0.82   | 0.82     | 1.6258   |
| QDA                         | 0.41                | 0.83      | 0.38   | 0.38     | 1.925    |
| ID3                         | 0.95                | 0.95      | 0.95   | 0.95     | 11.552   |
| AdaBoost                    | 0.94                | 0.94      | 0.94   | 0.94     | 144.166  |
| MLP                         | 0.79                | 0.815     | 0.84   | 0.84     | 51.799   |
| K Nearest Neighbours        | 0.97                | 0.97      | 0.97   | 0.97     | 1038.253 |
|                             |                     |           |        |          |          |

**Table 9.** Implementation of features obtained using Random Forest Regressor for All Dataset.

|                  | NaïveBayes |      |      | Randomforrest |       |      | KNN  |      |      | ID3   |       |      | Adaboost |      |       | MLP  |      |      | QDA   |       |      |      |       |       |       |      |      |       |       |       |      |       |       |       |       |
|------------------|------------|------|------|---------------|-------|------|------|------|------|-------|-------|------|----------|------|-------|------|------|------|-------|-------|------|------|-------|-------|-------|------|------|-------|-------|-------|------|-------|-------|-------|-------|
|                  | Acc        | F1   | Pr   | Rc            | Time  | Acc  | F1   | Pr   | Rc   | Time  | Acc   | F1   | Pr       | Rc   | Time  | Acc  | F1   | Pr   | Rc    | Time  | Acc  | F1   | Pr    | Rc    | Time  |      |      |       |       |       |      |       |       |       |       |
| Bot              | 0.56       | 0.55 | 0.82 | 0.56          | 0.003 | 0.97 | 0.97 | 0.97 | 0.97 | 0.030 | 0.96  | 0.96 | 0.95     | 0.96 | 0.014 | 0.97 | 0.97 | 0.97 | 0.008 | 0.98  | 0.98 | 0.98 | 0.170 | 0.69  | 0.62  | 0.61 | 0.69 | 0.125 | 0.68  | 0.84  | 0.68 | 0.003 |       |       |       |
| DDoS             | 0.77       | 0.76 | 0.76 | 0.77          | 0.045 | 0.96 | 0.96 | 0.97 | 0.96 | 0.430 | 0.93  | 0.93 | 0.93     | 0.93 | 1.361 | 0.96 | 0.96 | 0.97 | 0.96  | 0.200 | 0.96 | 0.96 | 0.96  | 0.28  | 0.818 | 0.77 | 0.74 | 0.76  | 0.77  | 4.962 | 0.42 | 0.35  | 0.80  | 0.42  | 0.054 |
| Dos GoldenEye    | 0.82       | 0.80 | 0.82 | 0.82          | 0.011 | 0.99 | 0.99 | 0.99 | 0.99 | 0.056 | 0.98  | 0.98 | 0.98     | 0.98 | 0.093 | 0.98 | 0.98 | 0.98 | 0.043 | 0.98  | 0.98 | 0.98 | 0.38  | 0.674 | 0.62  | 0.61 | 0.72 | 0.62  | 0.711 | 0.95  | 0.95 | 0.95  | 0.013 |       |       |
| Dos Hulk         | 0.34       | 0.23 | 0.80 | 0.34          | 0.298 | 0.94 | 0.93 | 0.94 | 0.94 | 3.738 | 0.96  | 0.96 | 0.95     | 0.96 | 2.244 | 0.96 | 0.96 | 0.96 | 0.96  | 0.009 | 0.96 | 0.96 | 0.96  | 0.26  | 22.16 | 0.94 | 0.94 | 0.94  | 0.94  | 25.63 | 0.41 | 0.36  | 0.81  | 0.41  | 0.319 |
| Dos Slowhttptest | 0.41       | 0.36 | 0.73 | 0.41          | 0.005 | 0.98 | 0.98 | 0.98 | 0.98 | 0.056 | 0.99  | 0.99 | 0.99     | 0.99 | 0.058 | 0.98 | 0.98 | 0.98 | 0.020 | 0.99  | 0.99 | 0.99 | 0.39  | 0.313 | 0.72  | 0.71 | 0.83 | 0.72  | 0.387 | 0.42  | 0.38 | 0.74  | 0.42  | 0.006 |       |
| Dos slowloris    | 0.42       | 0.36 | 0.80 | 0.42          | 0.005 | 0.95 | 0.94 | 0.94 | 0.94 | 0.055 | 0.95  | 0.95 | 0.95     | 0.95 | 0.035 | 0.95 | 0.95 | 0.95 | 0.022 | 0.95  | 0.95 | 0.95 | 0.372 | 0.77  | 0.76  | 0.80 | 0.77 | 0.494 | 0.48  | 0.46  | 0.79 | 0.48  | 0.008 |       |       |
| FTP-Patator      | 1.00       | 1.00 | 1.00 | 1.00          | 0.005 | 1.00 | 1.00 | 1.00 | 1.00 | 0.057 | -1.00 | 1.00 | 1.00     | 1.00 | 0.214 | 1.00 | 1.00 | 1.00 | 0.014 | 1.00  | 1.00 | 1.00 | 0.012 | 1.00  | 1.00  | 1.00 | 1.00 | 1.00  | 2.683 | 1.00  | 1.00 | 1.00  | 1.00  | 0.008 |       |
| Heartbleed       | 1.00       | 1.00 | 1.00 | 1.00          | 0.004 | 1.00 | 1.00 | 1.00 | 1.00 | 0.011 | -1.00 | 1.00 | 1.00     | 1.00 | 0.002 | 0.95 | 0.95 | 0.98 | 0.95  | 0.001 | 0.95 | 0.94 | 0.94  | 0.05  | 0.003 | 0.52 | 0.47 | 0.47  | 0.52  | 0.011 | 1.00 | 1.00  | 1.00  | 1.00  | 0.005 |
| Infiltration     | 0.32       | 0.79 | 0.82 | 0.82          | 0.002 | 0.93 | 0.93 | 0.95 | 0.95 | 0.011 | 0.85  | 0.86 | 0.87     | 0.85 | 0.003 | 0.91 | 0.91 | 0.94 | 0.91  | 0.002 | 0.90 | 0.90 | 0.93  | 0.90  | 0.051 | 0.59 | 0.53 | 0.53  | 0.59  | 0.008 | 0.83 | 0.82  | 0.85  | 0.83  | 0.002 |
| PortScan         | 0.44       | 0.39 | 0.80 | 0.44          | 0.185 | 1.00 | 1.00 | 1.00 | 1.00 | 2.554 | 1.00  | 1.00 | 1.00     | 1.00 | 54.72 | 1.00 | 1.00 | 1.00 | 1.00  | 0.784 | 1.00 | 1.00 | 1.00  | 1.00  | 15.01 | 0.72 | 0.61 | 0.63  | 0.72  | 13.77 | 0.84 | 0.84  | 0.89  | 0.84  | 0.205 |
| SSH-Patator      | 0.41       | 0.34 | 0.80 | 0.41          | 0.008 | 0.96 | 0.96 | 0.96 | 0.96 | 0.059 | 0.96  | 0.96 | 0.96     | 0.96 | 0.045 | 0.96 | 0.96 | 0.97 | 0.022 | 0.96  | 0.96 | 0.97 | 0.96  | 0.411 | 0.87  | 0.87 | 0.88 | 0.87  | 0.324 | 0.47  | 0.43 | 0.80  | 0.47  | 0.006 |       |
| Web Attack       | 0.73       | 0.75 | 0.86 | 0.74          | 0.005 | 0.97 | 0.97 | 0.97 | 0.97 | 0.029 | 0.93  | 0.94 | 0.94     | 0.94 | 0.014 | 0.96 | 0.96 | 0.96 | 0.009 | 0.97  | 0.96 | 0.96 | 0.009 | 0.97  | 0.170 | 0.69 | 0.64 | 0.68  | 0.69  | 0.111 | 0.83 | 0.84  | 0.89  | 0.84  | 0.003 |

Machine Learning Implementation Results (for Attack Files)

## Evaluation

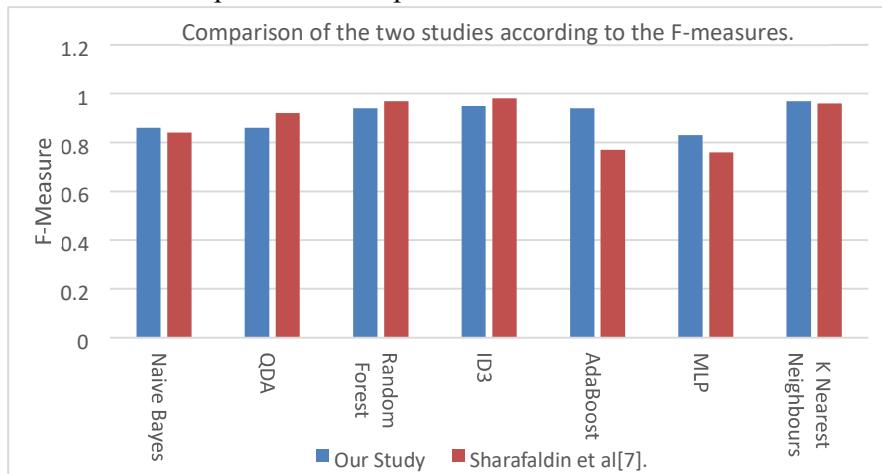
In this section, the final results of the implementation (See Table 11) are evaluated.

The F-measure was determined as the main evaluation criterion. However, in cases where equality prevailed in the results, Recall and Precision were assessed respectively. The timing of the execution was not included in the evaluation, considering that it may be misleading, even though it is shared in both works.

| Machine Learning Algorithms | Our Study   |           |        | Sharafaldin et al[4]. |           |                   |
|-----------------------------|-------------|-----------|--------|-----------------------|-----------|-------------------|
|                             | F-Measure   | Precision | Recall | F-Measure             | Precision | Recall            |
| Naive Bayes                 | <b>0.86</b> | 0.86      | 0.87   | 0.84 <sup>2</sup>     | 0.88      | 0.84 <sup>7</sup> |
| QDA                         | 0.86        | 0.87      | 0.88   | <b>0.92</b>           | 0.97      | 0.88              |
| Random Forest               | 0.94        | 0.94      | 0.94   | <b>0.97</b>           | 0.98      | 0.97              |
| ID3                         | 0.95        | 0.95      | 0.95   | <b>0.98</b>           | 0.98      | 0.98              |
| AdaBoost                    | <b>0.94</b> | 0.94      | 0.94   | 0.77                  | 0.77      | 0.84              |
| MLP                         | <b>0.83</b> | 0.82      | 0.87   | 0.76                  | 0.77      | 0.83              |
| K Nearest Neighbours        | <b>0.97</b> | 0.97      | 0.97   | 0.96                  | 0.96      | 0.96              |

**Table 12.** Comparison of the performance of two exercises against the evaluation criteria.

Table 12 and Figure 23 show the comparison of the results obtained. When the results obtained in them are examined, it is seen that Random Forest, ID3 and QDA algorithms of Sharafaldin and his colleagues [4] are 0.3,0.3,0.6 points higher than our study, respectively. This big difference in QDA is quite remarkable. Because, in our study, this algorithm has a maximum of 0.86 as a result of attempts to increase performance.



**Figure 23.** Comparison of the performance of the two studies compared to the F-measures.

On the other hand, the F-criterion score of our study is 0.1 and 0.2 points higher in the KNN and Naive Bayes algorithms, respectively. In the Adaboost and MLP methods, the difference is remarkably high (0.17 and 0.7) In both studies, MLP is the method with the lowest success rate. However, it is possible to increase performance by increasing the hidden layer size and depth in MLP. In our study, an MLP with a 3-step hidden layer and 13 neurons in each layer

was used, but in the other study, the comparison is not possible since no information is given on this topic.

The highest performance was in the KNN algorithm with 0.97 points in our Project whereas in the other study ID3 with 0.98 points. When speed is considered, in both studies Naive Bayes is the fastest algorithm, while the KNN is the slowest.

## 6 Conclusion and Future Work

### Conclusion

In this Project, it is aimed to detect network anomaly using machine learning methods. In this context, the CICIDS2017 has been used as dataset because of its up-to-datedness, wide attack diversity, and various network protocols (e.g. Mail services, SSH, FTP, HTTP, and HTTPS). This dataset contains more than 80 features that define the network flow. During the application, the importance weight calculation was made with the Random Forest Regressor algorithm to decide which of these features will be used in machine learning methods. Two approaches have been used when making these calculations. In the first place, importance weights are calculated separately for each attack type. In the second method, all the attacks are collected under a single group and the importance weights for this group are calculated. That is, the common properties that are important for all attacks are determined. Finally, seven machine learning algorithms, which are widely used and have different qualities, have been applied to this data. These algorithms and the achieved performance ratios according to F-measure are as follows (F-measure takes a value between 0 and 1): Naive Bayes: 0.86, QDA: 0.86, Random Forest: 0.94, ID3: 0.95, AdaBoost: 0.94, MLP: 0.83, and K Nearest Neighbours: 0.97.

### Future Work

In this project, various machine learning methods were applied independently from each other and experimental results were obtained. However, this method has a weak practical applicability in real life. In order to deal with this problem, a multi-layered / hierarchical machine learning structure can be designed.

## References

- [1] I. Sharafaldin, A. Gharib, A. H. Lashkari, and A. A. Ghorbani, "Towards a reliable intrusion detection benchmark dataset," *Software Networking*, vol. 1, no. 1, pp. 177200, 2017.
- [2] M. Ahmed, A. N. Mahmood, and J. Hu, "A survey of network anomaly detection techniques," *Journal of Network and Computer Applications*, vol. 60, pp. 19-31, 2016.
- [3] "NSL-KDD dataset," *Canadian Institute for Cybersecurity, University of New Brunswick*, [Online]. Available: <http://www.unb.ca/cic/datasets/nsl.html>. [Accessed 06 Aug 2018].
- [4] "Intrusion Detection Evaluation Dataset (CICIDS2017)," *Canadian Institute for Cybersecurity, University of New Brunswick*, [Online]. Available: <http://www.unb.ca/cic/datasets/ids-2017.html>. [Accessed 08 Aug 2018].

- [5] *J. P. Mueller and L. Massaron, "Machine Learning For Dummies Cheat Sheet," dummies. [Online]. Available: <https://www.dummies.com/programming/bigdata/data-science/machine-learning-dummies-cheat-sheet/>. [Accessed: 14 Aug 2018].*
- [6] "sklearn.preprocessing.LabelEncoder," *1.4. Support Vector Machines - scikit-learn 0.19.1 documentation*. [Online]. Available: <http://scikitlearn.org/stable/modules/generated/sklearn.preprocessing.LabelEncoder.html>. [Accessed: 19-Aug-2018].
- [7] "train\_test\_split," *scikit-learn 0.19.1 documentation*. [Online]. Available: [http://scikitlearn.org/stable/modules/generated/sklearn.model\\_selection.train\\_test\\_split.html](http://scikitlearn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html). [Accessed: 20-Aug-2018].
- [8] *J. Brownlee, "Feature Importance and Feature Selection With XGBoost in Python," Machine Learning Mastery, 10-Mar-2018. [Online]. Available: <https://machinelearningmastery.com/feature-importance-and-feature-selection-withxgboost-in-python/>. [Accessed: 19-Aug-2018].*
- [9] R. Alshammari and A. N. Zincir-Heywood, "A flow based approach for SSH traffic detection," in *Systems, Man and Cybernetics, 2007. ISIC. IEEE International Conference on*, 2007, pp. 296-301: IEEE.
- [10] H. Choi, H. Lee, and H. Kim, "Fast detection and visualization of network attacks on parallel coordinates," *computers & security*, vol. 28, no. 5, pp. 276-288, 2009.
- [11] "sklearn.ensemble.RandomForestRegressor," *scikit-learn*. [Online]. Available: <http://scikitlearn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>. [Accessed: 19-Aug-2018].