

Cyber Security Project: Network Attack prediction

MINI PROJECT REPORT

submitted to the faculty of the

INSTITUTE OF ENGINEERING AND MANAGEMENT, KOLKATA

in partial fulfillment of the requirements for the degree of Bachelor of Science

by – 66_Ankur Kumar Lal Das(12019002002096)

67_Ashutosh Kashyap (12019002002098)

23_Roushan Kumar Singh (12019002002032)

54_Rishav Raj (12019002002076)

Date: 16th June 2022

Keywords:

1. Computer Science Engineering
2. Cyber Security
3. Network Attack Prediction using different classifiers
4. Binary Label Classification
5. Supervised Learning

Approved By -

ABSTRACT

Every day 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 has the ability to 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. 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 Neighbours 97%.

PROBLEM STATEMENT

Select the CICIDS 2017 dataset and perform a binary classification of predicting attack types using the feature set. Also use different AI/ML models for classification and compare the performance of each classifier.

GOAL AND OBJECTIVE

The goals that are aimed to achieve at the end of this study 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.

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.

ANALYSIS

1. DATA EXPLORATION :

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 need, many datasets have been produced and continue to be produced. But as per the problem given, the suggested dataset was 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:

Flow Recording Day (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

NETWORK ATTACK PREDICTION

2. NETWORK ATTACK TYPES:

Network security tries to protect the network from attacks against these three principles: confidentiality, integrity, and availability.

Confidentiality: Information should be accessible only to legitimate users and unauthorized access should be prevented..

Integrity: adding, modifying and deleting information can only be done by the legitimate user. Unauthorized persons should not be able to modify information.

Accessibility: The system should always be accessible to the legitimate user. Network attacks are attempts to violate these 3 essential features. The attacks can be summarized in 4 headings.

Denial of Service (DoS): In this type of attack, an attacker abuses system resources, preventing the legitimate user from taking advantage of the service. The simplest example is to make it drop out of service by sending a large number of requests to a web server. Dos attacks can be divided into two subparts as bandwidth depletion and resource depletion. While bandwidth depletion aims to consume a victim's bandwidth by providing a very high data flow, resource depletion attacks usually try to consume resources such as the victim's memory and processor with a lot of packages.

Probe (Information Gathering): These attacks aim at collecting information about the target. Through this attack, attackers can get a lot of important information such as network structure, the operating system used, types and properties of networked devices. Although this attack does not directly affect the system, it is very important because it is preparing the ground for many attacks that can harm the system .

NETWORK ATTACK PREDICTION

U2R (User to Root): In this type of attack, an attacker tries to gain control of the administration account in order to access and steal important resources. the attacker may use system vulnerabilities or brute-force attacks to gain the Administrator account .

R2U / R2L (Remote to User / Remote to Local): In this attack, the attacker infiltrates the victim's network to gain a privilege to send packets from the victim computer. An attacker can use system vulnerabilities or brute-force attacks to gain this privilege.

Classifying the network attacks according to the anomalies they create, can be useful in detecting the attack. Each attack causes differentiation (an anomaly) in the network flow. For example, DoS attacks are known to increase the amount of data flowing and the number of packets in the network. During a DoS attack, a large number of abnormal streams can be mentioned. Thus, classify DoS attacks as collective anomalies. On the other hand, it would be more appropriate to classify U2R and R2U as contextual and point anomaly because the attack is for a particular user, certain port, and specific purpose. Additionally, the probe attack creates a dense packet traffic on the network, as well as it has a specific purpose. In this context, the probe attack can be defined as a contextual and collective anomaly.

ATTACKS

In this section, the types of attacks that the data set contains are examined in detail. All data in the dataset are 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. The names and numbers of these labels can be seen in Figure 3.

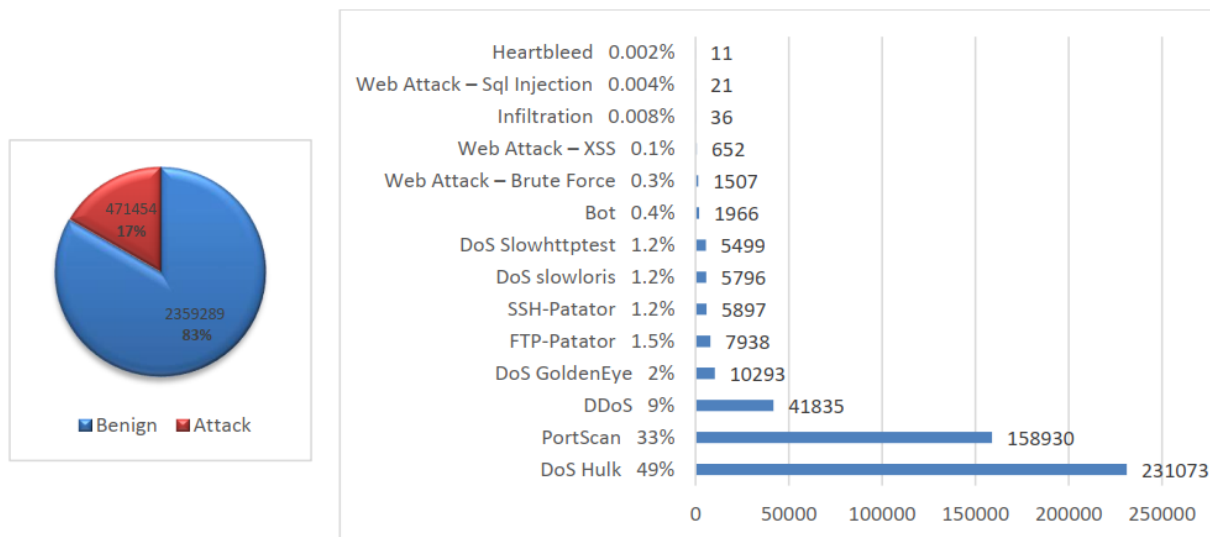


Figure 3. The distribution of data flow and attack types in the dataset.

When the numbers of these attacks are examined, it is immediately apparent that the numbers of some attacks are very high. For example, the DoS HULK attack is almost half of all attacks and the PortScan attack is one-third of all attacks. The reason for this imbalance in the distribution is the nature of these attacks. Both DoS and PortScan attacks cause too much data and packet flows during the attack. Therefore, during these attacks, it is completely natural to observe more intense traffic from normal usage and other types of attacks.

NETWORK ATTACK PREDICTION

The attacks in the dataset are given as follows:

- **DoS HULK**
- **TCP-SYN Flood**
- **HTTP-GET Flood**
- **PortScan**
- **SYN Scan**
- **TCP Connect Scan**
- **ACK Scan**
- **FIN Scan**
- **NULL Scan**
- **XMAS Scan**
- **UDP Scan**
- **Fragmentation Attack**
- **DDoS**
- **DoS Goldeneye**
- **FTP-Patator**
- **SSH-Patator**
- **DoS Slowloris**
- **DoS SlowHTTPTest**
- **Botnet**
- **Web Attack**
- **Infiltration**
- **Heartbleed**

MACHINE LEARNING

Machine learning is a science and art that enables the programmed computers to learn from the data given to them [40]. In the machine learning process, computers can be trained on the data (training set) given to them, and they can show their performance on a different data (test set). In this way, the problem is solved by minimizing human intervention. Machine learning is heavily used in many places where classical methods are inefficient. Areas of use can be listed as follows:

- It can interpret very large and complicated data.
- It can solve complex problems that traditional methods cannot find solutions.
- Machine Learning methods can find solutions to situations where existing solutions require too much external intervention/update without external intervention.
- It can work in variable environments. Machine learning methods can be applied to a new situation by examining the data. Machine learning algorithms are divided into 4 groups according to whether the training data are labelled or not and according to the training supervision they have received. These are supervised learning, unsupervised learning, semi-supervised learning and reinforcement learning.

Supervised learning: In this method, the training data had been correctly classified and labelled. For example, all flows in the dataset contain information (tags/labels) about their nature (such as normal or harmful). At the next stage, the test/prediction phase, these tags are compared to the results found by the algorithm, and the algorithm's success is calculated. The performance of the method is high. However, supervised learning is costly because it used external service (e.g. manual tagging) for labelling. Examples of such algorithms are Decision

NETWORK ATTACK PREDICTION

Trees, K-Nearest Neighbours, and Random Forests.

In this project, supervised learning methods will be used to take advantage of having a manually-labelled good dataset. In this respect, it is aimed to achieve high-performance advantage without getting cost disadvantage. Used machine learning methods used in the application phase are: Naive Bayes, QDA, Random Forest, ID3, AdaBoost, MLP, and K Nearest Neighbours. While choosing these methods, the focus is on bringing together popular algorithms with different characteristics. In this context, the algorithms used as follow.

- Naïve Bayes
- Decision Trees
- Random Forest
- K Nearest Neighbour
- AdaBoost
- MLP (Multi-Layer Perceptron)

METHDOLOGY

1. SOFTWARE PLATFORM:

Python, a free and open source object-oriented programming language, draws attention with its simple syntax and dynamic structure. It works in concert with many libraries which "machine learning" applications can be done.

Sklearn(Scikit-learn) is a machine learning library that can be used with the Python programming language. Sklearn offers a wide range of options to the user with its numerous machine learning algorithms.

NETWORK ATTACK PREDICTION

Pandas is a powerful data analysis library running on Python. When working with a large dataset, Pandas allows you to easily perform many operations such as filtering, bulk column / row deletion, addition, and replacement. Because of all these advantages, the Pandas library has been used.

Matplotlib is a library that runs on Python, allowing visualization of data. This library is used to create graphs used in the study.

NumPy, a Python library that allows you to perform mathematical and logical operations quickly and easily, has been used in calculations in this work.

Jupyter Notebook

2. PERFORMANCE EVALUATION METRIC

The results of this study are evaluated according to four criteria, namely accuracy, precision, f-measure, and recall. All these criteria take a value between 0 and 1. When it approaches 1, the performance increases, while when it approaches 0, it decreases.

Accuracy: The ratio of successfully categorized data to total data.

Recall (Sensitivity): The ratio of data classified as an attack to all attack data.

Precision: The ratio of successful classified data as the attack to all data classified as the attack.

F-measure (F-score/F1-score): The harmonic-mean of sensitivity and precision. This concept is used to express the overall success. so, in this

NETWORK ATTACK PREDICTION

study, when analysing the results, it will be focused, especially on the F1 Score.

In calculating these four items, the four values summarized below are used:

- **TP:** True Positive (Correct Detection) The attack data classified as attack
- **FP:** False Positive (Type-1 Error) The benign data classified as attack.
- **FN:** False Negative (Type-2 Error) The attack data classified as benign.
- **TN:** True Negative (Correct Rejection) The benign data classified as benign. This distribution is presented by visualizing Confusion matrix in the fig, also.

True Class	
Predicted Class	Attack
	Benign
Attack	TP True Positive (Desired)
Benign	FN False Negative (Undesired)

3. IMPLEMENTATION

In this section, various pre-processing and actual application 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.

NETWORK ATTACK PREDICTION

➤ DATA CLEANING

The dataset file contains 3119345 stream records. The distribution of these stream records can be seen from Table 2. When these records are examined, it can be seen that the 288602 record is incorrect / incomplete³. The first step in the pre-processing process will be to delete these unnecessary records.

Label Name	Number
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

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

NETWORK ATTACK PREDICTION

(column 62). Another change that needs to be made in the dataset is to convert the properties including the categorical and string values (Flow ID, Source IP, Destination IP, Timestamp, External IP) into numerical data to be used in machine learning algorithms. This can be done with `LabelEncoder()` from Sklearn classes. In this way, various string values that cannot be used in machine learning operations will get integer values between 0 and n-1 and will become more suitable for processing.

➤ CREATING TRAINING AND TESTING DATASET

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` 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 and this ratio is also preferred in this application.

➤ FEATURE ENGINEERING

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

- **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).

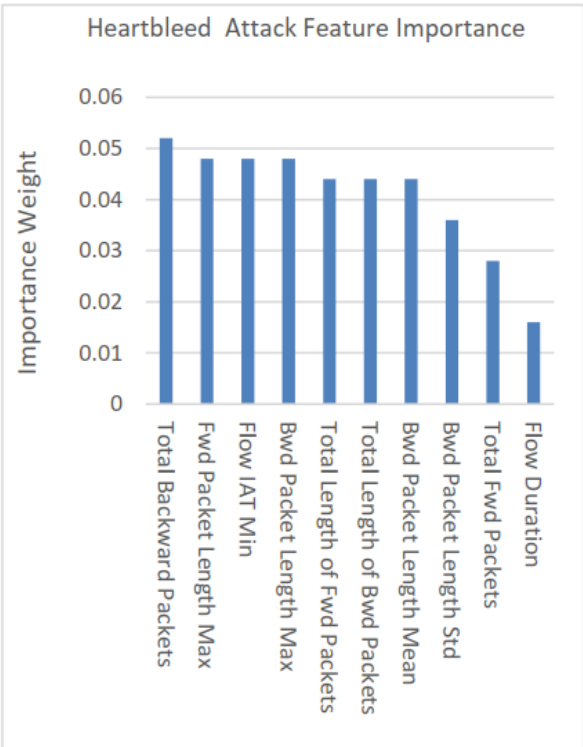
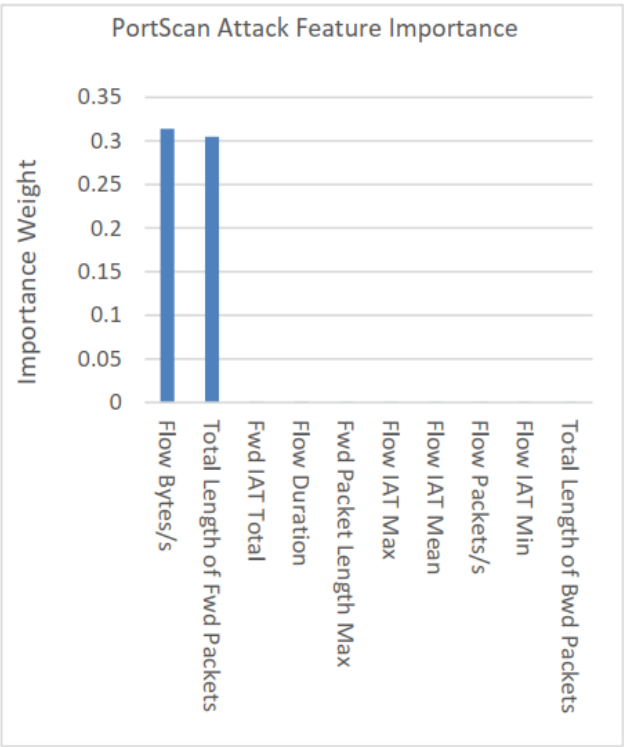
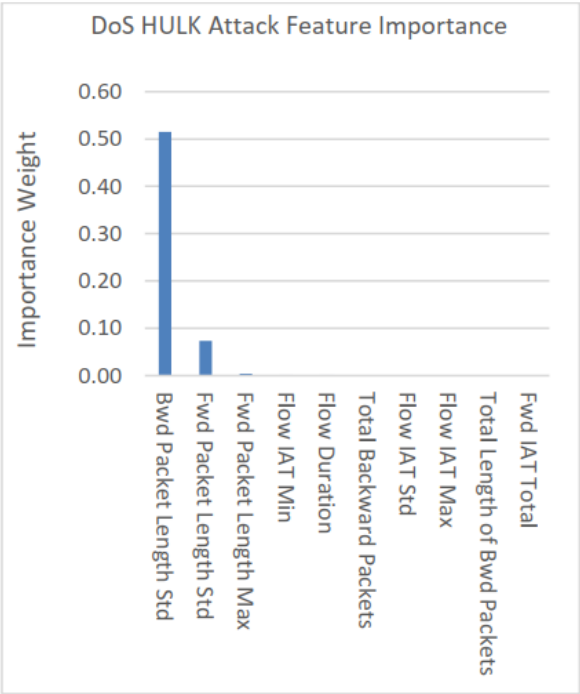
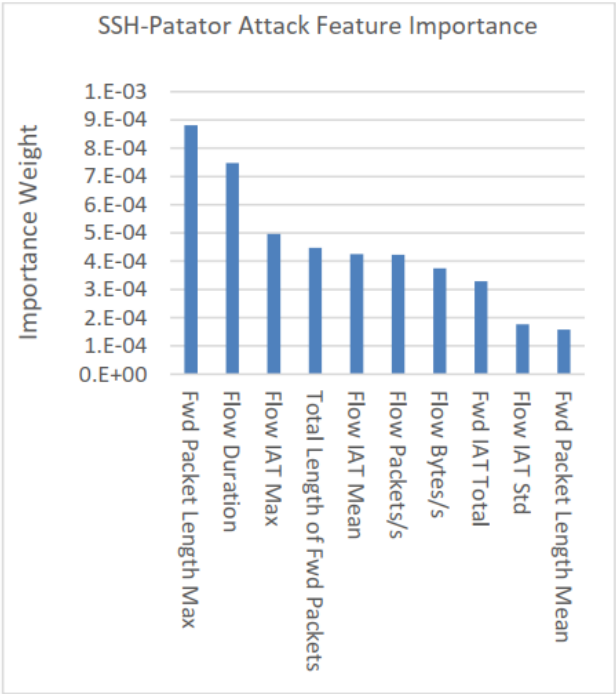
NETWORK ATTACK PREDICTION

While choosing the attribute importance of the attack, it will be much more effective to eliminate the misleading features such as IP address, Port number, Timestamp, use more generic and invariant attributes to define the attack. Because the shape of the data will give much more information about whether or not it is an attack.

Attack / Feature Name	Importance Weight	Attack / Feature Name	Importance Weight
Bot		FTP-Patator	
Bwd Packet Length Mean	0.304823	Fwd Packet Length Max	0.063671
Flow IAT Max	0.034495	Fwd Packet Length Std	0.022751
Flow IAT Std	0.019464	Fwd Packet Length Mean	0.002179
Flow Duration	0.010129	Total Length of Bwd Packets	0.000746
DDoS		Heartbleed	
Bwd Packet Length Std	0.468089	Bwd Packet Length Mean	0.064
Total Backward Packets	0.094926	Total Length of Bwd Packets	0.056
Fwd IAT Total	0.012066	Flow IAT Min	0.056
Total Length of Fwd Packets	0.006438	Bwd Packet Length Std	0.044
DoS GoldenEye		Infiltration	
Flow IAT Max	0.442727	Total Length of Fwd Packets	0.05238
Bwd Packet Length Std	0.091185	Flow IAT Max	0.036096
Flow IAT Min	0.053795	Flow Duration	0.016453
Total Backward Packets	0.041583	Flow IAT Min	0.015448
DoS Hulk		PortScan	
Bwd Packet Length Std	0.514306	Flow Bytes/s	0.313402
Fwd Packet Length Std	0.069838	Total Length of Fwd Packets	0.304917
Fwd Packet Length Max	0.008542	Flow Duration	0.000485
Flow IAT Min	0.001716	Fwd Packet Length Max	0.00013

DoS Slowhttptest		SSH-Patator	
Flow IAT Mean	0.64206	Flow Bytes/s	0.000846
Fwd Packet Length Min	0.075942	Total Length of Fwd Packets	0.000814
Fwd Packet Length Std	0.022194	Fwd Packet Length Max	0.000749
Bwd Packet Length Mean	0.020857	Flow IAT Mean	0.000734
DoS slowloris		Web Attack	
Flow IAT Mean	0.465561	Total Length of Fwd Packets	0.014697
Bwd Packet Length Mean	0.075633	Bwd Packet Length Std	0.00536
Total Length of Bwd Packets	0.049808	Flow Bytes/s	0.00257
Total Fwd Packets	0.01868	Bwd Packet Length Max	0.001922

NETWORK ATTACK PREDICTION



NETWORK ATTACK PREDICTION

- **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:

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

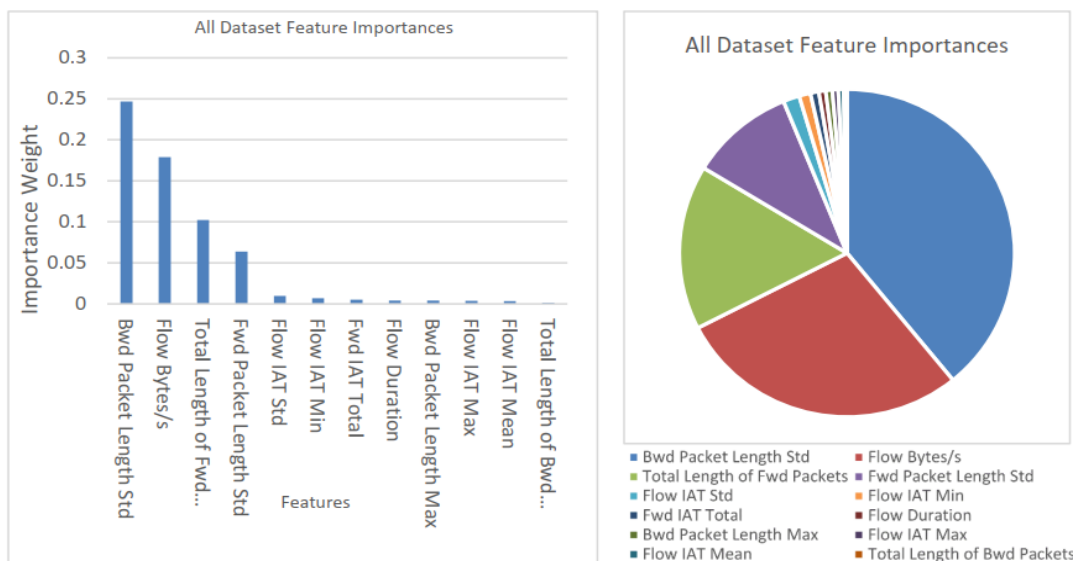


Figure 16. Graphs of Feature Importance Weights According to Attack and Benign Labels

NETWORK ATTACK PREDICTION

➤ MODELLING

Two different approaches have been used to apply machine learning algorithms⁴ 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 below:

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.

NETWORK ATTACK PREDICTION

Feature Name	Importance Weight	Percentage
Bwd Packet Length Std	0.246627	38.97%
Flow Bytes/s	0.178777	28.25%
Total Length of Fwd Packets	0.102417	16.18%
Fwd Packet Length Std	0.063889	10.10%
Flow IAT Std	0.009898	1.56%
Flow IAT Min	0.006946	1.10%
Fwd IAT Total	0.005121	0.8 %

The importance weights obtained in "Feature Selection According to Attack or Benign"

RESULTS

1. MODEL EVALUATION :

In this section, the results of the studies 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.

- **Using 12 Attack Types**

Seven different machine learning methods are applied to 12 different attack types and the obtained results are presented:

NETWORK ATTACK PREDICTION

Attack Names	F-Measures						
	NB	RF	KNN	ID3	AB	MLP	QDA
Bot	<u>0.54</u>	0.96	0.95	0.96	0.97	0.64	0.68
DDoS	0.77	0.96	0.92	0.96	0.96	0.76	<u>0.34</u>
DoS GoldenEye	0.81	0.99	0.98	0.99	0.99	<u>0.64</u>	0.71
DoS Hulk	<u>0.23</u>	0.93	0.96	0.96	0.96	0.95	0.36
DoS Slowhttptest	<u>0.35</u>	0.98	0.99	0.98	0.99	0.78	0.38
DoS slowloris	<u>0.37</u>	0.95	0.95	0.96	0.95	0.74	0.46
FTP-Patator	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Heartbleed	1.00	0.99	1.00	0.95	0.93	<u>0.66</u>	1.00
Infiltration	0.78	0.92	0.88	0.89	0.92	<u>0.52</u>	0.83
PortScan	<u>0.39</u>	1.00	1.00	1.00	1.00	0.61	0.85
SSH-Patator	<u>0.33</u>	0.96	0.95	0.96	0.96	0.83	0.41
Web Attack	0.74	0.97	0.93	0.97	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.

When these box and whisker graphics are examined, it appears that there is a balanced distribution for almost all attack types. The difference between the highest and lowest Fmeasure levels of all attack types is between 0%- 2%. However, this distribution is deteriorated in

NETWORK ATTACK PREDICTION

two types of attack (Infiltration and Heartbleed). The difference between the lowest and highest values in Heartbleed is 10%, while the difference between the highest and lowest values in Infiltration is 8%.

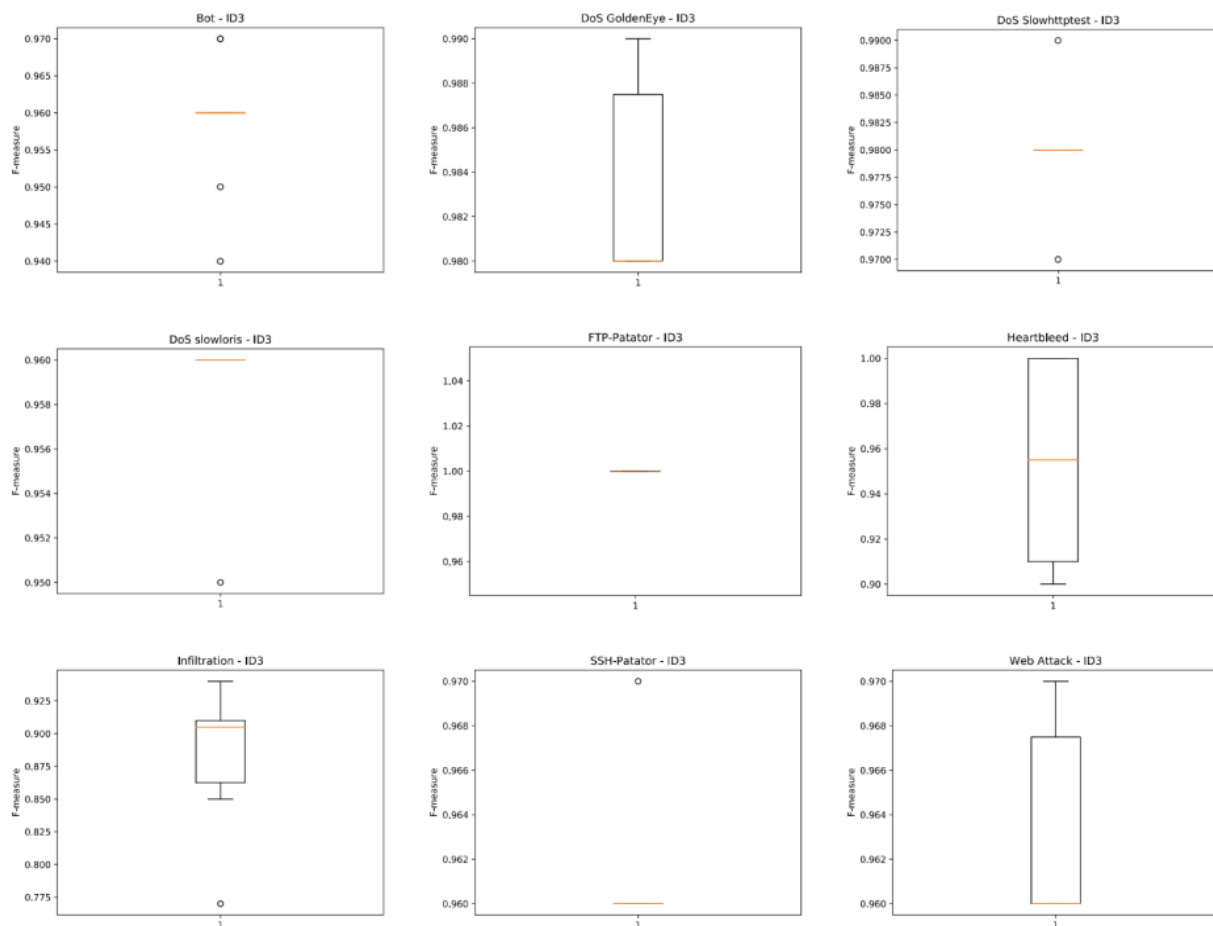


Figure 19. Box and whisker graphics containing the results of applying the ID3 algorithm to various attack types.

- **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, two methods will be used, the first one, the features created for attack files in approach 1 are used.

NETWORK ATTACK PREDICTION

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

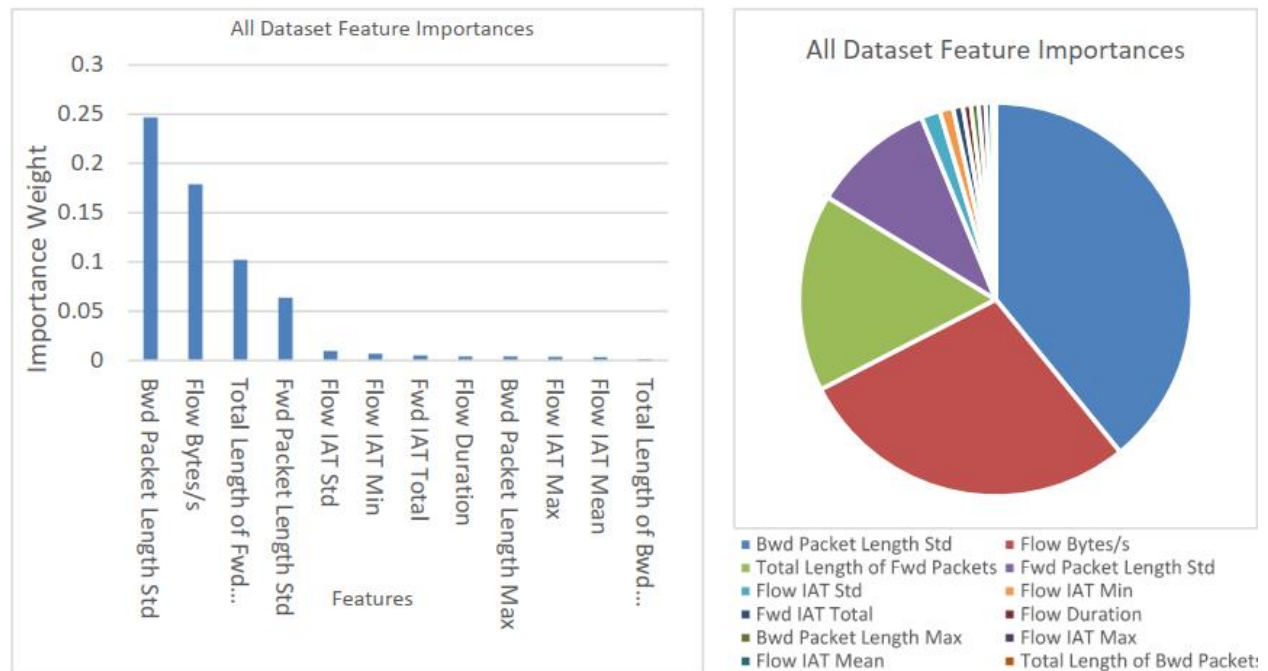


Figure 16. Graphs of Feature Importance Weights According to Attack and Benign Labels

2. EVALUATION AND DISCUSSION

- Using 12 Attack Types

Attack Names	F-Measures						
	NB	RF	KNN	ID3	AB	MLP	QDA
Bot	<u>0.54</u>	0.96	0.95	0.96	0.97	0.64	0.68
DDoS	0.77	0.96	0.92	0.96	0.96	0.76	<u>0.34</u>
DoS GoldenEye	0.81	0.99	0.98	0.99	0.99	<u>0.64</u>	0.71
DoS Hulk	<u>0.23</u>	0.93	0.96	0.96	0.96	0.95	0.36
DoS Slowhttptest	<u>0.35</u>	0.98	0.99	0.98	0.99	0.78	0.38
DoS slowloris	<u>0.37</u>	0.95	0.95	0.96	0.95	0.74	0.46
FTP-Patator	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Heartbleed	1.00	0.99	1.00	0.95	0.93	<u>0.66</u>	1.00
Infiltration	0.78	0.92	0.88	0.89	0.92	<u>0.52</u>	0.83
PortScan	<u>0.39</u>	1.00	1.00	1.00	1.00	0.61	0.85
SSH-Patator	<u>0.33</u>	0.96	0.95	0.96	0.96	0.83	0.41
Web Attack	0.74	0.97	0.93	0.97	0.97	<u>0.60</u>	0.84

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

- Using 2 Groups: Attack and Benign

Machine Learning Algorithms	Evaluation Criteria				
	F-Measure	Precision	Recall	Accuracy	Time ⁵
Naive Bayes	0.79	0.80	0.78	0.78	4.576
QDA	<u>0.30</u>	0.84	0.31	0.31	6.649
Random Forest	0.94	0.95	0.94	0.94	24.739
ID3	0.95	0.95	0.95	0.95	29.284
AdaBoost	0.95	0.95	0.95	0.95	391.804
MLP	0.79	0.81	0.84	0.84	81.668
K Nearest Neighbours	0.96	0.96	0.97	0.97	<u>1967.054</u>

Table 8. Application of the features obtained in the first approach.

NETWORK ATTACK PREDICTION

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
Random Forest	0.94	0.947	0.94	0.94	20.511
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.

Machine Learning Algorithms	Evaluation Criteria				
	F-Measure	Precision	Recall	Accuracy	Time
Naive Bayes	0.86	0.86	0.87	0.87	1.8255
QDA	0.86	0.87	0.88	0.88	2.3696
Random Forest	0.94	0.94	0.94	0.94	19.0899
ID3	0.95	0.95	0.95	0.95	9.5107
AdaBoost	0.94	0.94	0.94	0.94	135.2455
MLP	0.83	0.82	0.87	0.87	59.6933
K Nearest Neighbours	0.97	0.97	0.97	0.97	1626.833

Table 11. The Final Results - Implementation of features using Table 10.

NETWORK ATTACK PREDICTION

	NaiveBayes					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	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	0.96	0.96	0.014	0.97	0.97	0.97	0.97	0.008	0.98	0.98	0.98	0.98	0.170	0.69	0.62	0.61	0.69	0.125	0.68	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	0.96	2.818	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	0.096	0.98	0.98	0.98	0.98	0.093	0.98	0.98	0.98	0.98	0.043	0.98	0.98	0.98	0.98	0.674	0.62	0.61	0.72	0.62	0.711	0.95	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	0.96	0.96	254.4	0.96	0.96	0.96	0.96	0.909	0.96	0.96	0.96	0.96	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	0.006	0.98	0.98	0.98	0.98	0.056	0.99	0.99	0.99	0.99	0.058	0.98	0.98	0.99	0.98	0.020	0.99	0.99	0.99	0.99	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	0.006	0.95	0.94	0.94	0.94	0.055	0.95	0.95	0.95	0.95	0.035	0.96	0.96	0.96	0.96	0.022	0.95	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	0.006	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	1.00	0.014	1.00	1.00	1.00	1.00	0.412	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	0.95	0.003	0.52	0.47	0.47	0.52	0.011	1.00	1.00	1.00	1.00	0.005
Infiltration	0.82	0.79	0.82	0.82	0.002	0.93	0.93	0.95	0.93	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	0.069	0.96	0.96	0.96	0.96	0.045	0.96	0.96	0.97	0.96	0.027	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	0.96	0.009	0.97	0.96	0.96	0.97	0.170	0.69	0.64	0.68	0.69	0.111	0.83	0.84	0.89	0.84	0.003

CONCLUSION

In this study, it is aimed to detect network anomaly using machine learning methods. In this context, the CICIDS2017[16] 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)[4]. 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[63] 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

NETWORK ATTACK PREDICTION

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.

ABOUT US -

ANKUR KUMAR LAL DAS (12019002002096)

ASHUTOSH KASHYAP(12019002002098)

RAUSHAN KUMAR SINGH (12019002002032)

RISHAV RAJ (12019002002076)

Prof.

Mentor -

References

1. "Intro to Machine Learning | Udacity." Intro to Machine Learning | Udacity. <https://www.udacity.com/course/intro-to-machine-learning--ud120>.
2. K. Kostas, "Anomaly Detection in Networks Using Machine Learning," Research Proposal
3. "Internet Growth Statistics," *Miniwatts Marketing Group*, 2 Mar 2018. [Online]. Available: <https://www.internetworldstats.com/emarketing.htm>.
4. "Dataset" <https://www.unb.ca/cic/datasets/ids-2017.html>
5. "KDD Cup 1999 Data," *University of California, Irvine*, [Online]. Available: <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>.
6. 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.
7. "Intrusion detection evaluation dataset (ISCXIDS2012)," *Canadian Institute for Cybersecurity, University of New Brunswick*, [Online]. Available: <http://www.unb.ca/cic/datasets/ids.html>.
8. "A new DOS Perl Program," *GitHub*, 05 Nov 2013. [Online]. Available: <https://github.com/llaera/slowloris.pl>.
9. "slowhttpptest," *GitHub*. [Online]. Available: <https://github.com/shekyaan/slowhttpptest/wiki>.
10. "Ares," *GitHub*, 08-Dec-2017. [Online]. Available: <https://github.com/sweetsoftware/Ares>.
11. "antoinedelplace" Github . Available: <https://github.com/antoinedelplace/Cyberattack-Detection>