

Machine Learning-based Prevalent Model for Detecting DRDoS Attacks through Features Optimization Technique

Submitted as a research project in partial fulfillment of the
requirements for the Degree of Bachelor of Science in
Computer Science and Engineering

Submitted By

Pabon Shaha

ID: CE17009

Session: 2016-17

Supervised By

Dr. Mostofa Kamal Nasir

Professor,

Department of Computer Science & Engineering,
Mawlana Bhashani Science and Technology University



Department of Computer Science and Engineering (CSE)

Mawlana Bhashani Science and Technology University

Santosh, Tangail-1902,

Bangladesh

APPROVAL

The Research work entitled “Machine Learning based Prevalent Model for Detecting DRDoS Attack Through Feature Optimization Technique.” It has been accepted as satisfactory for the partial fulfillment of the requirements for the Degree of Bachelor of Science and Engineering in Computer Science and Engineering and approved as to its styles and contents. The submission was made by Pabon Shaha, ID: CE-17009, to the Department of Computer Science and Engineering, Mawlana Bhashani Science and Technology University, Santosh, Tangail-1902, Bangladesh.

.....
Dr. Mostofa Kamal Nasir

(Supervisor)

Professor,

Department of Computer Science & Engineering,

Mawlana Bhashani Science and Technology University

.....
(External Examiner)

DECLARATION

I, hereby, declare that this research work presented by me is the outcomes of the investigations under the supervision of Dr. Mostofa Kamal Nasir, Professor, Department of Computer Science and Engineering, Mawlana Bhashani Science and Technology University, Santosh, Tangail-1902, Bangladesh.

I want to reaffirm that there hasn't been any big financial backing for this effort that would have affected how it turned out. I further certify that no portion of this thesis or its components have been or are being considered for submission to another institution for the purpose of receiving a degree or diploma.

Pabon Shaha
CE-17009

ACKNOWLEDGEMENTS

I want to start by giving thanks to God. He gave me the capability, opportunity, and a helpful supervisor, which allowed me to successfully complete my task in this instance today.

I would like to convey my profound gratitude to my esteemed supervisor, professor Dr. Mostofa Kamal Nasir, for his invaluable advice, insight, encouragement, support, and reliance during the project. Even though he was constantly busy with other things, he left me plenty of time to complete this assignment. I much appreciated his advice and remarks as I was putting together my research paper.

I am also grateful to the rest of CSE Teachers, MBSTU who have helped us by giving valuable suggestions in different time.

I want to specially thank to my friends for encouraging me and proving me useful information during my work. I owe my loving thanks to my parents for their loving support.

Finally, my special thanks go to authors whose works have helped us during in this work.

ABSTRACT

Machine Learning based Prevalent Model for Detecting DRDoS Attack Through Feature Optimization Technique.

Pabon Shaha, Professor Dr. Mostofa Kamal Nasir (supervisor)

Internet users are becoming increasingly worried about different kinds of cyber security attacks in the modern eras. DRDoS attack is one of the spirited cyber-attack and it is also a type of DDoS attack. We need to detect this kind of attack quickly and take immediate steps to protect internet users from this kind of attack. In the previous year, several works are done on DRDoS attack detection and prevention. But the detection rate of DRDoS attacks is inadequate. So in this research, to increase the detection rate of DRDoS attacks four machine learning algorithms like Support Vector Machine(SVM), Decision Tree(DT), Random Forest (RT), Logistic Regression(LR), and also Principle Component Analysis (PCA) technique are used for detecting DRDoS attack. “CIC Bell DNS 2021” dataset from Canadian Institute for Cybersecurity (CIC) have used for this research work. The experimental results for the detection of DRDoS attacks appear that the DT and RF algorithms accomplish the highest accuracy of 100% respectively with a 1.00 F1 score in both without Feature Reduction and with Feature Reduction using PCA (except dataset 4) approach.

CONTENTS

| | Page |
|---|-------------|
| APPROVAL | 02 |
| DECLARATION | 03 |
| ACKNOWLEDGEMENTS | 04 |
| ABSTRACT | 05 |
| CONTENTS | 06 |
| LIST OF FIGURES | 09 |
| LIST OF TABLES | 10 |
| ABBREVIATIONS | 11 |
| | |
| Chapter 1: Introduction | 12 |
| 1.1 Distributed Reflection Denial of Service (DRDoS) attacks | 12 |
| 1.1.1 Domain Name Service (DNS) | 13 |
| 1.1.2 Network Basic Input / Output System(NetBIOS) | |
| 1.1.3 Network Time Protocol (NTP) | |
| 1.1.4 User Datagram Protocol(UDP) | |
| | |
| 1.2 Machine Learning | 15 |
| 1.4 Research questions and objectives | 15 |
| 1.5 Research synopsis | 16 |
| 1.6 Fundamental contributions | 16 |
| 1.7 Outline of the research | 17 |

Chapter 2: Literature Review

| | |
|---|----|
| 2.1 Background | 18 |
| 2.2 Exploration of Deep learning and Machine learning | 18 |
| 2.3 Summery table of all previous research work | 19 |
| 2.4 Research gap | 20 |

Chapter 3: Research Method

| | |
|---|----|
| 3.1 Dataset | 22 |
| 3.2 Proposed method | 23 |
| 3.2.1 Method-1 (Without Feature Reduction) | 23 |
| 3.2.2 Method-2 (With Feature Reduction using PCA) | 23 |
| 3.2.3 Feature optimization | 24 |
| 3.2.3.1. Principal Component Analysis(PCA) | |
| 3.3 Machine Learning Algorithm | 25 |
| 3.3.1 Decision Tree Algorithm(DT) | |
| 3.3.2 Support Vector Machine(SVM): | |
| 3.3.3 Random Forest(RF) | |
| 3.3.4 Logistic Regression(LR) | |
| 3.3.5 Performance parameters | |

Chapter 4: Result and Discussion

| | |
|--|----|
| | 29 |
| 4.1 Analysis the result for Dataset-1 | 30 |
| 4.2 Analysis the result for Dataset-2 | 32 |
| 4.3 Analysis the result for Dataset-3 | 34 |
| 4.4 Analysis the result for Dataset-4 | 36 |
| 4.5 Summary of all models performance for all Datasets | 38 |
| 4.5.1. Summary Table for Dataset-1 | |
| 4.5.2. Summary Table for Dataset-2 | |
| 4.5.3. Summary Table for Dataset-3 | |
| 4.5.4. Summary Table for Dataset-4 | |

| | |
|---|-----------|
| 4.6 ROC Curve | 49 |
| 4.7 Compare proposed model with existing models | 51 |
| Chapter 5: Conclusion and Future Scope | 53 |
| 5.1 Conclusion | 53 |
| 5.2 Limitation | 53 |
| 5.3 Future scope | 53 |
| References | 54 |

LIST OF FIGURES

| Serial | Figure Name | Page |
|-------------|--|------|
| Figure 1.1 | The DRDoS Attacks | 13 |
| Figure 3.1 | Proposed Model | 22 |
| Figure 3.2 | Covariance Curve | 22 |
| Figure 3.3 | Support Vector Machine(SVM) | 25 |
| Figure 3.4 | Random Forest (RF) | 24 |
| Figure 4.1 | Performance of all models for dataset 1 | 39 |
| Figure 4.2 | Performance of all models for dataset 2 | 39 |
| Figure 4.3 | Performance of all models for dataset 3 | 40 |
| Figure 4.4 | Performance of all models for dataset 4 | 40 |
| Figure 4.5 | ROC Curve of all models for Dataset-1(without Feature Reduction) | 42 |
| Figure 4.6 | ROC Curve of all models for Dataset-1(with Feature Reduction) | 43 |
| Figure 4.7 | ROC Curve of all models for Dataset-2(without Feature Reduction) | 44 |
| Figure 4.8 | ROC Curve of all models for Dataset-2(with Feature Reduction) | 45 |
| Figure 4.9 | ROC Curve of all models for Dataset-3(without Feature Reduction) | 46 |
| Figure 4.10 | ROC Curve of all models for Dataset-3(with Feature Reduction) | 47 |
| Figure 4.11 | ROC Curve of all models for Dataset-4(without Feature Reduction) | 48 |
| Figure 4.12 | ROC Curve of all models for Dataset-4(with Feature Reduction) | 49 |

LIST OF TABLE

| Serial | Table Name | Page |
|------------|---|------|
| Table 2.3 | Summery table of all previous research work | 20 |
| Table 3.1 | Dataset | 22 |
| Table 3.2 | Confusion Matrix table | 22 |
| Table 3.3 | Performance Measurement Parameters | 29 |
| Table 4.1 | Overall performance of ML model without Feature Reduction for Dataset-1 | 30 |
| Table 4.2 | Overall performance of ML model with Feature Reduction for Dataset-1 | 31 |
| Table 4.3 | Overall performance of ML model without Feature Reduction for Dataset-2 | 32 |
| Table 4.4 | Overall performance of ML model with Feature Reduction for Dataset-2 | 33 |
| Table 4.5 | Overall performance of ML model without Feature Reduction for Dataset-3 | 34 |
| Table 4.6 | Overall performance of ML model with Feature Reduction for Dataset-3 | 35 |
| Table 4.7 | Overall performance of ML model without Feature Reduction for Dataset-4 | 36 |
| Table 4.8 | Overall performance of ML model with Feature Reduction for Dataset-4 | 37 |
| Table 4.9 | Summary Table for Dataset-1 | 38 |
| Table 4.10 | Summary Table for Dataset-2 | 38 |
| Table 4.11 | Summary Table for Dataset-3 | 39 |
| Table 4.12 | Summary Table for Dataset-4 | 39 |
| Table 4.13 | Compare propose model with previous model | 50 |

ABBREVIATIONS

- ML - Machine Learning
- DoS - Denial of Service
- DDoS – Distributed Denial of Service
- DRDoS- Distributed Reflection Denial of Service
- DNS - Domain Name Service
- NetBIOS - Network Basic Input / Output System.
- NTP- Network Time Protocol
- UDP - User Datagram Protocol
- SVM – Support Vector Machine
- DT – Decision Tree
- LR – Logistic Regression
- RF – Random Forest
- DL - Deep Learning
- PCA - Principal Component Analysis
- CIC - Canadian Institute for Cybersecurity.

Chapter 1

Introduction

Cyber-attacks are a major concern for many internet users in this day and age. Denial of service (DoS), distributed denial of service (DDoS), and distributed reflection denial of service (DRDoS) attacks are three common types of cyber-attacks. A growing class of distributed denial of service (DDoS) attacks is known as distributed reflection denial of service (DRDoS). DRDoS attacks are harder to defend against because of their characteristics and attack methodology. [1].

1.1 Distributed Reflection Denial of Service (DRDoS) attacks:

An attacker uses the victim's faked source IP address to make outdated requests to numerous servers in a DRDoS attack. The servers will respond by sending messages to the affected PC. And these responses are frequently (many times) larger than the victim's requests. The attacks are sometimes known as amplification attacks in such instances. The DrDoS attack approach has two major advantages: anonymity and amplification.

The DRDoS tactic was employed by 39 percent of attackers in all Distributed Denial of Service (DDoS) attacks. For more than a decade, DrDoS attacks have been a reliable and effective sort of DDoS attack. The method appears to be indestructible. Its effectiveness and popularity continue to rise. [14]

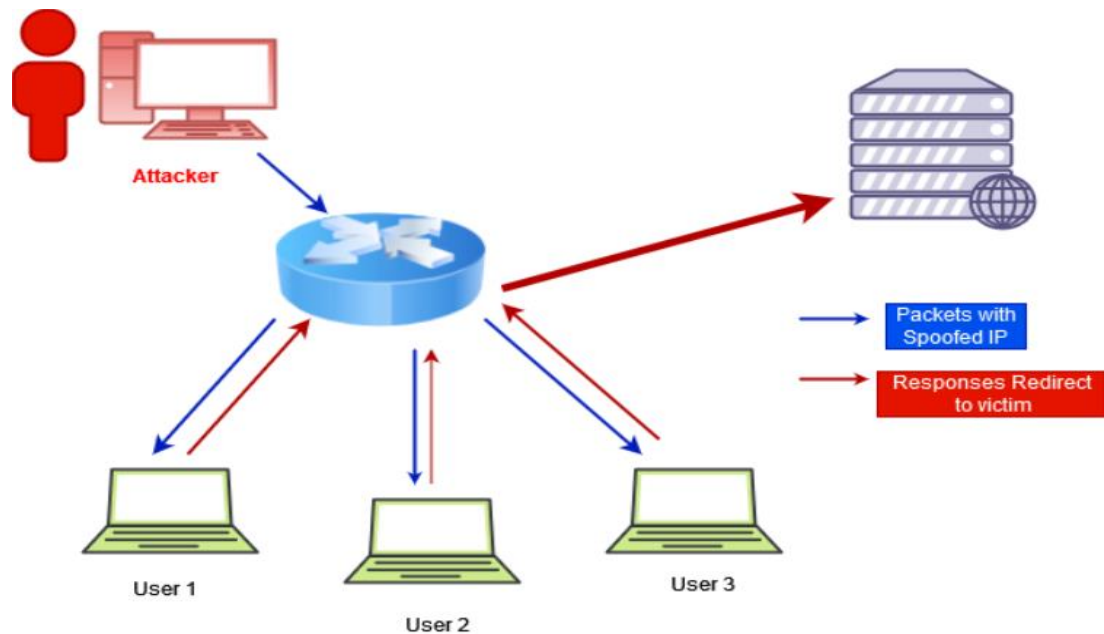


Fig. 1.1. The DRDoS Attack.

In order to launch assaults that result in traffic floods, take down websites, and disrupt networks at corporate targets, DRDoS attackers have been abusing the following protocols on Inter net-connected devices and servers:

1.1.1 Domain Name Service (DNS)

Users can type web addresses into browsers and have the IP address resolved by a DNS resolver to a DNS server thanks to the Domain Name Service (DNS), which transforms numerical IP addresses into domain names. Then DNS server response to the user with the IP address. The attacker has spoofed the user IP address first and then the attacker applies a DNS query by the spoofed IP to the DNS server. After that, the user gets some unwanted response from the DNS server which is called a DRDoS DNS attack.

1.1.2 NetBIOS

NetBIOS developed by IBM and then adopted by Microsoft is a network system providing some network services to the established connections between computers and allows software applications to communicate with each other on a LAN network. By using some distinct protocols named NetBIOS frames to make communication with network devices and transmission of data. NetBIOS stands for Network Basic Input / Output System that works at the Session layer (Layer 5) of the OSI model but it is often used over TCP/IP. Software applications use NetBIOS names with up to 16 characters. NetBIOS Session is started and terminated correspondingly by commands sent by clients. First of all, Cyber attackers send queries (or commands) to the victim's devices called NetBIOS queries. Then by spoofing, Cyber attackers make victim devices legitimate and also make them send a high number of requests to the target devices

1.1.3 Network Time Protocol (NTP)

The Network Time Protocol (NTP) is used by Internet servers to synchronize time. In DRDoS NTP attack, the attackers use the NTP server to send a MONLIST command to the victim machine. First attacker identifies the IP address of the victim. Then attacker will send a large amount of UDP packets to the NTP server. These packets support the monlist command, which gives the last number of IP addresses that use the NTP Server for their queries. The record of the queries will be sent to the source IP address which is spoofed by the attacker.

1.1.4. User Datagram Protocol(UDP)

User Datagram Protocol(UDP) is a connectionless protocol thus it's not required a host-to-host connection. UDP works at the transport layer by transmitting Datagrams over the network. Though UDP protocol is a connectionless protocol thus it does not inspect source IP addresses. That's why attackers choose this protocol. By forging IP packet Datagram, attackers include a free source IP address. The attacker copies the victim's IP address with a huge number of UDP packets then the target machine responds to the victim machine over the attacker.

1.2 Machine learning

A new method called machine learning enables computers to automatically learn from past data. In order to build mathematical models and generate predictions based on previously collected data or information, machine learning employs a range of methodologies. It is used for a wide range of operations, such as **speech recognition, email filtering, Facebook auto-tagging, recommender systems**, and more.

A subset of artificial intelligence known as "machine learning" is stated to focus on creating algorithms that enable computers to independently learn from data and past experiences. In 1959, Arthur Samuel coined the phrase "machine learning."

Machine learning can be used by cybersecurity systems to analyze trends and learn from them in order to assist prevent repeated attacks and respond to evolving behavior. It can let cybersecurity teams take proactive steps and respond to ongoing threats in real time with retaliation. It can speed up routine activities and aid businesses in making more efficient use of their resources.

In summary, machine learning has the potential to significantly improve cybersecurity by making it less complex, more proactive, and less expensive.

1.3 Research questions and Objectives

| *Why DRDoS attack detection is necessary?*

Now a day the most vital cyber security attack is the DRDoS attack. It can cause server outages and monetary loss. It also responsible for excessive stress on resources and the backbone of a network. So,It is necessary to protect internet users and network system from DRDoS attacker. 2nd fundamental question about this research is- *How can we detect the DRDoS attack?*

| *How can we detect the DRDoS attack?*

To detect DRDoS attack, we can use several machine learning algorithms to create a best classification model. There are so many machine learning classifiers use for classification problem like Logistic Regression(LR), Support Vector Machine(SVM), Decision Tree(DT), Random Forest(RF) etc. We can easily use this type of classifiers.

Which Classifier model is best to detect DRDoS attack detection?

We can use performance evaluation parameter (Accuracy, Precision, Recall, F1-score) to select the best model to detect the DRDoS attack.

Why use PCA for feature optimization in this research?

PCA means Principle Component Analysis. We have used PCA for dimensionality Reduction in our research. PCA is a technique of feature optimization.

1.4 Research synopsis

This research aims to make an efficient model to detect DRDoS attack. This efficient model can quickly and more accurately detect DRDoS attack. The internet users take necessary steps, if they detect DRDoS attack earlier.

1.5 Fundamental contributions

So we need to detect this type of attack to protect internet users. This paper discusses the application of the SVM, DT, RT, and LR for detecting DRDoS attacks without feature reduction and with feature reduction using Principal Component Analysis(PCA). The following are the key contributions of this paper:

- In this research, Used four datasets (DRDoS-DNS, DRDoS-NetBIOS, DRDoS-UDP, DRDoS-NTP) from Canadian Institute for Cybersecurity (CIC).
- four Machine Learning algorithms like SVM, LR, RF, and DT are used to detect the DRDoS attack.

- Apply SVM, LR, RF, and DT algorithms in two ways.
 - Firstly, use these four algorithms without considering the features reduction technique.
 - Secondly, use these four algorithms with considering the features reduction technique.
- Accuracy, Precision, F1-score, and Recall are used to gauge how well the suggested algorithm performs.

In The rest of the paper in section II: we discuss a similar type of work of DRDoS attack. In section III: we have to describe the method of detecting DRDoS attacks. In section IV: Show the experimental result with the table. Finally, give the conclusion in section V and the references in section VI.

1.6 Outline of the research

The outline of this research is partitioned into five chapters as follows. Firstly, chapter 1 gives a brief introduction to the research topic. The same chapter presents the motivation of this work and the related research questions and objectives to fulfill these questions. Secondly, chapter 2 provides a review of previous works on Distributed Reflection Denial of Service (DRDoS) Attacks detection and their pros and cons. Also, different research findings are also included in chapter 2. Thirdly, chapter 3 demonstrates the methodological approach proposed in this research with a brief overview of different techniques applied in this method. Fourthly, chapter 4 presents the outcomes of this research with discussion and comparison with existing methods. Finally, the conclusion and future scopes are discussed in chapter 5. At the end of this book, complete information sources are provided in the Reference section.

Chapter 2

Literature Review

2.1 Background

Several works have been proposed to develop a computerized approach to DRDoS attack detection. Most of the techniques use classical Machine Learning (ML) as well as Deep Learning (DL) to develop their scheme. Therefore, by taking existing works in mind this research has tried to find an improved way for the execution of accurate DRDoS attack detection than existing works.

By analyzing all the existing works, we have tried to get a solution for our work by applying the various Machine learning approach. Finding appropriate research gap is necessary to establish the core mechanism of any research so, for the purpose of better understanding and to get the appropriate research gap, all analyzed existing works are summarized in this chapter.

2.2 Exploration of Deep learning and Machine learning

Several computer-based methods for DRDoS attack identification are investigated and we have found most of them used DL and ML to develop their scheme. This section described the most recent DRDoS attack identification techniques.

The authors of the paper [1] proposed to detect DRDoS assaults, a novel proactive feature selection model based on better optimization techniques has been developed. They developed a model based on proactive feature selection (PFS) that used machine learning methods such as k-nearest neighbor (KNN), random forest (RF), and support vector machine (SVM) to predict DRDoS assaults. With the proposed model, they were able to attain an accuracy of 89.59 percent.

In the study [2] the authors wrote a review article titled "Distributed reflection denial of service attack: A critical review." They noted and discussed the distinctions between DRDoS attacks based on TCP and UDP in this article.

The authors [3] proposed an Improved Source IP Address Validation Method for 5G Networks to Prevent DRDoS Attacks. They devised a strategy for defeating DRDoS assaults by increasing the 5G core network's User Plane Function (UPF). They

demonstrated that the model may effectively deter DRDoS attacks by increasing the packet inspection rate (PIR).

The authors of reference [4] proposed a Deep Forest-Based DRDoS Detection and Defense Method in a Big Data Environment Deep Forest, IoT, Big Data, and other approaches were applied. They demonstrated that their technology is both successful for DRDoS detection and defense, with a greater detection rate and lower false alarm rate.

The authors of the paper[5] proposed Cloud-Base Defense against DRDoS Attacks. They published a paper to prevent DRDoS attacks by applying a cloud.

The authors[6] suggested in an essay titled "Visualization of Actionable Knowledge to Mitigate DRDoS Attacks," they provided an approach and technology for Internet Service Providers (ISPs) to put an end to DRDoS attacks.

The author of the paper [7] suggested a machine learning-based method for identifying and assessing DRDoS attacks. Calculations were made for the overall success rate, detection rate, and false positive rate.

The authors of the paper [8] proposed for DRDoS attacks, the Protocol Independent Detection and Classification (PIDC) System is used. They used data mining and machine learning algorithms like the C4.5 classification algorithm in the PIDC system. This system had a true positive rate of 99 percent and a false positive rate of less than 1%.

The authors of the paper [9] proposed, "A simple response packet confirmation system was developed as a model for detecting DRDoS assaults" . They proposed a simple and reliable approach with little processing costs.

In the study of the paper [10], the authors proposed a method to detect RA (Reflection Amplification) attacks based on UDP Protocol by classifying and analyzing the RA attacks. They also used detection algorithm and found the reliabilities of these algorithms.

The authors of the paper [11] proposed "CARD (Continuous and Random Dropping) based DRDOS Attack Detection and Prevention Techniques in MANET", they published a mechanism for detecting and preventing DrDoS attacks based on CARD (Continuous and Random Dropping) technology in MANET network.

The authors of the paper [12] proposed "An Integrated Approach for DRDoS Attacks Using E-RED and ANT Classification Methods." They released a technique to recognize and categorize DRDoS attacks using the Enhanced-Random Early Detection (E-RED) algorithm and the Application-based Network Traffic (ANT) classification method, which

recognizes 99 percent true positives and 1 percent false positives while categorizing attacks with 98 percent accuracy.

2.3. Summery table of all previous research work

| Ref. | Model | Algorithm /Method | Major findings |
|------|---|---|--|
| 1 | To detect DRDoS assaults, a novel proactive feature selection model based on better optimization techniques has been developed. | nature-inspired optimization algorithm, (KNN), (RF), and (SVM). | 89.59% accuracy. |
| 4 | In the Big Data Environment, A Deep Forest-Based DRDoS Detection and Defense Method | Deep Forest, IoT, Big Data, and other techniques. | A higher rate of detection and false alarms. |
| 7 | An Evaluation of a Machine Learning Approach for Detecting DRDoS Attacks | Support Vector Machine algorithm (SVM).. | The overall success rate, as well as the detection and false positive rates. |
| 8 | For DRDoS attacks, the Protocol Independent Detection and Classification (PIDC) System is used. | Data mining concept and machine learning algorithms like C4.5 classification algorithm. | True positive rates of 99 percent and false positive rates of less than 1% (1 percent). |
| 10 | Based on the UDP Protocol, detect the Reflection Amplification Attack. | detection algorithm | -- |
| 11 | DRDOS Attack Detection and Prevention Methods Using CARD in the MANET (Continuous and Random Dropping). | CARD(Continuous and Random Dropping) technology. | -- |
| 12 | DRDoS Attack Classification Methods Using an Integrated Approach of E-RED and ANT. | Enhanced-Random Application-based Network Traffic (ANT) classification method Early. | Detects 99 percent of true positives and 1% of false positives, and classifies attacks with a classification accuracy of 98 percent. |

2.4 Research gap

By analyzing several existing works, it seems that most of the computerized methods of detecting DRDoS attack are based on classical ML and DL and it also seems that no work is used feature optimization techniques. The machine learning or deep learning model using feature optimization technique might be the best approach to quickly detect DRDoS attack.

Chapter 3

Research Method

This chapter presents the principal mechanism of our research along with the related data acquiring process. In this research, several numbers of dataset are used to detect the DRDoS attack detection. Firstly, all the datasets need to be pre-processed. We have used the label encoding technique for convert the string data to numerical data. We have done this whole work by following two methods. Method-1, we detect the DRDoS attack without using the feature reduction technique, and Method-2, we detect the DRDoS attack by using feature reduction by PCA technique.

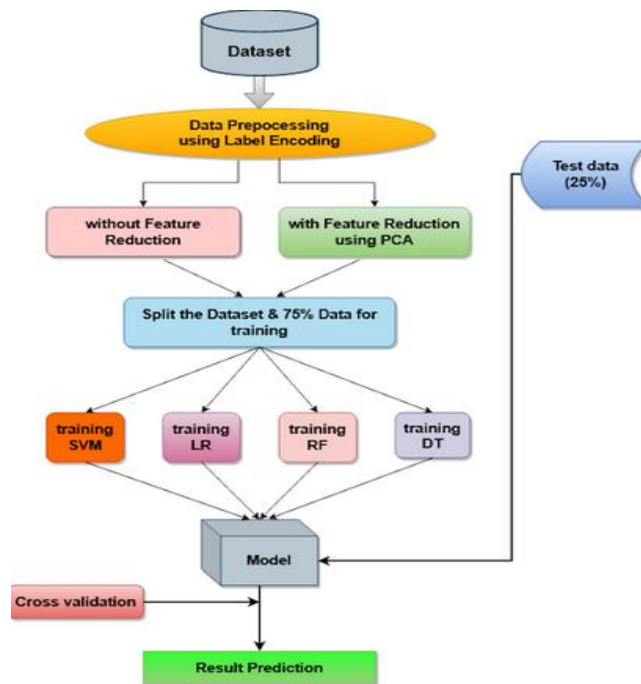


Figure 3.1: Proposed Model.

3.1 Dataset

The datasets have used in this work are publicly available in Canadian Institute for Cybersecurity (CIC). We have used Four datasets with 100000 data.

Table 3.1: Dataset.

| Dataset Number | Dataset Name | Feature Size |
|----------------|---------------|--------------|
| Dataset-1 [27] | DRDoS-DNS | 88 |
| Dataset-2 [27] | DRDoS-NetBIOS | 88 |
| Dataset-3 [27] | DRDoS-UDP | 89 |
| Dataset-4 [27] | DRDoS-NTP | 89 |

3.2 Proposed method

The full framework of our plan is shown in Figure 3.1, and sections 3.2.1 to 3.2.5 of that figure cover the specific working methods.

3.2.1 Method-1 (Without Feature Reduction):

After completing the pre-processing part, split the dataset for training and testing. We have to use 75% data of a dataset for training and 25% of data for testing. 75% of data have used to training by the four machine learning algorithms (SVM, LR, RF, DT) to make the individual model. Then use test data for testing the performance of the individual model. In this method, we have not used the feature reduction technique.

3.2.2 Method-2 (With Feature Reduction using PCA):

We have to use Principal Component Analysis (PCA) technique for feature reduction. In PCA, some features select from the dataset which is 95% important to making the decision. After that, 75% data have used for training by the four machine learning algorithms (SVM, LR, RF, DT) for making the individual model. Then use test data for testing the performance of the individual model.

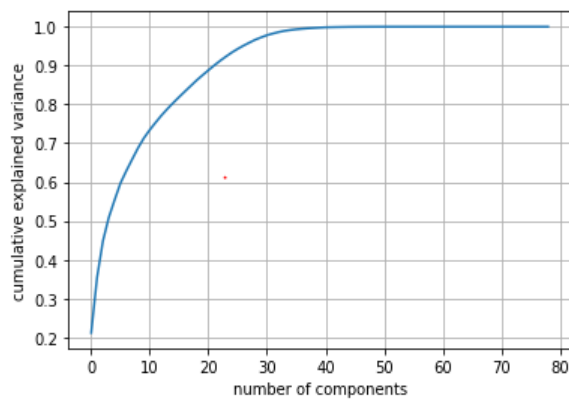
3.2.3. Feature Optimization Technique:

3.2.3.1. *Principal Component Analysis(PCA):*

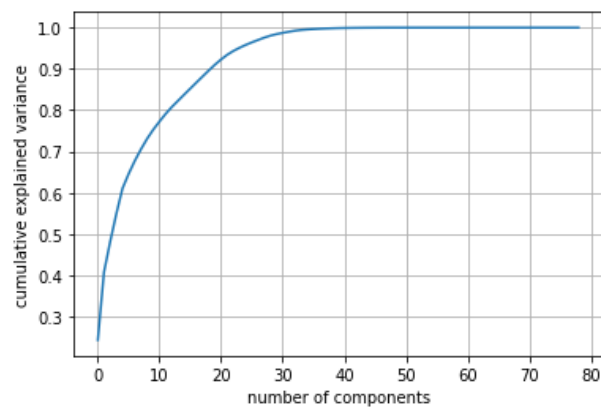
Algorithm for unsupervised machine learning using principal components (PCA).It is employed to lessen a dataset's dimensionality.It is a statistical method that uses orthogonal transformation to turn observations of correlated features into a collection of linearly uncorrelated data.The Principal Components are the recently modified features.

Following steps for PCA:

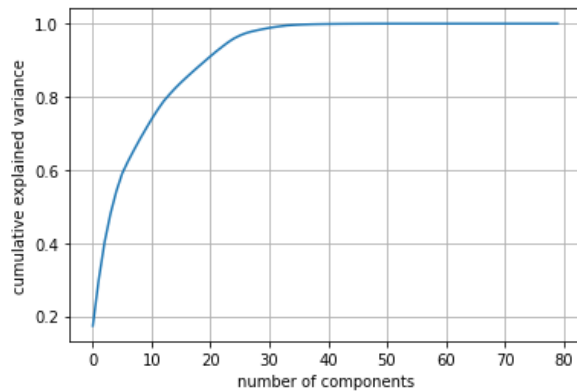
- Acquiring the data set.
- structuring data representation
- data standardization
- determining the covariance
- Making the Eigen Values and Eigen Vectors calculations
- Classifying the Eigen Vectors
- determining the new characteristics or principal components
- Eliminate less significant or irrelevant features from the new dataset.



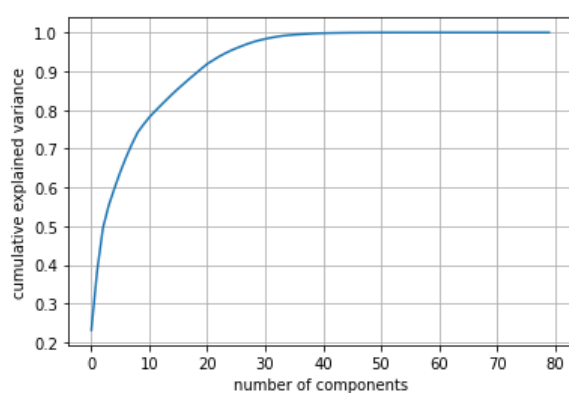
Dataset-1



Dataset-2



Dataset-3



Dataset-4

Figure 3.2: Covariance Curve for four Dataset

3.3. Machine Learning Algorithm:

Machine learning algorithms make it easy to detect the DRDoS attack. There are four machine learning algorithms like SVM, DT, RF, and LR are used to detect the DRDoS attack. We discuss all the algorithm one by one.

3.3.1. Decision Tree Algorithm(DT):

A well-known machine learning algorithm is the decision tree. It is employed in the classification of data. It's a tree-structured algorithm in which the features of a database are specified by internal nodes and branches that represent decision rules, with each leaf node indicating the outcome. We create a Decision tree by following the steps below.

- Select the desired attribute.
- Determine Information Gain (I.G.) for the desired attribute.
- Use the formula below to calculate the entropy of the other properties:

$$\text{Entropy}(s) = \sum_{i=1}^n -(P_i \log_2 P_i) \quad (1)$$

- To determine the Gain, deduct the Entropy(s) from the Information Gain (I.G) of each characteristic (G)

$$\text{Gain}(S, A) = \text{Entropy}(s) - \sum_{i=1}^n \frac{s_v}{s} \times \text{Entropy}(S_v) \quad (2)$$

3.3.2. Support Vector Machine(SVM):

This machine learning approach can be used to do both classification and regression tasks. The SVM method attempts to partition an n-dimensional area into categories by constructing near-perfect lines or decision boundaries. As a result, certain data points can be quickly assigned to the correct category in the future. Support vectors are the locations closest to the decision boundaries, and they aid in calculating the decision plane's margin.

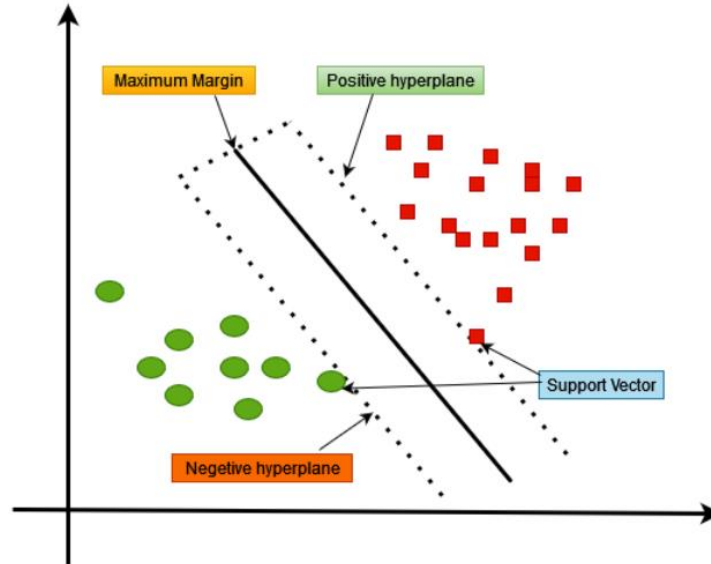


Figure 3.3 Support Vector Machine(SVM)

3.3.3. Random Forest(RF):

RF is a supervise learning algorithm in machine learning. In this model, different decision trees are trained based on the dataset. This model is called an ensemble of decision trees because a large number of decision trees present to make the decision.

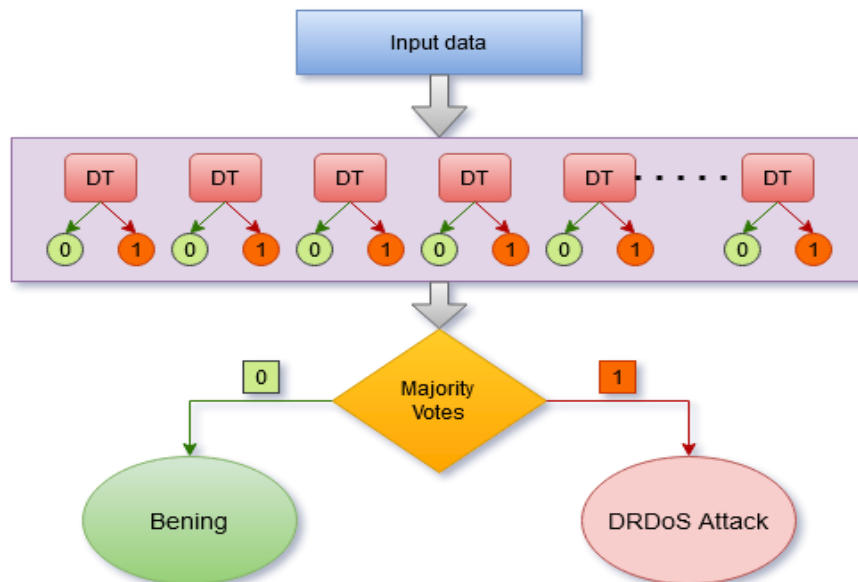


Figure 3.4. Random Forest (RF)

3.3.4. Logistic Regression(LR):

LR is a machine learning technique that employs classification to determine if an input is benign or not. The maximum-entropy classifier (MaxEnt), logit regression, and log-linear classifier are other names for the same algorithm. It makes use of the sigmoid function. It calculates the linear equation denoted by equation 4 for the data values $A=A_1, A_2, A_3, \dots A_n$. The logit function is used to squeeze the outcome of a linear equation between 0 and 1 using equation 5. Equation 3 calculates the regression coefficient W^T maximum likelihood estimation.

$$W^T = \max \sum_{j=1}^n (Y_j \times W_j A_j) \quad (3)$$

$$r = Y_j \times W^T A_j \quad (4)$$

$$P(r) = \frac{1}{1 + \exp^{-r}} \quad (5)$$

When $P(r) > 0.5$, the probability of the input instance being a DRDoS assault is indicated, and when $P(r) < 0.5$, the probability of the input instance being benign is shown.

The data points are denoted by Y_j .

3.3.5. Performance parameters:

We have used four data set and all the dataset have been trained and tested by the different machine learning models. We measured the performance of each model by calculating Accuracy, Precision, Recall and F1 score. To calculate the performance parameters, we have used the confusion matrix. Confusion matrix is not a performance matrix but to build the other performance matrix we need this confusion matrix. We have showed the confusion matrix in figure.

Table 3.2: Confusion matrix table

| Actual Data | | | |
|----------------|------------|----------------------------|----------------------------|
| Predicted Data | | (Positive) | (Negative) |
| | (Positive) | True Positive (TP) | False Positive (FP) |
| | (Negative) | False Negative (FN) | True Negative (TN) |

In a confusion matrix there are four parameters. Now we describe it below:

- **True Positive:** When algorithm predict that is true and if it is actually true then it is called True Positive.
- **False Positive:** When algorithm predict that is true but if it is not actually true then it is called False Positive.
- **False Negative:** When algorithm predict that is false but if it is not actually false then it is called False Negative.
- **True Negative:** When algorithm predict that is false and if it is actually false then it is called True Negative.

The important parameters for performance measure of all the models are Accuracy, Precision, Recall, F1-score. In Table [2], we have described all of this parameters.

Table 3.3: Performance Measurement Parameters.

| Metrics | Formula | Meaning |
|-----------|--|--|
| Accuracy | $\frac{TP + TN}{(TP + FP) + (FN + TN)} \times 100$ | Rate of overall correct prediction |
| Precision | $\frac{TP}{(TP + FP)} \times 100$ | Correct Edible prediction cases out of all Edible cases predicted. |
| Recall | $\frac{TP}{(TP + FN)} \times 100$ | Percentage of Edible prediction cases that are correct. |
| F1-score | $2 \times \frac{Recall \times Precision}{Recall + Precision} \times 100$ | Indicates how accurate and reliable a classifier is. |

Chapter 4

Result and Discussion

In this section, we have discussed about the result of all four models for all the datasets. First we have discussed about the result with five folds of all models without feature reduction and with feature reduction using PCA for individual dataset. Second We show the summary of all results in Table [10].

4.1 Analysis the result for Dataset-1:

In table-4.1, shows that the performance of ML models using without Feature Reduction for Dataset-1. In this observation, the Random Forest(RF) and the Decision Tree(DT) both classifiers show the best result with 100% accuracy using without feature reduction. But the other two classifiers Support Vector Machine(SVM) and Logistic Regression(LR) both show 98% and 99% accuracy using without Feature Reduction and RF and DT show 100% F1-score only using without feature reduction.

Table 4.1: Overall performance of ML model without Feature Reduction for Dataset-1.

| Dataset-1 | | | | | |
|-----------|-------------|-------------|--------------|--------------|--------------|
| Model | Fold | Accuracy | Precision | Recall | F1 score |
| SVM | Fold 1 | 98% | 50% | 50% | 49% |
| | Fold 2 | 98% | 49% | 50% | 50% |
| | Fold 3 | 98% | 49% | 50% | 50% |
| | Fold 4 | 98% | 99% | 50% | 50% |
| | Fold 5 | 98% | 49% | 50% | 50% |
| | Mean | 98% | 59.2% | 50% | 49.8% |
| LR | Fold 1 | 99% | 88% | 87% | 88% |
| | Fold 2 | 99% | 88% | 87% | 87% |
| | Fold 3 | 99% | 91% | 88% | 89% |
| | Fold 4 | 99% | 89% | 87% | 88% |
| | Fold 5 | 99% | 90% | 89% | 89% |
| | Mean | 99% | 89.4% | 87.6% | 88.2% |
| RF | Fold 1 | 100% | 100% | 100% | 100% |
| | Fold 2 | 100% | 100% | 100% | 100% |
| | Fold 3 | 100% | 100% | 100% | 100% |
| | Fold 4 | 100% | 100% | 100% | 100% |
| | Fold 5 | 100% | 100% | 100% | 100% |
| | Mean | 100% | 100% | 100% | 100% |
| | Fold 1 | 100% | 100% | 100% | 100% |
| | Fold 2 | 100% | 100% | 100% | 100% |

| | | | | | |
|----|-------------|-------------|-------------|-------------|-------------|
| DT | Fold 3 | 100% | 100% | 100% | 100% |
| | Fold 4 | 100% | 100% | 100% | 100% |
| | Fold 5 | 100% | 100% | 100% | 100% |
| | Mean | 100% | 100% | 100% | 100% |

In table-4.2, shows that the performance of ML models with Feature Reduction using PCA for Dataset-1. In this observation, the Random Forest(RF) and the Decision Tree(DT) both classifiers show the best result with 100% accuracy using this method. But the other two classifier Support Vector Machine(SVM) and Logistic Regression(LR) both show 99% and 98% accuracy using with feature reduction.

Table 4.2: Overall performance of ML model using PCA for Dataset-1.

| Dataset-1 | | | | | |
|-----------|-------------|-------------|--------------|--------------|--------------|
| Model | Fold | Accuracy | Precision | Recall | F1 score |
| SVM | Fold 1 | 99% | 97% | 59% | 65% |
| | Fold 2 | 99% | 99% | 61% | 68% |
| | Fold 3 | 99% | 99% | 60% | 67% |
| | Fold 4 | 99% | 99% | 61% | 67% |
| | Fold 5 | 99% | 99% | 60% | 66% |
| | Mean | 99% | 98.6% | 60.2% | 66.6% |
| LR | Fold 1 | 98% | 93% | 56% | 60% |
| | Fold 2 | 98% | 95% | 56% | 61% |
| | Fold 3 | 98% | 95% | 56% | 61% |
| | Fold 4 | 98% | 94% | 57% | 64% |
| | Fold 5 | 98% | 91% | 57% | 61% |
| | Mean | 98% | 93.6% | 56.4% | 61.4% |
| RF | Fold 1 | 100% | 98% | 96% | 97% |
| | Fold 2 | 100% | 98% | 97% | 97% |
| | Fold 3 | 100% | 99% | 97% | 98% |
| | Fold 4 | 100% | 99% | 96% | 97% |
| | Fold 5 | 100% | 98% | 96% | 97% |
| | Mean | 100% | 98.4% | 96.4% | 97.2% |
| DT | Fold 1 | 100% | 98% | 96% | 97% |
| | Fold 2 | 100% | 97% | 97% | 97% |
| | Fold 3 | 100% | 99% | 97% | 98% |
| | Fold 4 | 100% | 98% | 96% | 97% |
| | Fold 5 | 100% | 97% | 96% | 96% |
| | Mean | 100% | 97.8% | 96.4% | 97% |

4.2. Analysis the result for Dataset-2:

In Table-4.3, Shows that the performance of ML model using without Feature Reduction for Dataset-2. In this observation, the Random Forest(RF) and the Decision Tree(DT) both classifiers show the best result with 100% Accuracy using without feature reduction. But the other two classifier Support Vector Machine(SVM) and Logistic Regression(LR) both show 98% and 99% accuracy using without Feature Reduction and RF and DT show 100% F1-score only using without feature reduction.

Table 4.3: Overall performance of ML model without Feature Reduction for Dataset-2.

| Dataset-2 | | | | | |
|-----------|-------------|---------------|--------------|--------------|--------------|
| Model | Fold | Accuracy | Precision | Recall | F1 score |
| SVM | Fold 1 | 99% | 50% | 51% | 100% |
| | Fold 2 | 99.50% | 50% | 50% | 50% |
| | Fold 3 | 99% | 50% | 50% | 50% |
| | Fold 4 | 99% | 100% | 51% | 51% |
| | Fold 5 | 99.46% | 50% | 50% | 50% |
| | Mean | 99.19% | 60% | 50.4% | 60.2% |
| LR | Fold 1 | 100% | 79% | 73% | 92% |
| | Fold 2 | 99.75% | 89% | 85% | 87% |
| | Fold 3 | 100% | 95% | 76% | 83% |
| | Fold 4 | 100% | 87% | 85% | 86% |
| | Fold 5 | 99.70% | 89% | 81% | 84% |
| | Mean | 99.89% | 87.8% | 80% | 86.4% |
| RF | Fold 1 | 100% | 100% | 100% | 100% |
| | Fold 2 | 99.99% | 100% | 100% | 100% |
| | Fold 3 | 100% | 100% | 100% | 100% |
| | Fold 4 | 100% | 100% | 100% | 100% |
| | Fold 5 | 100% | 100% | 100% | 100% |
| | Mean | 100% | 100% | 100% | 100% |
| DT | Fold 1 | 100% | 100% | 100% | 100% |
| | Fold 2 | 100% | 100% | 100% | 100% |
| | Fold 3 | 100% | 100% | 100% | 100% |
| | Fold 4 | 100% | 100% | 100% | 100% |
| | Fold 5 | 100% | 100% | 100% | 100% |
| | Mean | 100% | 100% | 100% | 100% |

In table-4.4, also shows that the performance of ML models with Feature Reduction using PCA for Dataset-2. In this observation, the Random Forest(RF) and the Decision Tree(DT) both classifiers show the best result with 100% accuracy using this method. But the other two classifier Support Vector Machine(SVM) and Logistic Regression(LR) both show 99.48% and 98.50% accuracy using with feature reduction.

Table 4.4: Overall performance of ML model using PCA for Dataset-2.

| Dataset-2 | | | | | |
|------------------|-------------|-----------------|------------------|---------------|-----------------|
| Model | Fold | Accuracy | Precision | Recall | F1 score |
| SVM | Fold 1 | 99.45% | 100% | 54% | 57% |
| | Fold 2 | 99.43% | 100% | 54% | 58% |
| | Fold 3 | 99.46% | 100% | 55% | 59% |
| | Fold 4 | 99.55% | 100% | 57% | 62% |
| | Fold 5 | 99.53% | 100% | 55% | 60% |
| | Mean | 99.48% | 100% | 55% | 59.2% |
| | | | | | |
| LR | Fold 1 | 99.48% | 92% | 58% | 63% |
| | Fold 2 | 99.47% | 98% | 58% | 64% |
| | Fold 3 | 99.49% | 94% | 58% | 64% |
| | Fold 4 | 99.55% | 97% | 58% | 63% |
| | Fold 5 | 99.55% | 97% | 58% | 63% |
| | Mean | 99.50% | 95.6% | 58% | 63.40% |
| | | | | | |
| RF | Fold 1 | 100% | 95% | 96% | 96% |
| | Fold 2 | 100% | 98% | 98% | 98% |
| | Fold 3 | 100% | 99% | 97% | 98% |
| | Fold 4 | 100% | 97% | 98% | 97% |
| | Fold 5 | 100% | 98% | 98% | 98% |
| | Mean | 100% | 97.4% | 97.4% | 97.4% |
| | | | | | |
| DT | Fold 1 | 100% | 96% | 96% | 96% |
| | Fold 2 | 100% | 96% | 98% | 97% |
| | Fold 3 | 100% | 97% | 96% | 97% |
| | Fold 4 | 100% | 96% | 97% | 97% |
| | Fold 5 | 100% | 96% | 97% | 97% |
| | Mean | 100% | 96.20% | 96.80% | 96.8% |

4.3. Analysis the result for Dataset-3:

In Table-4.5, again shows that the Random Forest(RF) and the Decision Tree(DT) both classifiers show the best result with 100% accuracy without feature reduction. But the other two classifier Support Vector Machine(SVM) and Logistic Regression(LR) both show 99.09% and 99.40% accuracy using without Feature Reduction.

Table 4.5: Overall performance of ML model without Feature Reduction for Dataset-3.

| Dataset-3 | | | | | |
|------------------|-------------|-----------------|------------------|---------------|-----------------|
| Model | Fold | Accuracy | Precision | Recall | F1 score |
| SVM | Fold 1 | 99% | 50% | 50% | 50% |
| | Fold 2 | 99.12% | 100% | 50% | 50% |
| | Fold 3 | 99.07% | 100% | 50% | 50% |
| | Fold 4 | 99.17% | 100% | 50% | 50% |
| | Fold 5 | 99.11% | 100% | 50% | 51% |
| | Mean | 99.09% | 90% | 50% | 50.2% |
| | | | | | |
| LR | Fold 1 | 99% | 82% | 76% | 91% |
| | Fold 2 | 99.57% | 92% | 81% | 86% |
| | Fold 3 | 99.47% | 92% | 77% | 83% |
| | Fold 4 | 99.49% | 93% | 78% | 85% |
| | Fold 5 | 99.50% | 93% | 77% | 83% |
| | Mean | 99.40% | 90.4% | 77.8% | 85.2% |
| | | | | | |
| RF | Fold 1 | 100% | 100% | 100% | 100% |
| | Fold 2 | 100% | 100% | 100% | 100% |
| | Fold 3 | 100% | 100% | 100% | 100% |
| | Fold 4 | 100% | 100% | 100% | 100% |
| | Fold 5 | 100% | 100% | 100% | 100% |
| | Mean | 100% | 100% | 100% | 100% |
| | | | | | |
| DT | Fold 1 | 100% | 100% | 100% | 100% |
| | Fold 2 | 100% | 100% | 100% | 100% |
| | Fold 3 | 100% | 100% | 100% | 100% |
| | Fold 4 | 100% | 100% | 100% | 100% |
| | Fold 5 | 100% | 100% | 100% | 100% |
| | Mean | 100% | 100% | 100% | 100% |

In Table-4.6, the Random Forest(RF) and the Decision Tree(DT) both classifiers also show the best result with 100% accuracy without feature reduction. But the other two classifier Support Vector Machine(SVM) and Logistic Regression(LR) achieve 99.13% and 99.04% accuracy using with Feature Reduction.

Table 4.6: Overall performance of ML model using PCA for Dataset-3.

| Dataset-3 | | | | | |
|------------------|-------------|-----------------|------------------|---------------|-----------------|
| Model | Fold | Accuracy | Precision | Recall | F1 score |
| SVM | Fold 1 | 99.12% | 100% | 51% | 52% |
| | Fold 2 | 99.09% | 100% | 51% | 52% |
| | Fold 3 | 99.14% | 100% | 51% | 52% |
| | Fold 4 | 99.16% | 100% | 52% | 53% |
| | Fold 5 | 99.16% | 100% | 51% | 53% |
| | Mean | 99.13% | 100% | 51.2% | 52.4% |
| | | | | | |
| LR | Fold 1 | 99.02% | 71% | 61% | 64% |
| | Fold 2 | 98.99% | 71% | 60% | 64% |
| | Fold 3 | 99.07% | 73% | 63% | 67% |
| | Fold 4 | 99.07% | 71% | 61% | 65% |
| | Fold 5 | 99.05% | 70% | 62% | 65% |
| | Mean | 99.04% | 71.2% | 61.4% | 65% |
| | | | | | |
| RF | Fold 1 | 100% | 97% | 93% | 94% |
| | Fold 2 | 100% | 97% | 92% | 95% |
| | Fold 3 | 100% | 96% | 92% | 94% |
| | Fold 4 | 100% | 97% | 93% | 95% |
| | Fold 5 | 100% | 96% | 92% | 94% |
| | Mean | 100% | 96.6% | 92.4% | 94.6% |
| | | | | | |
| DT | Fold 1 | 100% | 92% | 93% | 93% |
| | Fold 2 | 100% | 93% | 92% | 93% |
| | Fold 3 | 100% | 92% | 92% | 92% |
| | Fold 4 | 100% | 93% | 93% | 93% |
| | Fold 5 | 100% | 93% | 93% | 93% |
| | Mean | 100% | 92.6% | 92.6% | 92.8% |

4.4. Analysis the result for Dataset-4:

Finally, in table-4.7, the Random Forest(RF) and the Decision Tree(DT) both classifiers show the best result with 100% accuracy and the other two classifier Support Vector Machine(SVM) and Logistic Regression(LR) both shows 87.60% and 98.75% accuracy using without Feature Reduction.

Table 4.7: Overall performance of ML model without Feature Reduction for Dataset-4.

| Dataset-4 | | | | | |
|-----------|--------|----------|-----------|--------|----------|
| Model | Fold | Accuracy | Precision | Recall | F1 score |
| SVM | Fold 1 | 88% | 97% | 50% | 47% |
| | Fold 2 | 87.57% | 94% | 50% | 47% |
| | Fold 3 | 87.50% | 94% | 50% | 47% |
| | Fold 4 | 87.46% | 94% | 50% | 47% |
| | Fold 5 | 87.48% | 95% | 50% | 47% |
| | Mean | 87.60% | 94.4% | 50% | 47% |
| LR | Fold 1 | 99% | 97% | 98% | 95% |
| | Fold 2 | 98.67% | 96% | 98% | 97% |
| | Fold 3 | 98.63% | 96% | 98% | 97% |
| | Fold 4 | 98.71% | 97% | 98% | 96% |
| | Fold 5 | 98.76% | 96% | 98% | 97% |
| | Mean | 98.75% | 96.40% | 98% | 96.4% |
| RF | Fold 1 | 100% | 100% | 100% | 100% |
| | Fold 2 | 100% | 100% | 100% | 100% |
| | Fold 3 | 100% | 100% | 100% | 100% |
| | Fold 4 | 100% | 100% | 100% | 100% |
| | Fold 5 | 100% | 100% | 100% | 100% |
| | Mean | 100% | 100% | 100% | 100% |
| DT | Fold 1 | 100% | 100% | 100% | 100% |
| | Fold 2 | 100% | 100% | 100% | 100% |
| | Fold 3 | 100% | 100% | 100% | 100% |
| | Fold 4 | 100% | 100% | 100% | 100% |
| | Fold 5 | 100% | 100% | 100% | 100% |
| | Mean | 100% | 100% | 100% | 100% |

But in table 4.8, the Random Forest(RF) and the Decision Tree(DT) both classifiers show 98.79% and 98.50% accuracy and other two classifier Support Vector Machine(SVM) and Logistic Regression(LR) both shows 93.61 % and 88.80% accuracy with Feature Reduction using PCA technique respectively.

Table 4.8: Overall performance of ML model using PCA for Dataset-4.

| Dataset-4 | | | | | |
|------------------|-------------|-----------------|------------------|---------------|-----------------|
| Model | Fold | Accuracy | Precision | Recall | F1 score |
| SVM | Fold 1 | 91.33% | 95% | 65% | 71% |
| | Fold 2 | 94.10% | 96% | 77% | 83% |
| | Fold 3 | 94.18% | 96% | 77% | 83% |
| | Fold 4 | 94.43% | 96% | 78% | 84% |
| | Fold 5 | 94.04% | 96% | 77% | 83% |
| | Mean | 93.61% | 95.8% | 74.8% | 80.8% |
| | | | | | |
| LR | Fold 1 | 88.86% | 92% | 56% | 57% |
| | Fold 2 | 88.86% | 92% | 56% | 57% |
| | Fold 3 | 88.94% | 92% | 56% | 58% |
| | Fold 4 | 88.87% | 91% | 55% | 57% |
| | Fold 5 | 88.48% | 91% | 56% | 57% |
| | Mean | 88.80% | 91.6% | 55.8% | 57.8% |
| | | | | | |
| RF | Fold 1 | 99.03% | 98% | 97% | 98% |
| | Fold 2 | 98.66% | 97% | 97% | 97% |
| | Fold 3 | 98.79% | 97% | 97% | 97% |
| | Fold 4 | 98.76% | 97% | 97% | 97% |
| | Fold 5 | 98.74% | 98% | 97% | 97% |
| | Mean | 98.79% | 97.4% | 97% | 97.8% |
| | | | | | |
| DT | Fold 1 | 98.82% | 97% | 97% | 97% |
| | Fold 2 | 98.34% | 96% | 96% | 96% |
| | Fold 3 | 98.47% | 96% | 97% | 96% |
| | Fold 4 | 98.46% | 96% | 97% | 96% |
| | Fold 5 | 98.45% | 97% | 97% | 97% |
| | Mean | 98.50% | 96.4% | 96.8% | 96% |

4.5. Summary of all models performance for all Datasets:

Table 4.9. Summary Table for Dataset-1:

| Dataset-1 | | | | | |
|-----------|----------------------------------|-------------|-------------|-------------|-------------|
| Model | | Accuracy | Precision | Recall | F1 score |
| SVM | Without Feature Reduction | 98% | 59.2% | 50% | 49.8% |
| | With Feature Reduction Using PCA | 99% | 98.6% | 60.2% | 66.6% |
| LR | Without Feature Reduction | 99% | 89.4% | 87.6% | 88.2% |
| | With Feature Reduction Using PCA | 98% | 93.6% | 56.4% | 61.4% |
| RF | Without Feature Reduction | 100% | 100% | 100% | 100% |
| | With Feature Reduction Using PCA | 100% | 98.4% | 96.4% | 97.2% |
| DT | Without Feature Reduction | 100% | 100% | 100% | 100% |
| | With Feature Reduction Using PCA | 100% | 97.8% | 96.4% | 97% |

Table 4.10. Summary Table for Dataset-2:

| Dataset-2 | | | | | |
|-----------|----------------------------------|-------------|-------------|-------------|-------------|
| Model | | Accuracy | Precision | Recall | F1 score |
| SVM | Without Feature Reduction | 99.19% | 60% | 50.4% | 60.2% |
| | With Feature Reduction Using PCA | 99.48% | 100% | 55% | 59.2% |
| LR | Without Feature Reduction | 99.89% | 87.8% | 80% | 86.4% |
| | With Feature Reduction Using PCA | 99.50% | 95.6% | 58% | 63.40% |
| RF | Without Feature Reduction | 100% | 100% | 100% | 100% |
| | With Feature Reduction Using PCA | 100% | 97.4% | 97.4% | 97.4% |
| DT | Without Feature Reduction | 100% | 100% | 100% | 100% |
| | With Feature Reduction Using PCA | 100% | 96.20% | 96.80% | 96.8% |

Table 4.11. *Summary Table for Dataset-3:*

| Dataset-3 | | | | | |
|------------------|----------------------------------|-----------------|------------------|---------------|-----------------|
| Model | | Accuracy | Precision | Recall | F1 score |
| SVM | Without Feature Reduction | 99.09% | 90% | 50% | 50.2% |
| | With Feature Reduction Using PCA | 99.13% | 100% | 51.2% | 52.4% |
| LR | Without Feature Reduction | 99.40% | 90.4% | 77.8% | 85.2% |
| | With Feature Reduction Using PCA | 99.04% | 71.2% | 61.4% | 65% |
| RF | Without Feature Reduction | 100% | 100% | 100% | 100% |
| | With Feature Reduction Using PCA | 100% | 96.6% | 92.4% | 94.6% |
| DT | Without Feature Reduction | 100% | 100% | 100% | 100% |
| | With Feature Reduction Using PCA | 100% | 92.6% | 92.6% | 92.8% |

Table 4.12. *Summary Table for Dataset-4:*

| Dataset-4 | | | | | |
|------------------|----------------------------------|-----------------|-----------------|---------------|------------------|
| Model | | Accuracy | F1 score | Recall | Precision |
| SVM | Without Feature Reduction | 87.60% | 94.4% | 50% | 47% |
| | With Feature Reduction Using PCA | 93.61% | 95.8% | 74.8% | 80.8% |
| LR | Without Feature Reduction | 98.75% | 96.40% | 98% | 96.4% |
| | With Feature Reduction Using PCA | 88.80% | 91.6% | 55.8% | 57.8% |
| RF | Without Feature Reduction | 100% | 100% | 100% | 100% |
| | With Feature Reduction Using PCA | 98.79% | 97.4% | 97% | 97.8% |
| DT | Without Feature Reduction | 100% | 100% | 100% | 100% |
| | With Feature Reduction Using PCA | 98.50% | 96.4% | 96.8% | 96% |

4.5. Bar Chart Representation

Now we have to show the result of all models of machine learning of all datasets by bar chart. The figure is given below:

From the bar chart, The Random Forest (RF) and Decision Tree(DT) achieve the highest accuracy for Dataset-1 and Dataset-2.

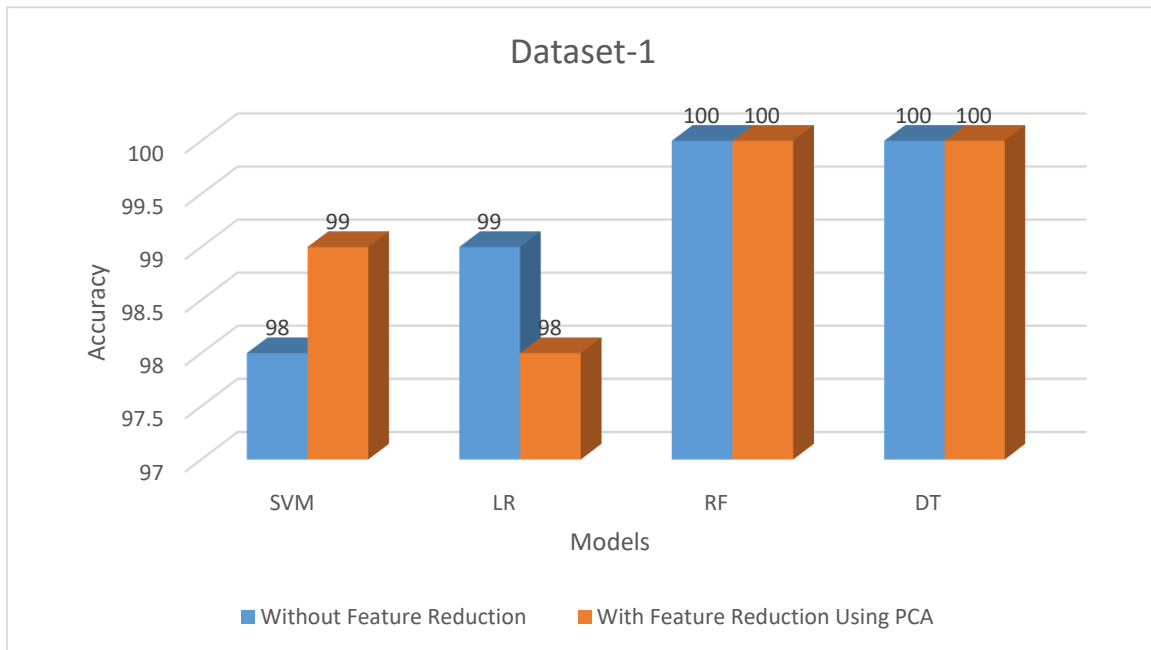


Figure 4.1: Performance of all models for dataset 1.

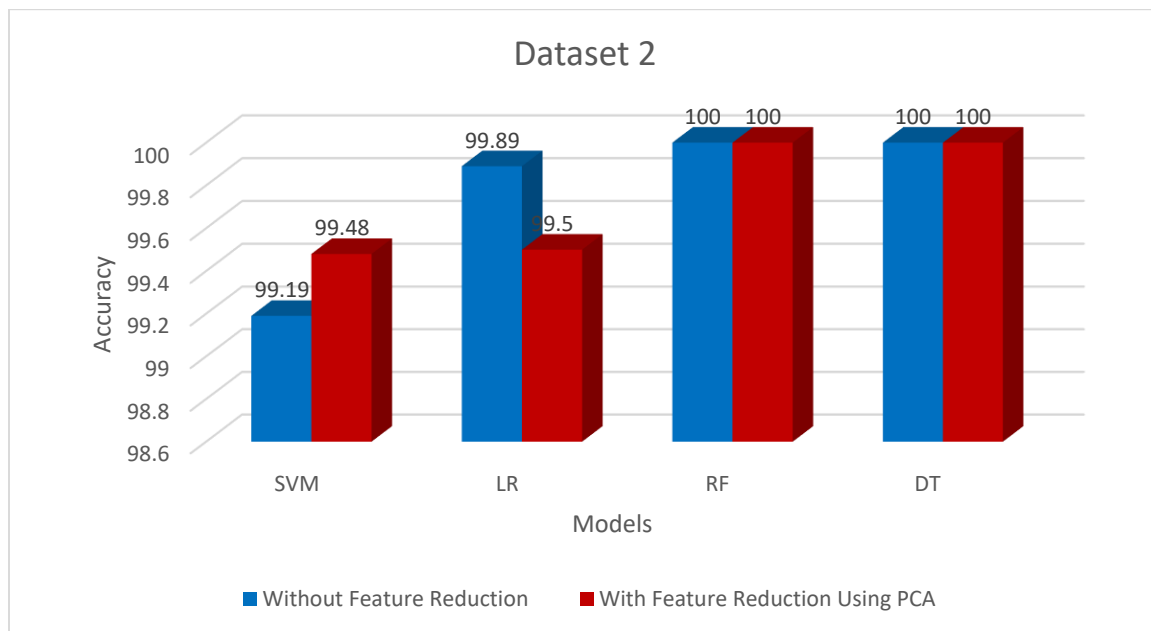


Figure 4.2: Performance of all models for dataset 2.

From the bar chart, the Random Forest (RF) and Decision Tree(DT) again achieve the highest accuracy for Dataset-3 and Dataset-4.

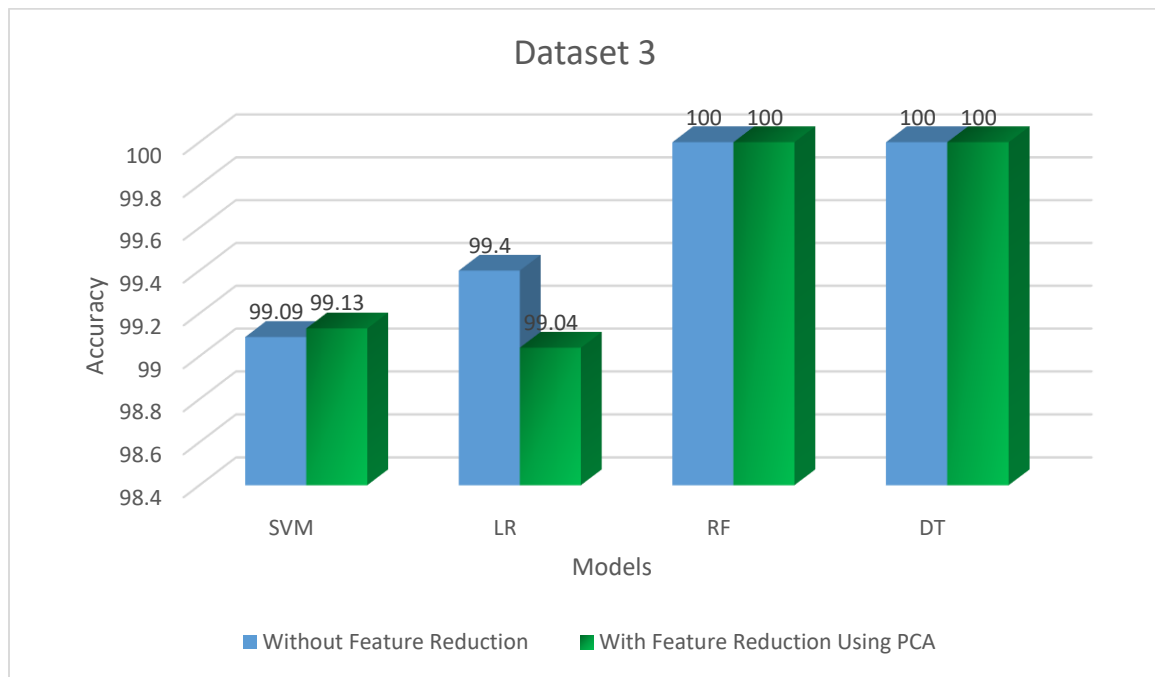


Figure 4.3: Performance of all models for dataset 3.

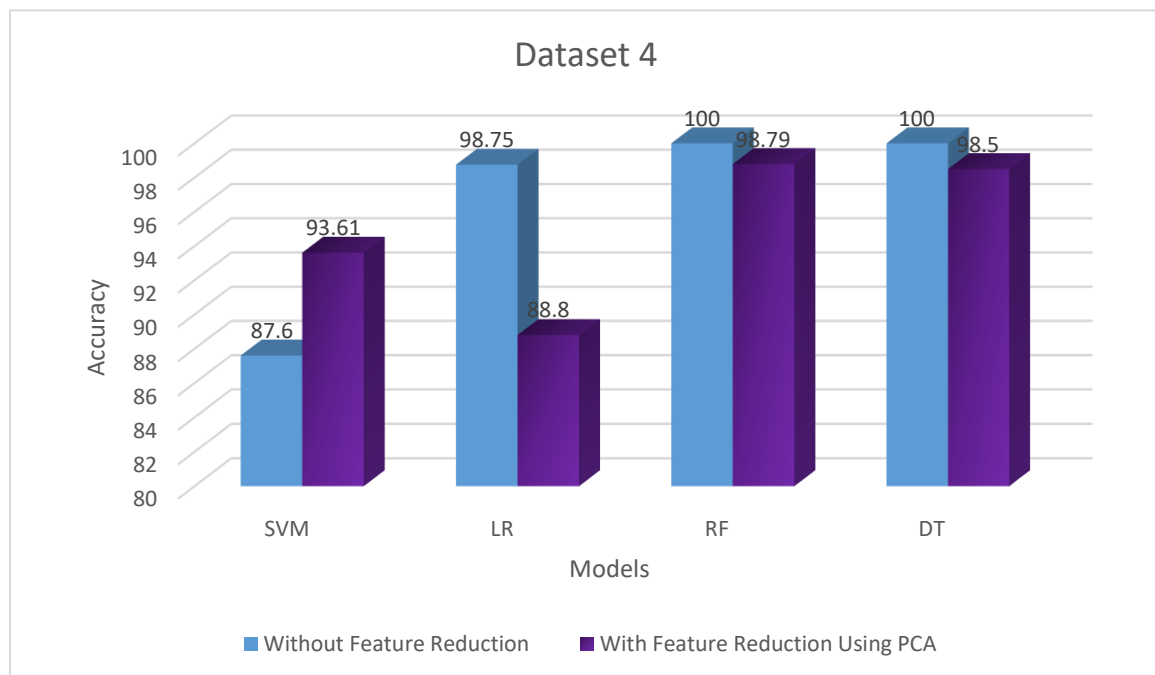


Figure 4.4: Performance of all models for dataset 4.

4.6 ROC Curve:

How effectively a classification model performs across all levels of categorization is shown on a graph called the receiver operating characteristic curve (ROC curve). This graph shows the two parameters:

- True Positive Rate(TPR)
- False Positive Rate(FPR)

True Positive Rate (TPR) is a synonym for recall and is therefore defined as follows

$$TPR = \frac{TP}{(TP+FN)}$$

False Positive Rate (FPR) is defined as follows

$$FPR = \frac{FP}{(FP+TN)}$$

4.6.1 ROC Curve for Dataset-1:

The receiver operating characteristic curve (ROC curve) for Dataset-1 shows how effectively a classification model works at every level of categorization.

The Random Forest (RF) and Decision Tree (DT) models perform better than the Support Vector Machine (SVM) and Logistic Regression models, according to the ROC curve.

4.6.1.1 ROC Curve for Dataset-1 (Without Feature Reduction):

The Random Forest (RF) and Decision Tree (DT) models shows best performance rather than the Support Vector Machine (SVM) and Logistic Regression models for Dataset-1 (Without Feature Reduction).

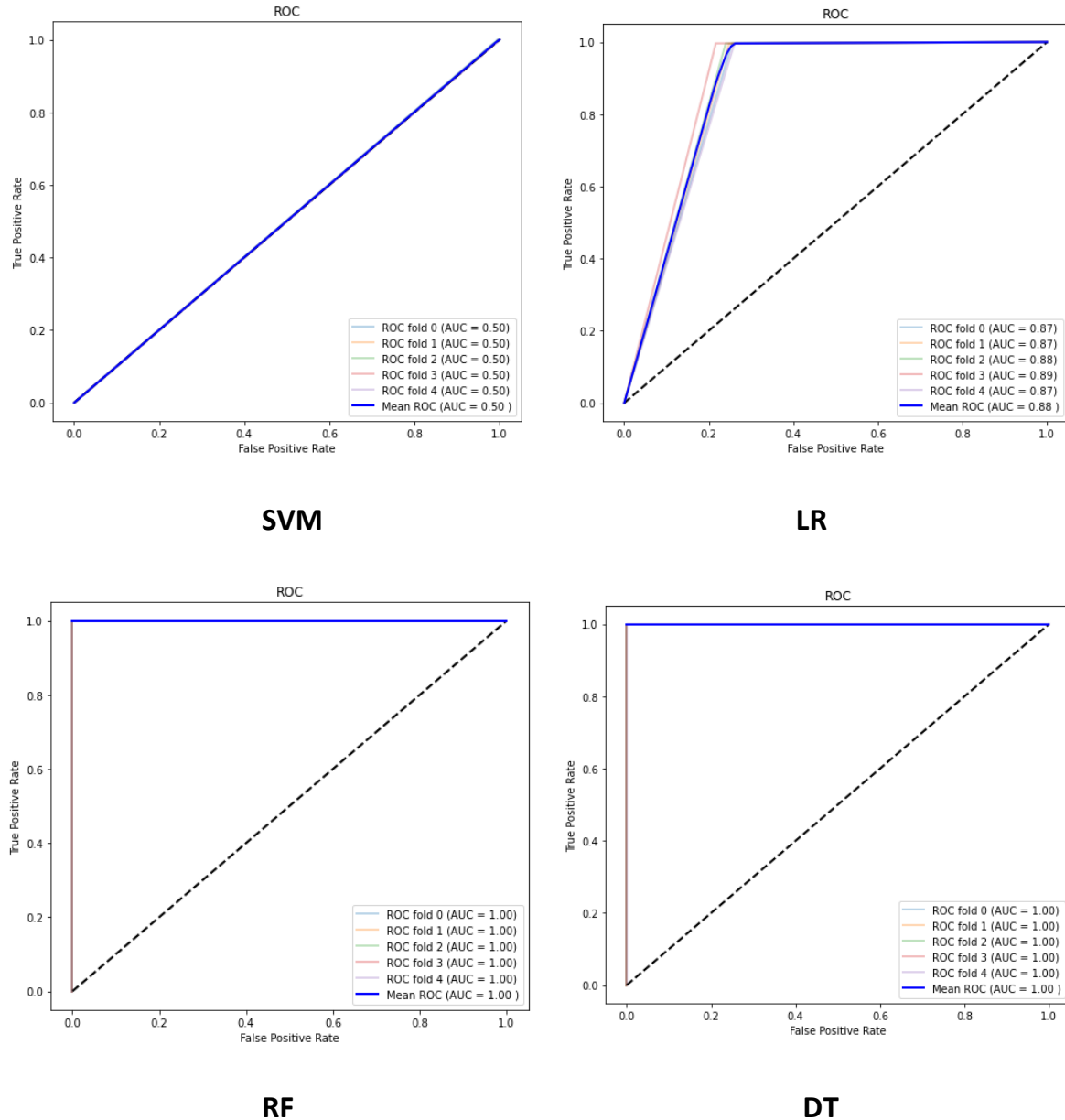


Figure 4.5: ROC Curve of all models for Dataset-1(without Feature Reduction)

4.6.1.2 ROC Curve for Dataset-1 (With Feature Reduction):

The Random Forest (RF) and Decision Tree (DT) models show best performance than the Support Vector Machine (SVM) and Logistic Regression models for Dataset-1 (With Feature Reduction).

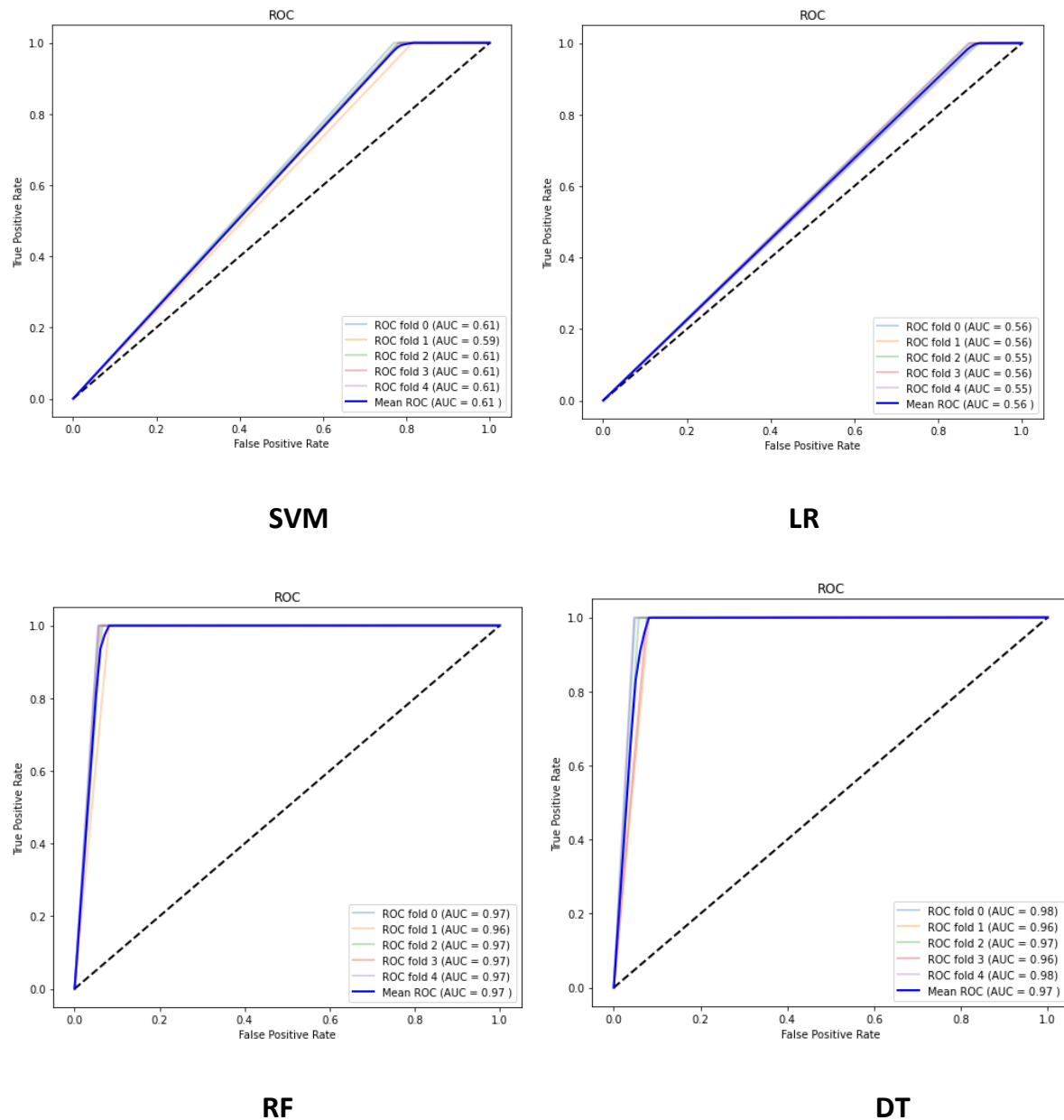


Figure 4.6: ROC Curve of all models for Dataset-1(with Feature Reduction)

4.6.2 ROC Curve for Dataset-2:

The ROC Curve indicates that, for Dataset 2, the Random Forest (RF) and Decision Tree (DT) models perform better, as opposed to the Support Vector Machine (SVM) and Logistic Regression (RF).

4.6.2.1 ROC Curve for Dataset-2 (Without Feature Reduction):

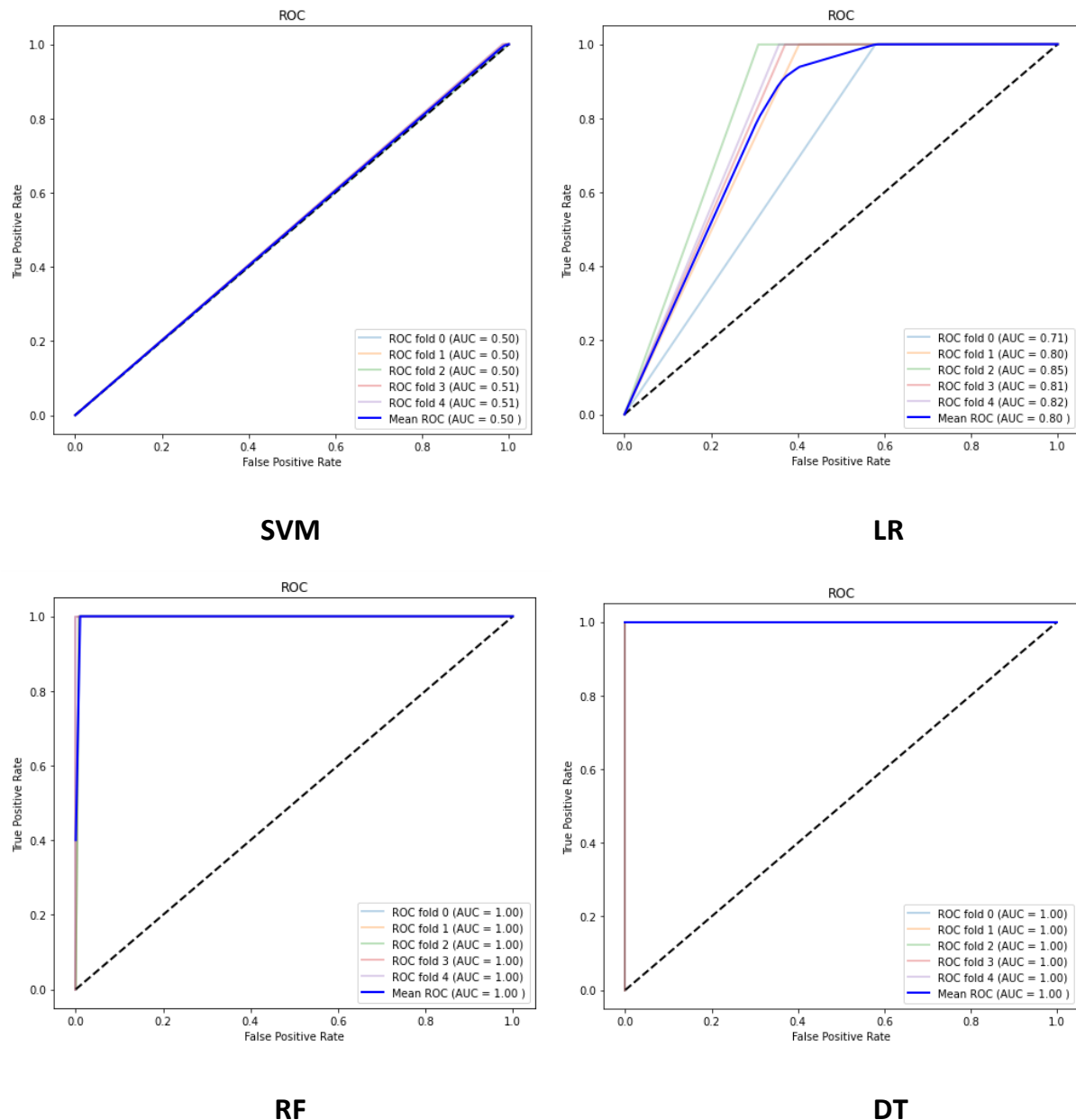


Figure 4.7: ROC Curve of all models for Dataset-2(without Feature Reduction)

4.6.2.2 ROC Curve for Dataset-2 (With Feature Reduction):

For Dataset-2 (With Feature Reduction), The Random Forest (RF) and Decision Tree (DT) models show best performance.

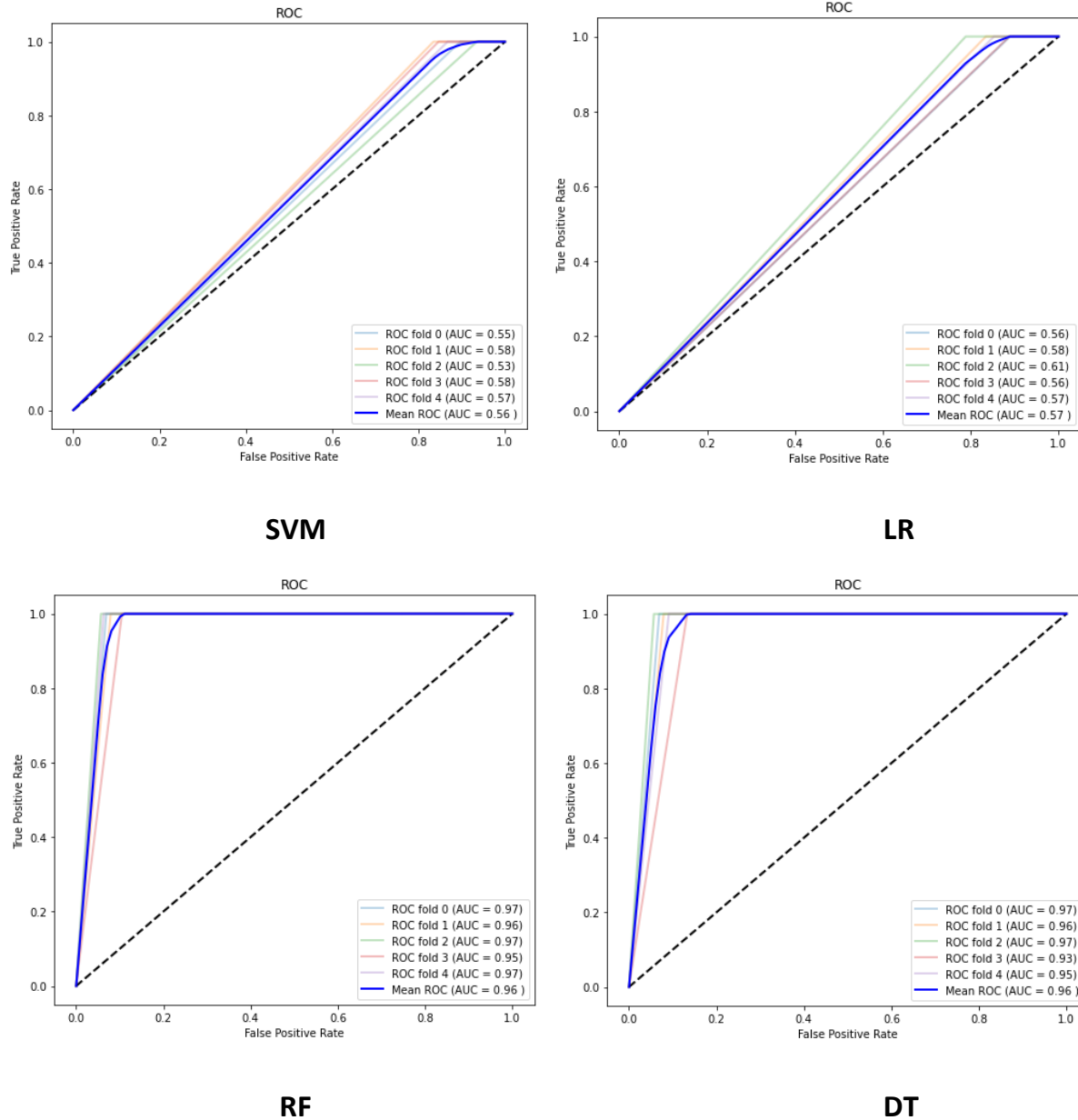


Figure 4.8: ROC Curve of all models for Dataset-2(with Feature Reduction)

4.6.3 ROC Curve for Dataset-3:

The Random Forest (RF) and Decision Tree (DT) models exhibit the best performance for Dataset 3 according to the ROC Curve below rather than the other two models.

4.6.3.1 ROC Curve for Dataset-3 (Without Feature Reduction):

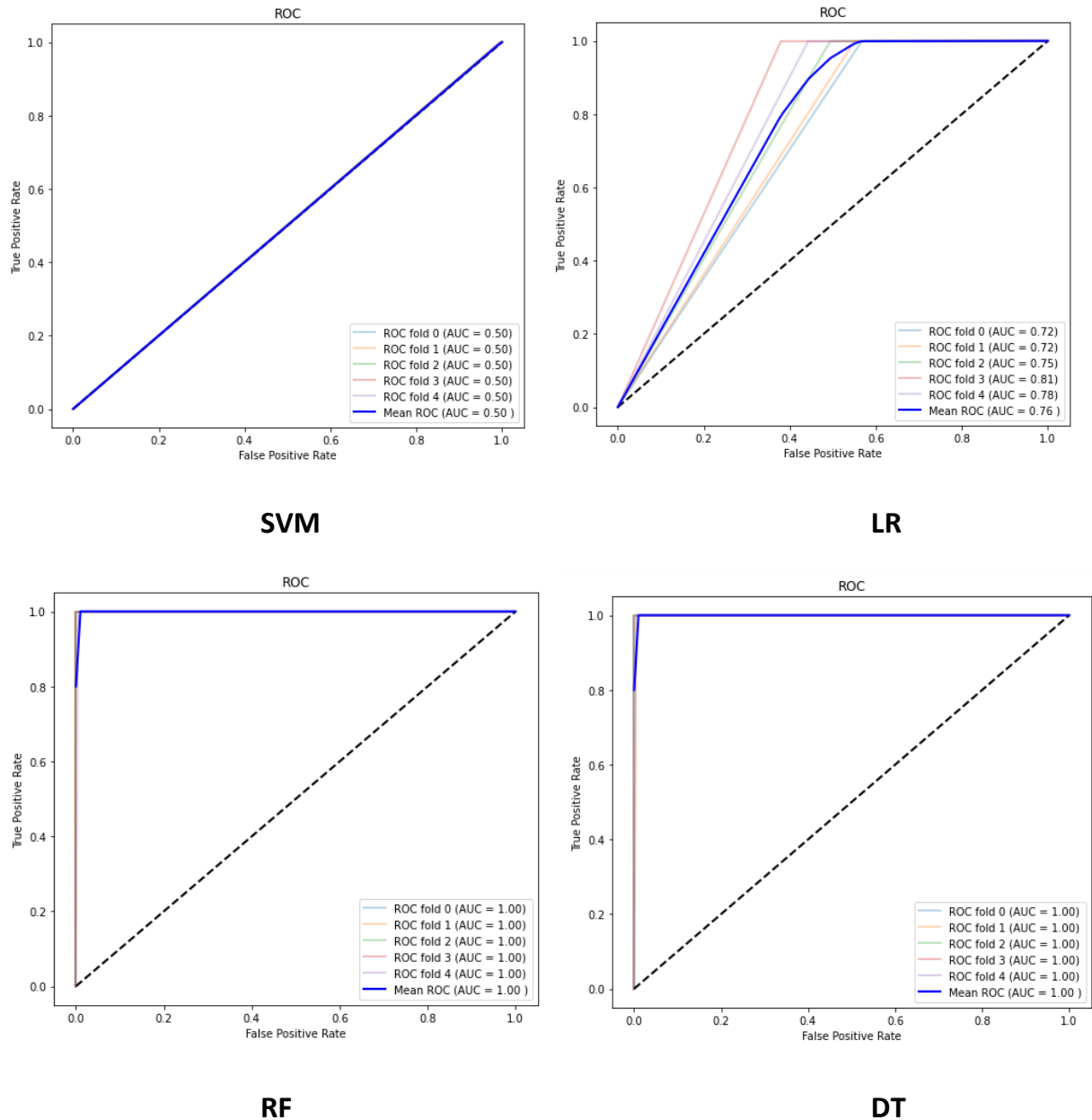


Figure 4.9: ROC Curve of all models for Dataset-3(without Feature Reduction)

4.6.3.2 ROC Curve for Dataset-3 (With Feature Reduction):

For Dataset-3 (Without Feature Reduction), The Random Forest (RF) and Decision Tree (DT) models also show best performance.

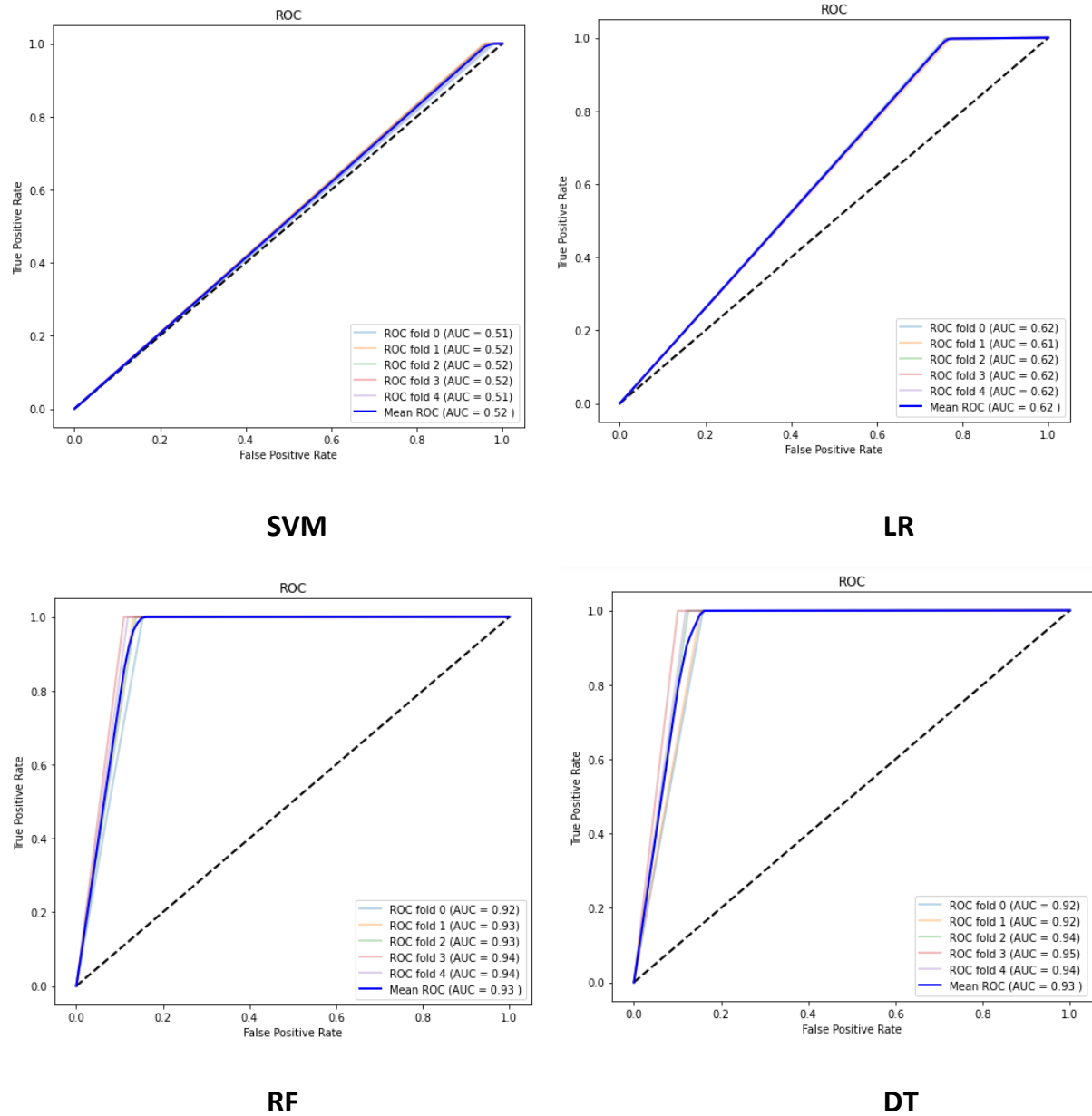


Figure 4.10: ROC Curve of all models for Dataset-3(with Feature Reduction)

4.6.2 ROC Curve for Dataset-4:

Last but not least, the ROC Curve for Dataset-4 reveals that the Random Forest (RF) and Decision Tree (DT) models perform better than the other two models.

4.6.4.1 ROC Curve for Dataset-4 (Without Feature Reduction):

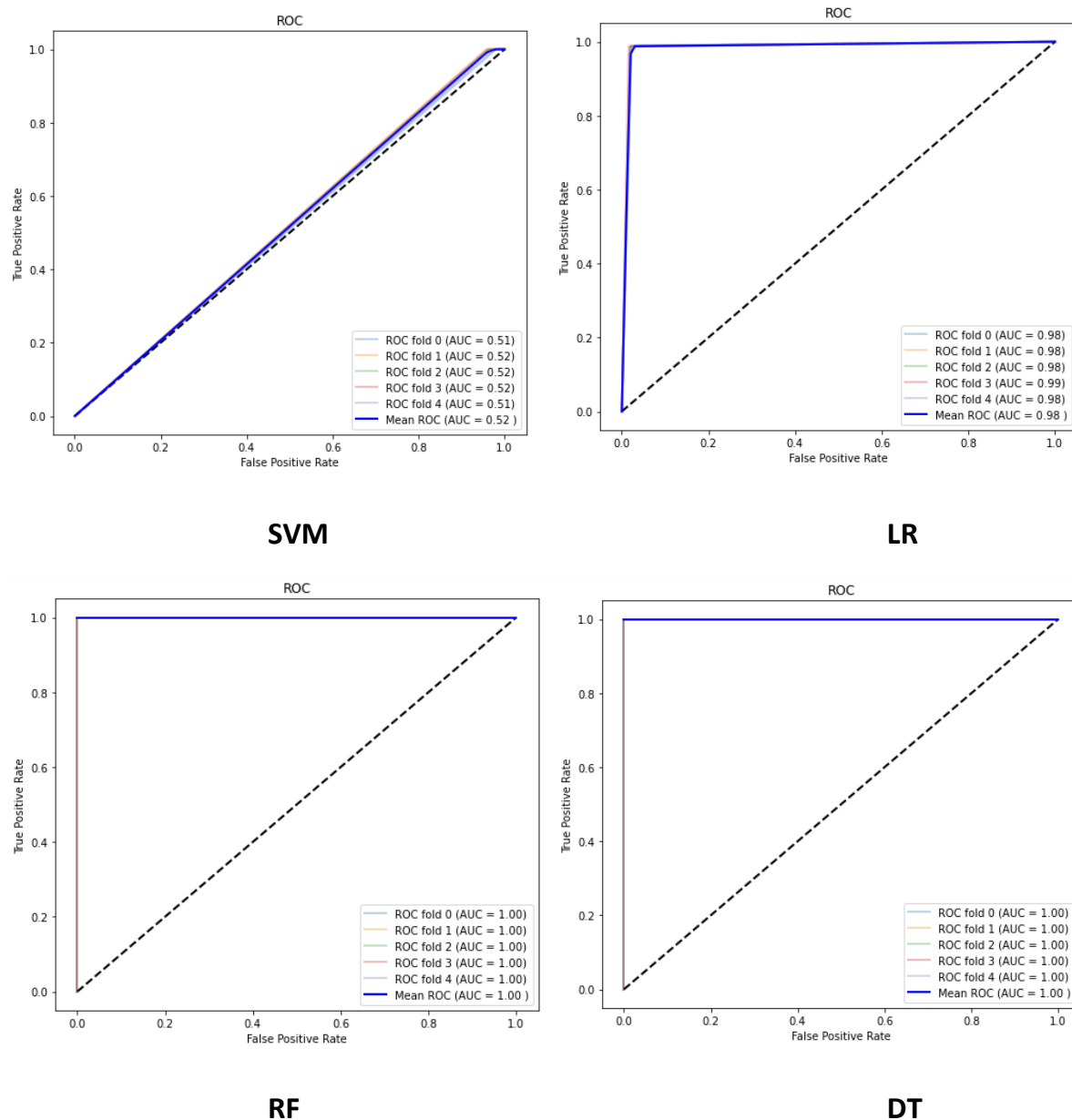


Figure 4.11: ROC Curve of all models for Dataset-4(without Feature Reduction)

4.6.4.2 ROC Curve for Dataset-4 (With Feature Reduction):

Finally, for Dataset-4 (Without Feature Reduction), The Random Forest (RF) and Decision Tree (DT) models show best performance.

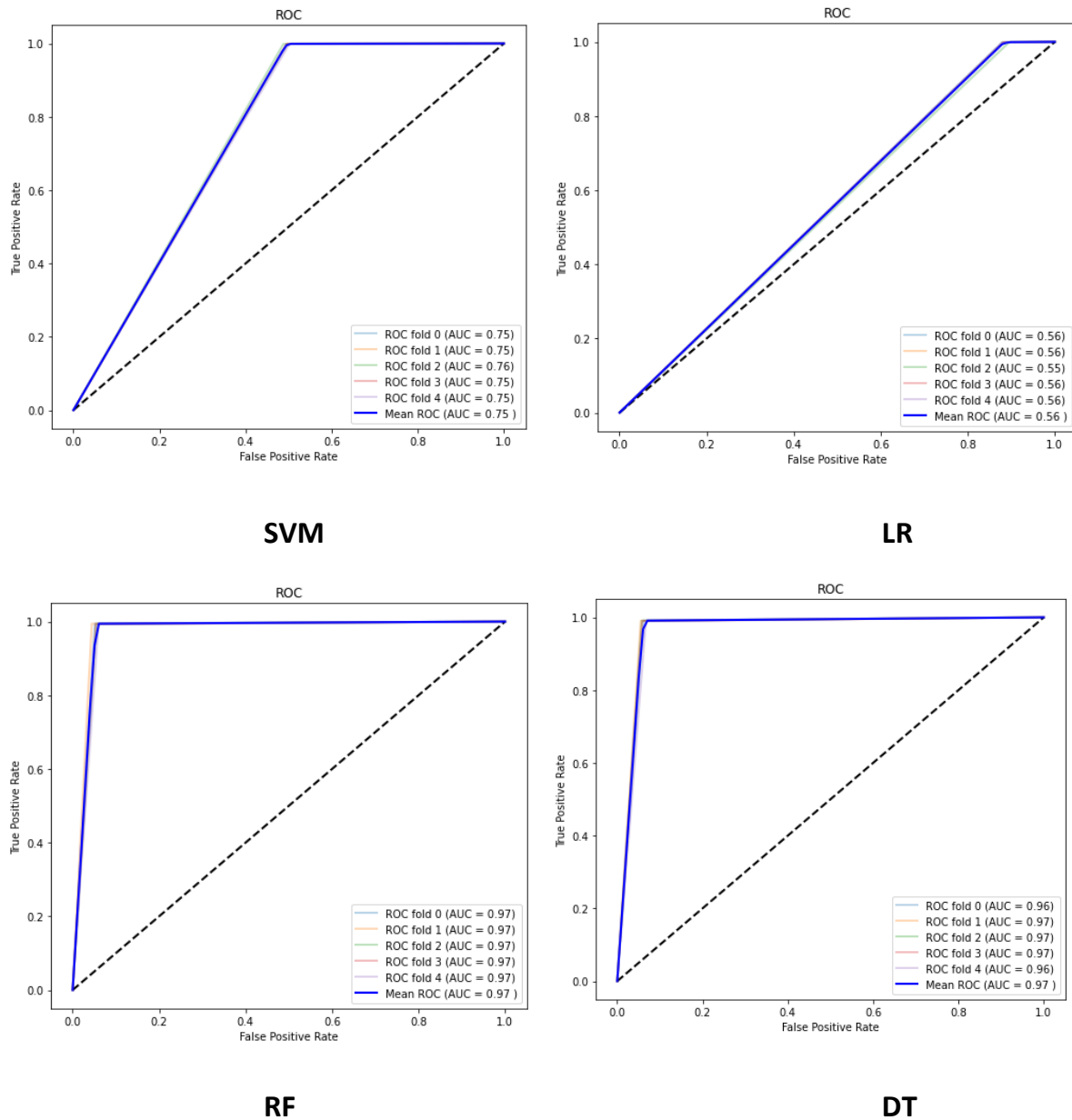


Figure 4.12: ROC Curve of all models for Dataset-4(with Feature Reduction)

4.7 Compare proposed model with existing models

Table 4.13: Compare propose model with previous model.

| Ref. | Model / Dataset | Algorithm / other method used | Accuracy / Rates. |
|------|---|--|--|
| 1 | To detect DRDoS assaults, a novel proactive feature selection model based on better optimization techniques has been developed. | nature-inspired optimization algorithm, (KNN), (RF), and (SVM). | 89.59% accuracy. |
| 4 | In the Big Data Environment, A Deep Forest-Based DRDoS Detection and Defense Method | Deep Forest, IoT, Big Data and other techniques. | A higher rate of detection and false alarms. |
| 7 | An Evaluation of a Machine Learning Approach for Detecting DRDoS Attacks | Support Vector Machine algorithm (SVM).. | The overall success rate, as well as the detection and false positive rates. |
| 8 | For DRDoS attacks, the Protocol Independent Detection and Classification (PIDC) System is used. | Data mining concept and machine learning algorithms like C4.5 classification algorithm. | True positive rates of 99 percent and false positive rates of less than 1% (1 percent). |
| 10 | Based on the UDP Protocol, detect the Reflection Amplification Attack. | detection algorithm | -- |
| 11 | DRDOS Attack Detection and Prevention Methods Using CARD in the MANET (Continuous and Random Dropping). | CARD(Continuous and Random Dropping) technology. | -- |
| 12 | DRDoS Attack Classification Methods Using an Integrated Approach of E-RED and ANT. | Enhanced-Random Application-based Network Traffic (ANT) classification method Early Detection (E-RED) algorithm. | Detects 99 percent of true positives and 1% of false positives, and classifies attacks with a classification |

| | | | | | |
|-----------------------|--|------------------|----------------------------------|-----------------|-------------------------|
| | | | | | accuracy of 98 percent. |
| Proposed Model | Feature Optimization Technique and Machine Learning based DRDoS attack detection model and it's Performance Measurement | Dataset-1 | With Feature Reduction | SVM , LR | 98% , 99% |
| | | | | RF , DT | 100% ,100% |
| | | | Without Feature Reduction | SVM , LR | 99% , 98% |
| | | | | RF , DT | 100% ,100% |
| | | Dataset-2 | With Feature Reduction | SVM , LR | 99.19% , 99.48 % |
| | | | | RF , DT | 100% ,100% |
| | | | Without Feature Reduction | SVM , LR | 99.89% , 99.50% |
| | | | | RF , DT | 100% ,100% |
| | | Dataset-3 | With Feature Reduction | SVM , LR | 99.09% , 99.40% |
| | | | | RF , DT | 100% ,100% |
| | | | Without Feature Reduction | SVM , LR | 99.13% , 99.04% |
| | | | | RF , DT | 100% ,100% |
| | | Dataset-4 | With Feature Reduction | SVM , LR | 87.60% , 98.75% |
| | | | | RF , DT | 100% ,100% |
| | | | Without Feature Reduction | SVM , LR | 99.61% , 88.80% |
| | | | | RF , DT | 98.79% , 98.50% |

Chapter 5

Conclusion and Future Scope

5.1 Conclusion

This paper has presented four machine learning algorithms for detecting the DRDoS attack. The four algorithms like DT, SVM, RF and LR are utilized in this research and also use Principal Component Analysis (PCA) technique for Feature Reduction. The authors have applied without Feature Reduction and with Feature Reduction technique using PCA for the best result and also for time optimization. Accuracy, precision, f1-score, and recall are some of the evaluation criteria that are used to assess each model's performance. Experiment results show that both the RF and DT algorithms achieved the best results compared to SVM and LR algorithms with a detection accuracy of 100% and the F1 score is also the same for both without Feature Reduction and with Feature Reduction technique using PCA (except Dataset-4) method. The authors of the paper only focus on detection DRDoS attack, but they do not give any suggestion for prevent this attack. Future research will examine more machine learning algorithms for DRDoS attack detection utilizing the same datasets with the abundant data and suggested prevention approach.

5.2 Limitation

Major limitations of this work are summarized below:

- The dataset which is used in this research is little bit imbalanced.
- The datasets have huge number of data. The authors cannot work for all the huge data of this datasets.

5.3 Future scope

In future, researchers try to balance the dataset. After that they try to apply machine learning models on their balance dataset. It might be show more accurate result.

References

1. R. R. Nuijaa, S. Manickam, A. H. Alsaeedi, E. S. Alomari, "A new proactive feature selection model based on the enhanced optimization algorithms to detect DRDoS attacks", in International Journal of Electrical and Computer Engineering (IJECE), Vol. 12, No. 2, p. 1869, 2022.
2. R. R. Nuijaa, S. Manickam, A. H. Alsaeedi, E. S. Alomari, "Distributed reflection denial of service attack: A critical review". International Journal of Electrical and Computer Engineering (IJECE), Vol. 11, No. 6, pp. 5327~5341, December 2021.
3. X. Chen, W. Feng, Y. Ma, N. Ge, X. Wang, "Preventing DRDoS Attacks in 5G Networks: a New Source IP Address Validation Approach", Conference: GLOBECOM 2020 - 2020 IEEE Global Communications Conference, December 2020, DOI: 10.1109/GLOBECOM42002.2020.9322314.
4. R. Xu 1, J. Cheng , F. Wang, X. Tang and J. Xu, "A DRDoS Detection and Defense Method Based on Deep Forest in the Big Data Environment", *Symmetry* 11(1):78,2019,Doi:<https://doi.org/10.3390/sym11010078>.
5. H. Fujinoki. "Cloud-Base Defense against DRDoS Attacks", IEEE International Conference on Consumer Electronics-Taiwan (ICCE-TW), 2018, DOI: [10.1109/ICCE-China.2018.8448533](https://doi.org/10.1109/ICCE-China.2018.8448533).
6. M. Aupetit, Y. Zhauniarovich, G Vasiliadis, M. Dacier, Y. Boshmaf, "Visualization of Actionable Knowledge to Mitigate DRDoS Attacks", In IEEE SYMPOSIUM ON VISUALIZATION FOR CYBER SECURITY (VIZSEC), 2016.
7. Y. Gao, Y. Feng, J. Kawamoto, K. Sakurai, "A Machine Learning Based Approach for Detecting DRDoS Attacks and Its Performance Evaluation", in IEEE 2016 DOI:10.1109/AsiaJCIS.2016.24.
8. P. M. Priya, V.Akilandeswari, Dr.S. M. Shalinie, V.Lavanya, M.S. Priya, "The Protocol Independent Detection and Classification (PIDC) System for DRDoS Attack", in IEEE International Conference on Recent Trends in Information Technology, 2014, DOI: [10.1109/ICRTIT.2014.6996154](https://doi.org/10.1109/ICRTIT.2014.6996154).
9. H. Tsunoda, K. Ohta, A. Yamamoto, N. Ansari, Y. Waizumi, Y. Nemoto, "Detecting DRDoS attacks by a simple response packet confirmation mechanism" . *Computer Communications* 31(14):3299-3306, 2008.
10. C. Liu, G. Xiong, J. Liu, G. Gou, "Detect the Reflection Amplification Attack Based on UDP Protocol", in IEEE 10th International Conference on Communications and

- Networking in China (China Com), 2015, DOI: 10.1109/CHINACOM.2015.7497948
11. R. Rani, A. K. Vatsa, “CARD (Continuous and Random Dropping) based DRDOS Attack Detection and Prevention Techniques in MANET”, International Journal of Engineering and Technology (IJET) – Volume 2 No. 8, August, 2012.
 12. M. Priya, V. Akilandeswari, G. Akilarasu, and S. M. Shalinie, “An Integrated Approach of E-RED and ANT Classification Methods for DRDoS Attacks”, Communications in Computer and Information Science· September 2014, DOI: 10.1007/978-3-662-44966-0_29.
 13. M. H. H. Khairi, S. H. S. Ariffin, N. M. A. Latiff, A. S. Abdullah, and M. K. Hassan, “A review of anomaly detection techniques and distributed denial of service (DDoS) on software defined network (SDN),” Eng., Technol. Appl. Sci. Res., vol. 8, no. 2, pp. 2724–2730, Apr. 2018.
 14. <https://www.artur.ai/en/support/faqs/drdoS-attacks>.
 15. F. J. Ryba, M. Orlinski, M. Wahlisch, C. Rossow, T. C. Schmidt, “Amplification and DRDoS Attack Defense – A Survey and New Perspectives”, arXiv, 2015. [arXiv:1505.07892v3](https://arxiv.org/abs/1505.07892v3). <https://doi.org/10.48550/arXiv.1505.07892>
 16. Y. A. Bekeneva, A. V. Shorov, “Simulation of DRDoS-attacks and Protection Systems against them”, 2017 XX IEEE International Conference on Soft Computing and Measurements (SCM), DOI: [10.1109/SCM.2017.7970527](https://doi.org/10.1109/SCM.2017.7970527)
 17. Zhijun, W.; Qing, X.; Jingjie, W.; Meng, Y.; Liang, L. Low-Rate DDoS Attack Detection Based on Factorization Machine in Software Defined Network. IEEE Access 2020, 8, 17404–17418.
 18. Y. Bekeneva, N. Shipilov, and A. Shorov, “Investigation of Protection Mechanisms Against DRDoS Attacks Using a Simulation Approach”, 2016 International Conference on Next Generation Wired/Wireless Networking Conference on Internet of Things and Smart Spaces, DOI: [10.1007/978-3-319-46301-8_26](https://doi.org/10.1007/978-3-319-46301-8_26).
 19. A. D. Olaniyi, R. Christoph, S. A. Simon, A. A. Taofeek, and B. S. Badmus, “Resolving DRDoS Attack in Cloud Database Service Using Common Source IP and Incremental Replacement Strategy”, 2018 Proceedings of SAI Intelligent Systems Conference, DOI: [10.1007/978-3-319-56991-8_52](https://doi.org/10.1007/978-3-319-56991-8_52)
 20. I. M. TAS, B. G. UNSALVER, AND S. BAKTIR, “A Novel SIP Based Distributed Reflection Denial-of-Service Attack and an Effective Defense Mechanism”, June 2020 [IEEE Access](https://doi.org/10.1109/Access.2020.2999999) PP(99).

21. H. TSUNODA and Y. NEMOTO, "A Simple Response Packet Confirmation Method for DRDOS Detection", Source [IEEE Xplore](#), Advanced Communication Technology, 2006. ICACT 2006. The 8th International Conference Volume: 3, DOI: [10.1109/ICACT.2006.206282](#), 2006.
22. K. Subramani, R. Perdisci, M. Konte, "IXmon: Detecting and Analyzing DRDoS Attacks at Internet Exchange Points", Published 22 June 2020. Computer Science. ArXiv.
23. L. Berti-Equille, Y. Zhauniarovich, "Profiling DRDoS Attacks with Data Analytics Pipeline", Conference: the 2017 ACM, DOI: [10.1145/3132847.3133155](#).
24. A. A. NASSER A, W. Chen, "NTP DRDoS Attack Vulnerability and Mitigation", September 2014 [Applied Mechanics and Materials](#) 644-650:2875-2880, DOI: [10.4028/www.scientific.net/AMM.644-650.2875](#).
25. B. A. Sassani (Sarrafpour), C. Abarro, I. Pitton, C. Young and F. Mehdipour, "Analysis of NTP DRDoS Attacks' Performance Effects and Mitigation Techniques", 2016 14th Annual Conference on Privacy, Security and Trust (PST), DOI: [10.1109/PST.2016.7906966](#)
26. S. MahdaviFar, N. Maleki, A. H. Lashkari, M. Broda, Amir H. Razavi, "Classifying Malicious Domains using DNS Traffic Analysis", The 19th IEEE International Conference on Dependable, Autonomic, and Secure Computing (DASC), Oct. 25-28, 2021, Calgary, Canada.
27. <http://205.174.165.80/CICDataset/CICBellDNS2021/>.