

Network Attack Analysis Using Machine Learning

Abstract

Cyber threats result in compromising the confidential information of a particular person or the organization. Hackers may leak necessary information to the unknown which may be hazardous. Cyber-attacks are not only limited to bigger organizations but also individuals. So, the hacker can target anyone for the information. The tendency of large-scale deployment of essential corporate applications in cloud centers has risen dramatically as cloud computing has progressed.

To solve this problem a new concept called Threat Intelligence was introduced by the research developers. Due to its capacity to evaluate threat intelligence data from various online assaults, threat intelligence will become an essential tool for resolving sophisticated network attacks, real-time threat time indications, and assault monitoring. This project's major objective is to employ various machine learning techniques to examine the threat intelligence datasets in predicting network attacks coming either from a corporate network or from a specific cluster.

Introduction

- The network has become the main tool or weapon equipment for some malicious organizations to do evil, and the types of network threats are more complex and changeable.
- Threat Intelligence is a subset of intelligence focused on information security.

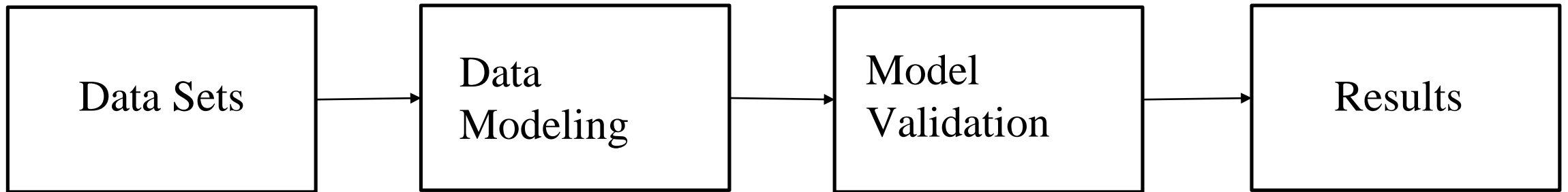


Figure 1: Overview of the proposed model

Cont..

- Security intelligence has become a lightweight important means to solve complex network attacks, threat early warning and security traceability.
- Therefore, it is very important to reasonably use threat intelligence, share security intelligence information in real time, and establish a set of border real-time warning and defense system driven by security intelligence.
- The main objective is to not to identify which algorithm is best, but to review the main dataset publicly available used to train and test security solutions that employ classification/regression algorithm. Hence KDD datasets were taken into considerations.
- Training on infrastructure specific datasets is more efficient from the security point of view than training on old attack signatures.

Literature Survey

Sl. No	TITLE	AUTHOR	YEAR OF PUBLICATION	LIMITATIONS
1	Cyber Threat Intelligence Model: An Evaluation of Taxonomies, Sharing Standards, and Ontologies within Cyber Threat Intelligence	Vasileios Mavroeidis, Siri Bromander	Feb 2021	Hurdles to achieve include a lack of dedicated ontological cyber threat intelligence efforts that can account for the strategic, operational, and tactical levels; ambiguity in defined concepts that prevents ontology integration and adoption
2	Detecting Vulnerabilities in Critical Infrastructures by Classifying Exposed Industrial Control Systems Using Deep Learning	Pablo Blanco-Medina , Eduardo Fidalgo , Enrique Alegre , Roberto A. Vasco-Carofilis , Francisco Jañez-Martino and Victor Fidalgo Villar	2018	Further improving the proposed solution, increasing our training images, and adding new layers on top of the given architectures based on VGG16 results.

Literature Survey

Sl. No	TITLE	AUTHOR	YEAR OF PUBLICATION	LIMITATIONS
3	A Cloud Computing Based Network Monitoring and Threat Detection System for Critical Infrastructures	ZhijiangChen, GuobinXu, VivekMahalingam, LinqiangGe, JamesNguyen, WeiYu*, ChaoLu	2 November 2015	The researchers did not put the computational cost into consideration in the choice of which machine learning technique to use for spam mail filtering. Their main focus is performance in terms of classification accuracy.
4	Analysis of adversary activities using cloud-based web services to enhance cyber threat intelligence	Hamad Al-Mohannadi1 · Irfan Awan1 · Jassim Al Hamar	2020	Attacker tried to get access to root but failed. In a UNIX-based system, getting access to the root gives attackers full control to the system. Since the authentication was not confirmed, the attacker could not get the access.

Cont..

- Disadvantages
 - Lack of precision in attack detection.
 - It is difficult to control the access of network in real time.
 - Less security and privacy.
- Advantage of Proposed system
 - Predicts the precision about insider and outsider attacks.
 - Improves reliable data access control and can handle huge number of data.

Problem Identification

- With the development of the Internet and the application of information technology, the general trend is opener. Important data and important business applications will face more serious security threats.
- The network has become the main tool or weapon equipment for some malicious organizations to do evil, and the types of network threats are more complex and changeable.
- The existing security solutions are difficult to cope with the high level, continuous, group and weaponization threats.
- However, important industries and key information infrastructure security network units mainly support apt detection and defense equipment with high application threshold and general landing effect.

Objectives

- To analyze the security given from the insider and outsider attacks. Applications employ machine learning techniques to recognize and respond to assaults. Larger data sets of security events can be analyzed to find correlations in harmful activity and help with this. When similar actions are found, ML makes it if the trained ML model can intelligently handle them.
- It is crucial to investigate and use advanced techniques for the detection and reporting of such assaults because there are more and more cyberattacks targeting crucial networked resources that cannot be detected by conventional monitoring tools. Hence the main objective is to improve the accuracy of attack detection.

Proposed System

- As the previous sections explain the existing problems, the proposed system is being developed for analyzing the dataset using a Graphical User Interface. It mainly analyses the insider or outsider attacks to enhance the accuracy of attack detection. It improves reliable data access control and can handle huge amounts of data.
- Threat Intelligence is the process of identifying the threats and analyzing them in other words this feature informs the user about the current threats and the threats which are about to happen in the current scenario of the organization.
- Hence several machine learning algorithms are employed just to check the efficiency of the network attack blocking method. Threat intelligence feeds can help with this approach by detecting real-time indications of the breach and advising preventative measures.
- The results are examined utilizing network attack data sets, demonstrating that it is one of the best competitor data sets for threat intelligence analysis and testing.

Requirement Analysis

- Software Requirements
 - Anaconda
 - Python Programming language
- Software Packages
 - Keras
 - TensorFlow
- Software And Hardware Requirements
 - Operating system used is the Windows 8/10/11.
 - Processor – Ryzen 7@ 3.20 GHz
 - RAM – 8GB and above
 - Hard disk –20 GB is required to store the data set obtained.

Threat Intelligence

1. Threat intelligence, or cyber threat intelligence, is information an organization uses to understand the threats that have, will, or are currently targeting the organization. This info is used to prepare, prevent, and identify cyber threats looking to take advantage of valuable resources.
 - Ensure to stay up to date with the often overwhelming volume of threats, including methods, vulnerabilities, targets and bad actors.
 - Help one to become more proactive about future cybersecurity threats.
 - Keep leaders, stakeholders and users informed about the latest threats and repercussions they could have on the business.

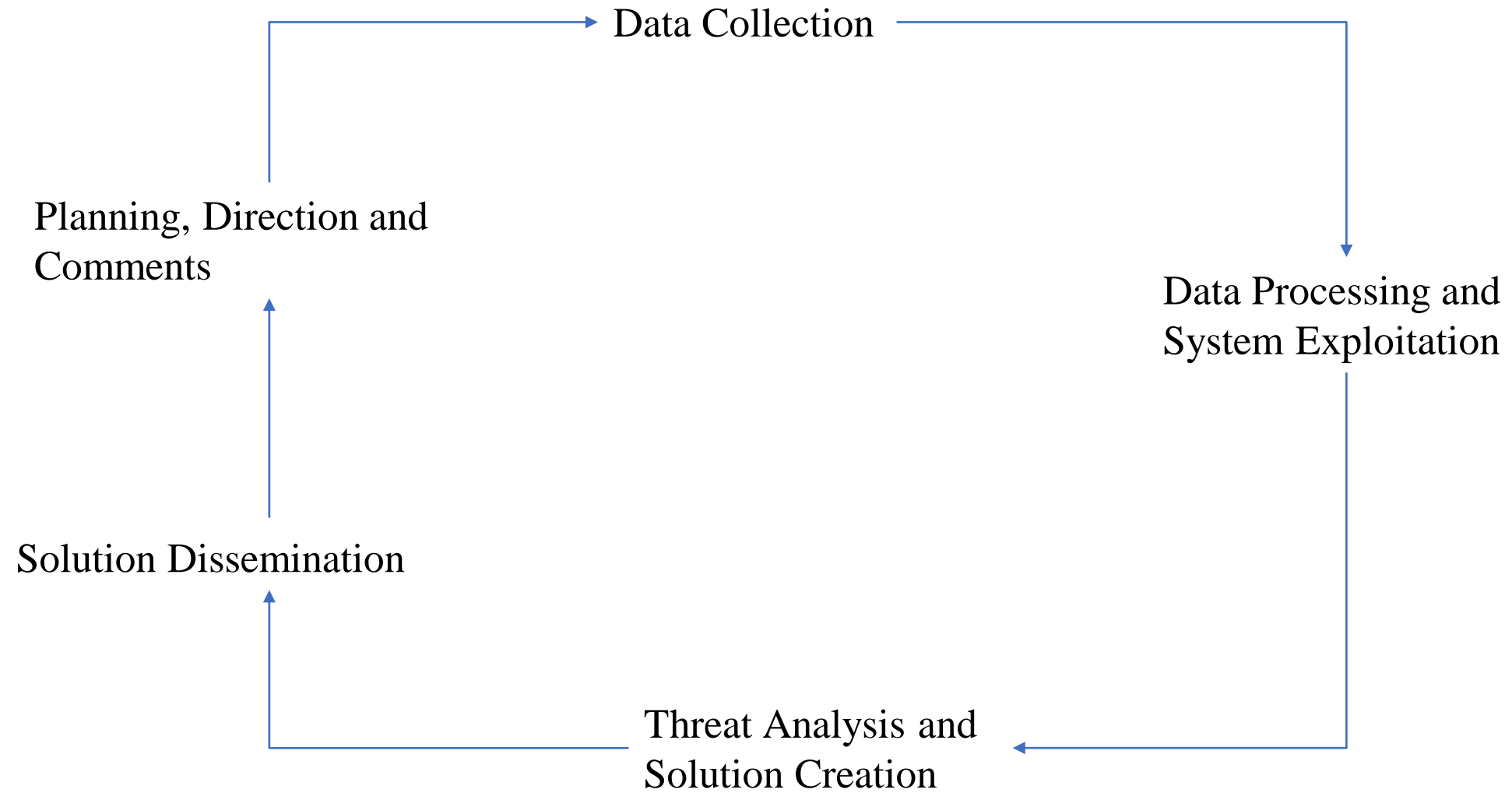


Figure 2: Life cycle of Threat intelligence

System Workflow

- The System Workflow description is briefly explained where the first step in creating a predictive model is to create a dataset. The dataset can be thought of as practically being an MN matrix, where M stands for the columns (features) and N for the rows (samples). One might divide columns into X and Y for example. The method of picking features, whether automatic mode or not, depends on which components are most important to the target variable or desired outcome.
- Algorithms' performance was accurately depicted after being evaluated against the aforementioned criteria. Some experts were just concerned with the algorithm's correctness, some were worried about how well the system performed when tested against every matrix parameter.

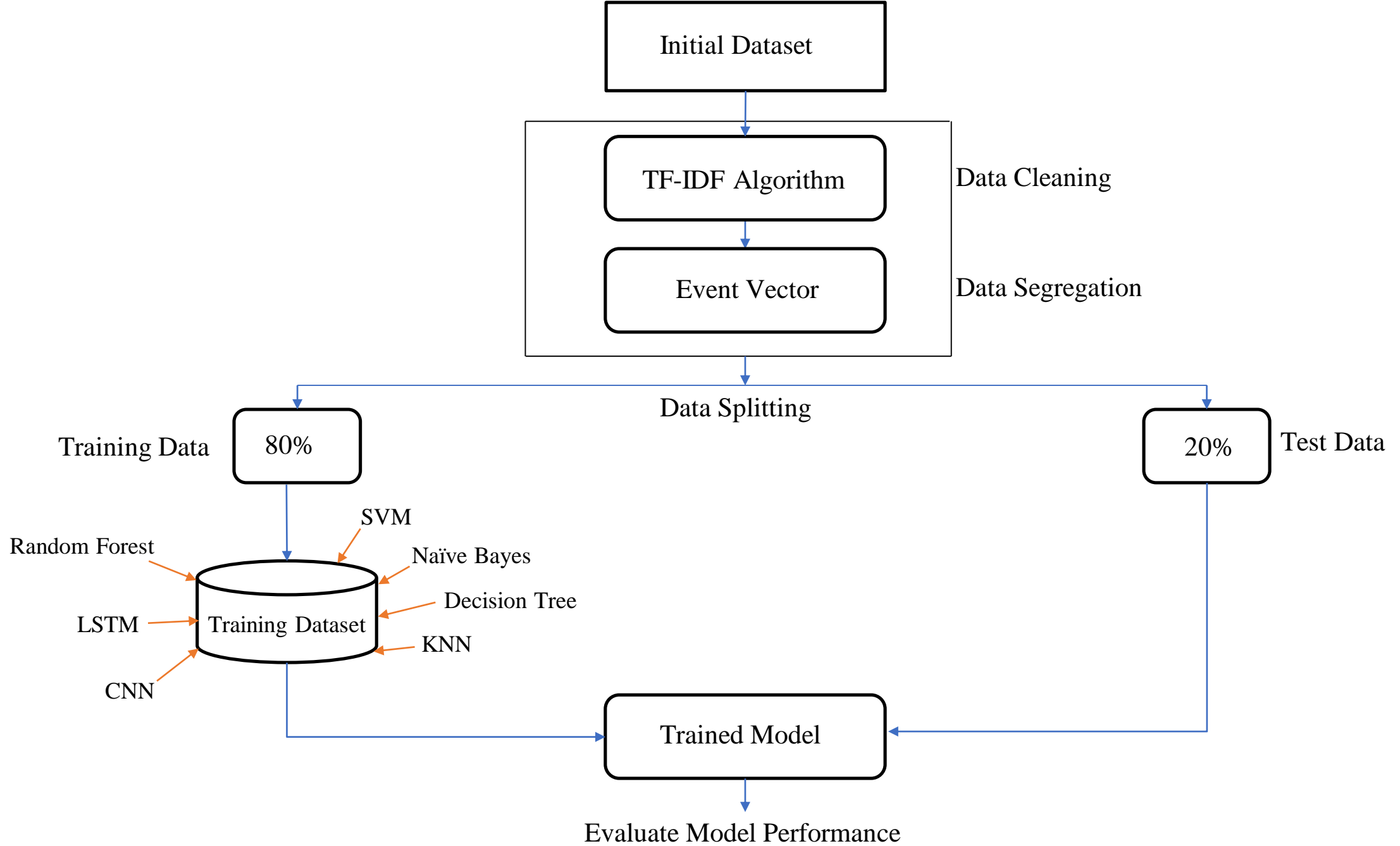


Figure 3: System Workflow of the proposed system

Sequence Diagram

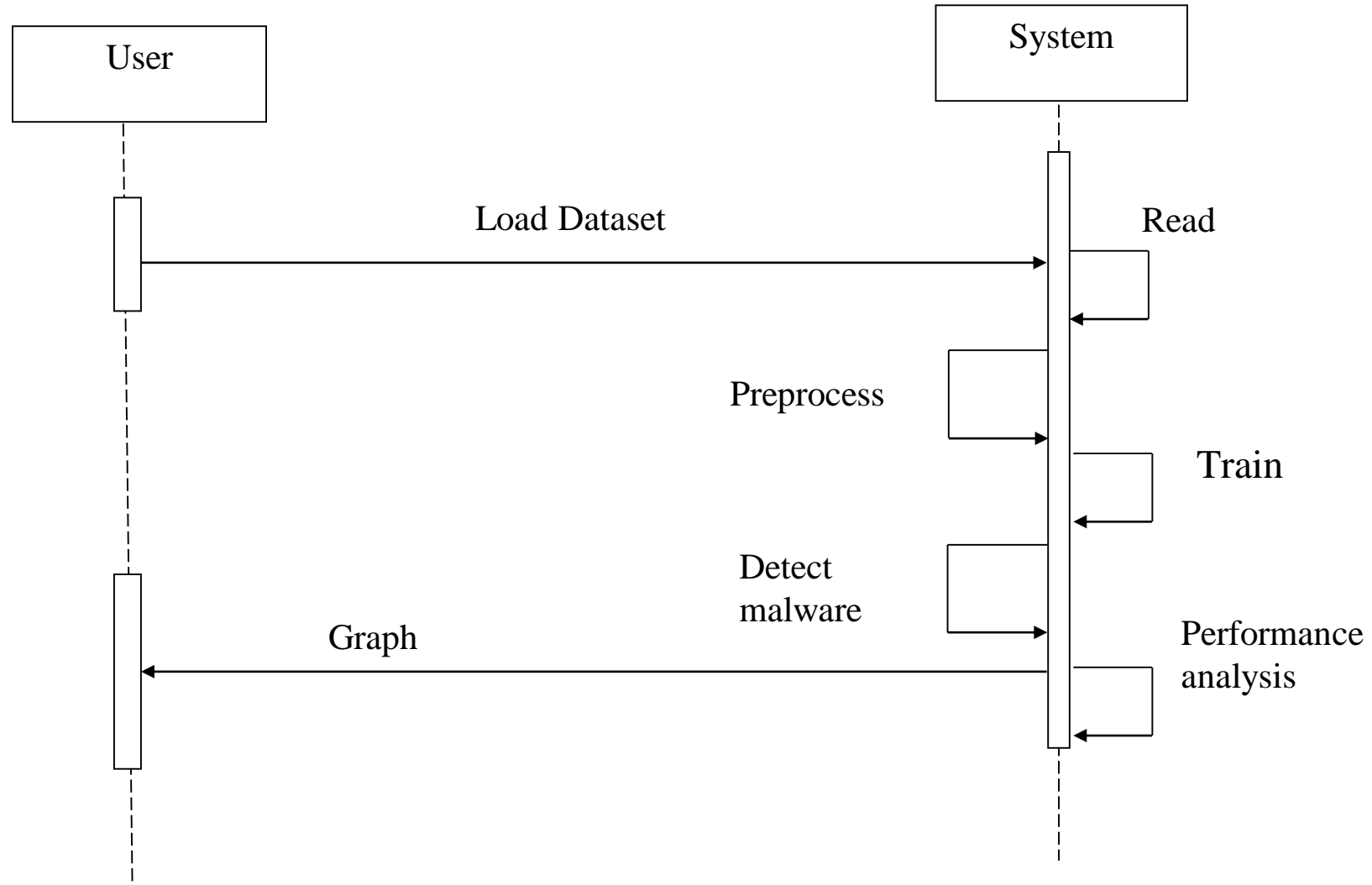


Figure 4: Sequence Diagram

Implementation

Dataset description

- The dataset is collected from the [kaggle.com](https://www.kaggle.com)
- The dataset contains day wise packet transaction detail attributes like Destination Port, Flow Duration, Total Fwd Packets, Total Backward Packets, Total Length of Fwd Packets, Total Length of Bwd Packets, Fwd Packet Length Max, Fwd Packet Length Min, Fwd Packet Length Mean, Fwd Packet Length Std, Bwd Packet Length Max, Bwd Packet Length Min, Bwd Packet Length Mean Bwd, total 76 attributes and 1 label class.

Cont..

Friday-WorkingHours-Afternoon-PortScan - Excel

Abhishek Abhishek

FileHomeInsertPage LayoutFormulasDataReviewViewHelpTell me what you want to do

<

Figure 5: Snapshot of Dataset

Cont..

- **Destination Port:** Destination of the packets in the network.
- **Flow Duration:** Time taken by the packet to reach from source to destination
- **Total Fwd Packets:** Total number of packets that flowed from source to destination
- **Total Backward Packets:** Total number of packets that flowed from destination to source
- **Total Length of Fwd Packets:** Total length of the Forwarded packet.
- **Total Length of Bwd Packets:** Total length of the backward packets
- **Fwd Packet Length Max:** Max length of the forward packet
- **Fwd Packet Length Min:** Min length of the forward packets
- **Fwd Packet Length Mean:** Mean of the forward packets which specifies the mean value
- **Fwd Packet, Length Std, Bwd Packet Length Max, Bwd Packet Length Min, Bwd Packet Length Mean Bwd,**
total 76 attributes and 1 label class.

Preprocessing

- Getting the dataset.
- Importing libraries.
- Importing datasets.
- Finding Missing Data.
- Removing the duplicate data.
- Fixing Structural errors.

Algorithms used

- **Random forests** are ensemble learning methods that are used for either classification or regression purposes. Random forests are composed of several decision trees that are combined together to make a unanimous decision or classification. Random forest are better than just regular decision trees because they do not cause overfitting of the data.

Algorithm 1

Step-1: Select random K data points from the training set.

Step-2: Build the decision trees associated with the selected data points (Subsets).

Step-3: Choose the number N for decision trees that you want to build.

Step-4: Repeat Step 1 & 2.

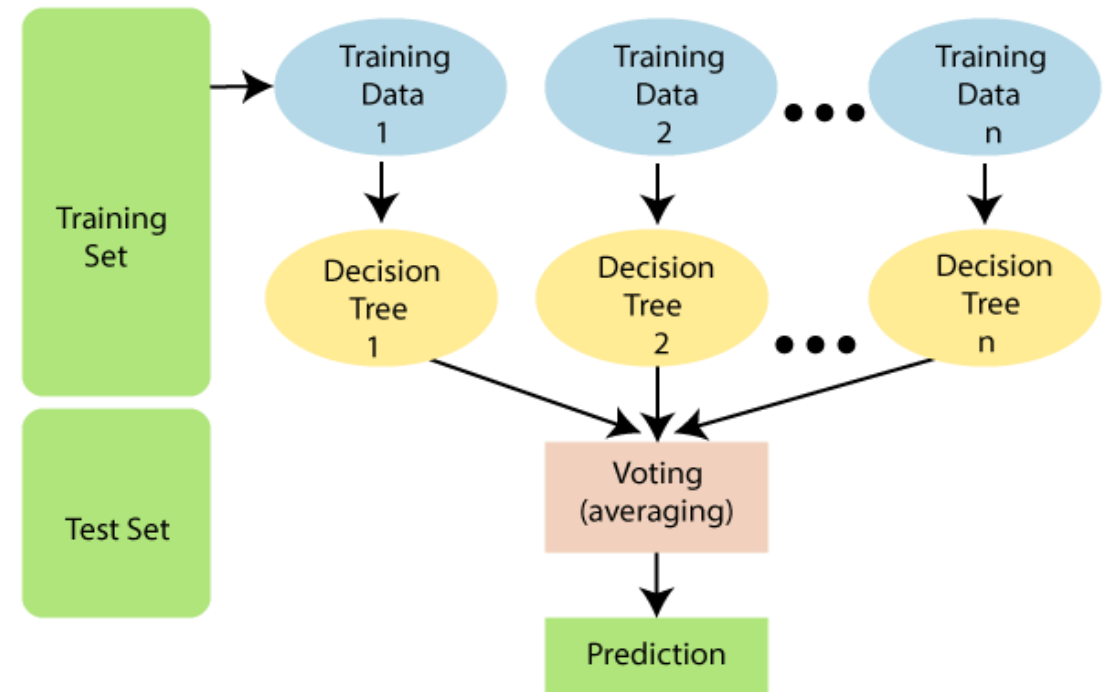


Figure 6: Random forest tree

Cont..

- **KNN(k-nearest neighbor)** is a lazy learning, non-parametric algorithm. It uses data with several classes to predict the classification of the new sample point. KNN is non-parametric since it doesn't make any assumptions on the data being studied, i.e., the model is distributed from the data.

Algorithm 2

Step 1 – For implementing any algorithm, we need dataset. So during the first step of KNN is to load the training as well as test data.

Step 2 – Next, we need to choose the value of K i.e. the nearest data points. K can be any integer.

Step 3 – Calculate the distance between test data and each row of training data with the help of any of the method namely: Euclidean, Manhattan or Hamming distance. The most commonly used method to calculate distance is Euclidean.

Step 4 – End

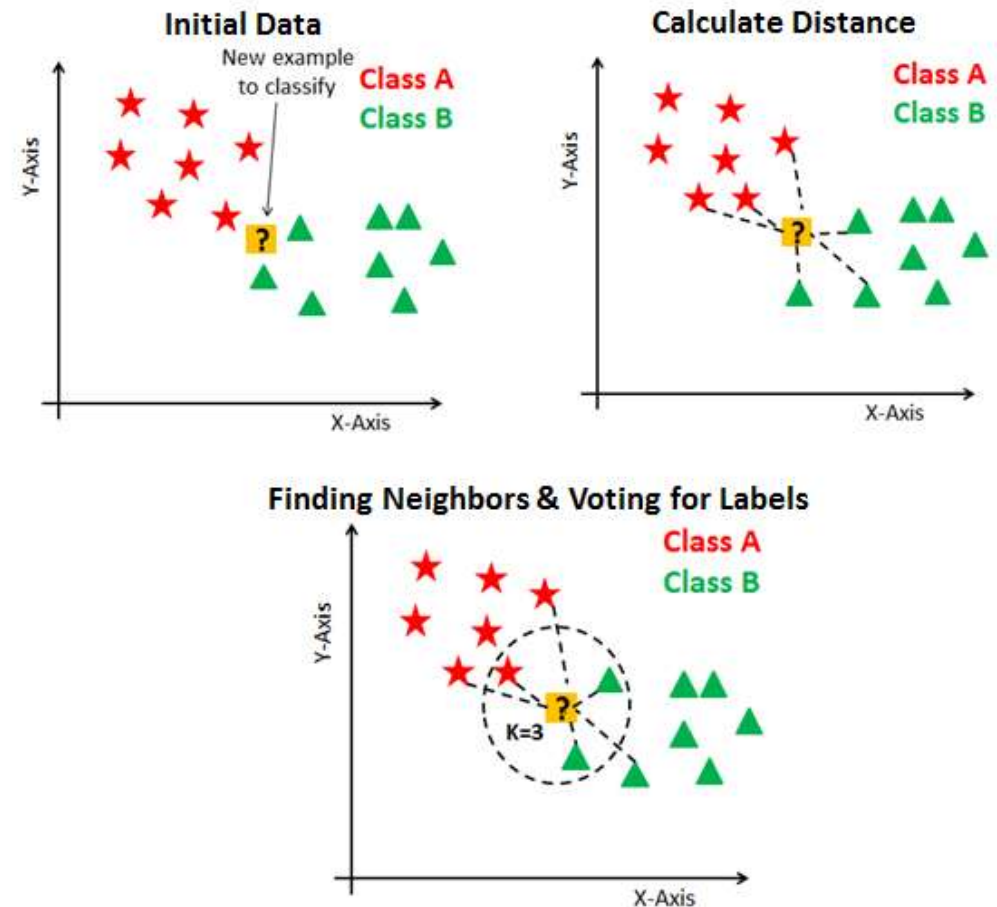


Figure 7: KNN Algorithm

Cont..

- **Naive Bayes** classifiers are a collection of classification algorithms based on Bayes' Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle, i.e. every pair of features being classified is independent of each other.

Algorithm 3

1. Convert the given dataset into frequency tables.
2. Generate Likelihood table by finding the probabilities of given features.
3. Now, use Bayes theorem to calculate the posterior probability.

A handwritten diagram illustrating Bayes' Theorem. The central equation is $P(A|B) = \frac{P(B|A) P(A)}{P(B)}$. Annotations with arrows explain each term:
 - Above the equation, an arrow points from "THE PROBABILITY OF 'B' BEING TRUE GIVEN THAT 'A' IS TRUE" to $P(B|A)$.
 - To the right, an arrow points from "THE PROBABILITY OF 'A' BEING TRUE" to $P(A)$.
 - Below the equation, an arrow points from "THE PROBABILITY OF 'A' BEING TRUE GIVEN THAT 'B' IS TRUE" to $P(A|B)$.
 - Below the denominator, an arrow points from "THE PROBABILITY OF 'B' BEING TRUE" to $P(B)$.

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

Figure 8: Naïve Bayes Algorithm

Cont..

- **Decision tree** are simple to implement and equally easy to interpret.
- Decision Trees are a non-parametric supervised learning method used for both classification and regression tasks. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features

Algorithm 3

Step-1: Begin the tree with the root node, says S, which contains the complete dataset.

Step-2: Find the best attribute in the dataset using **Attribute Selection Measure (ASM)**.

Step-3: Divide the S into subsets that contains possible values for the best attributes.

Step-4: Generate the decision tree node, which contains the best attribute.

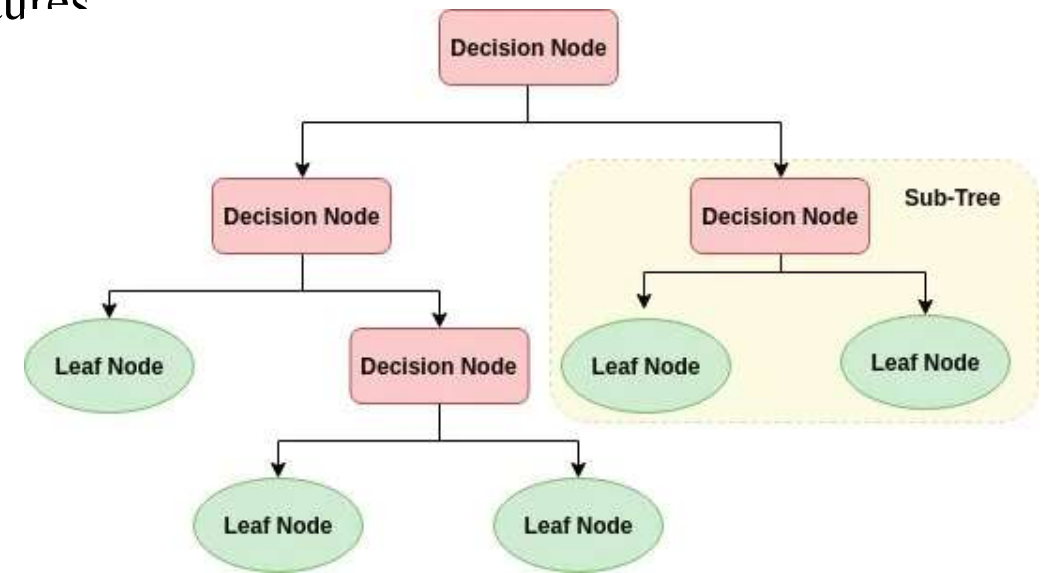


Figure 9: Decision Tree Algorithm

Cont..

- **SVM(Support Vector Machine)** is a supervised machine learning algorithm used for both classification and regression. The objective of SVM algorithm is to find a hyperplane in an N-dimensional space that distinctly classifies the data points. The dimension of the hyperplane depends upon the number of features. The objective of the support vector machine algorithm is to find a hyperplane in an N-dimensional space (N — the number of features) that distinctly classifies the data points.
- Hyperplanes are decision boundaries that help classify the data points. Data points falling on either side of the hyperplane can be attributed to different classes.

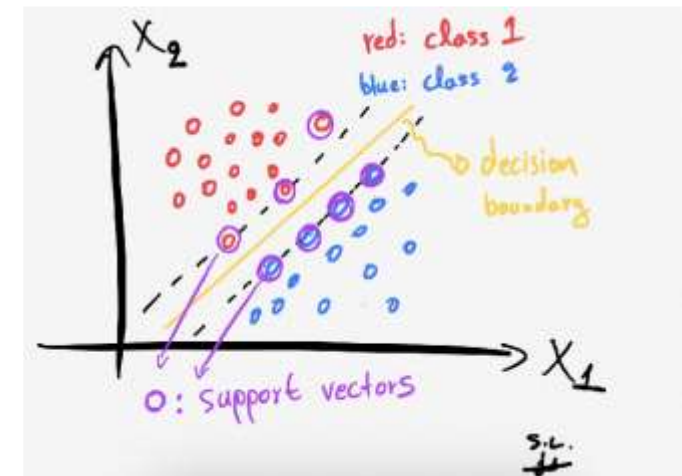


Figure 10: SVM Algorithm

Cont..

Neural Networks also known as Artificial Neural Networks (ANNs) or Simulated Neural Networks (SNNs), are a subset of machine learning and are at the heart of deep learning algorithms. Their name and structure are inspired by the human brain, mimicking the way that biological neurons signal to one another.

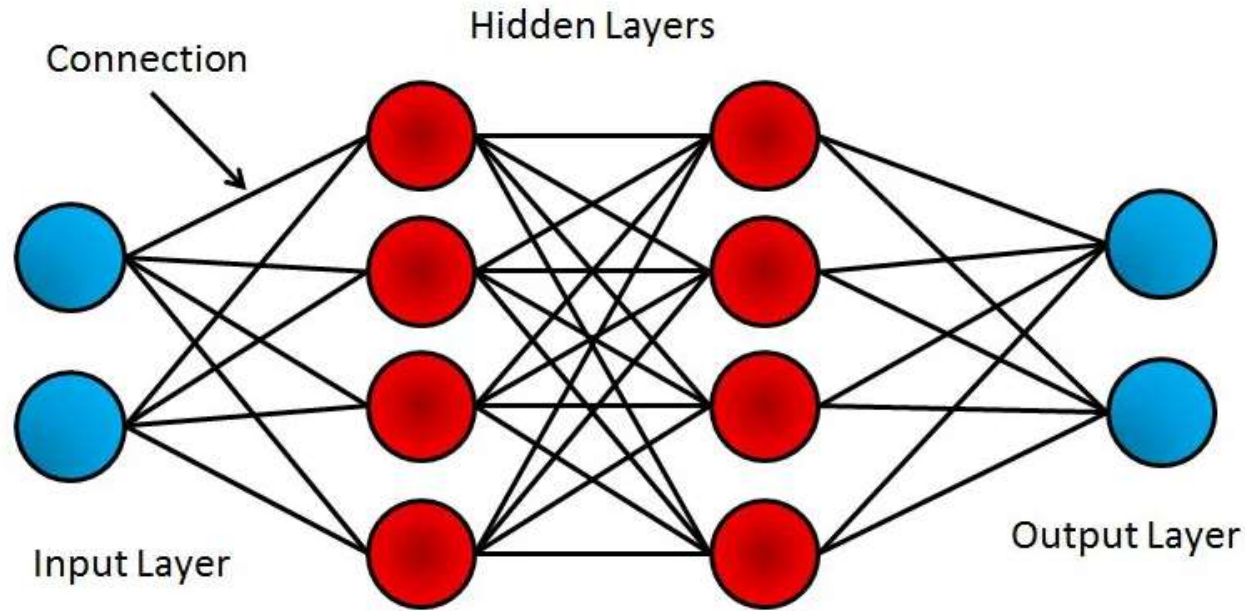


Figure 11: Neural Network

Results and Analysis

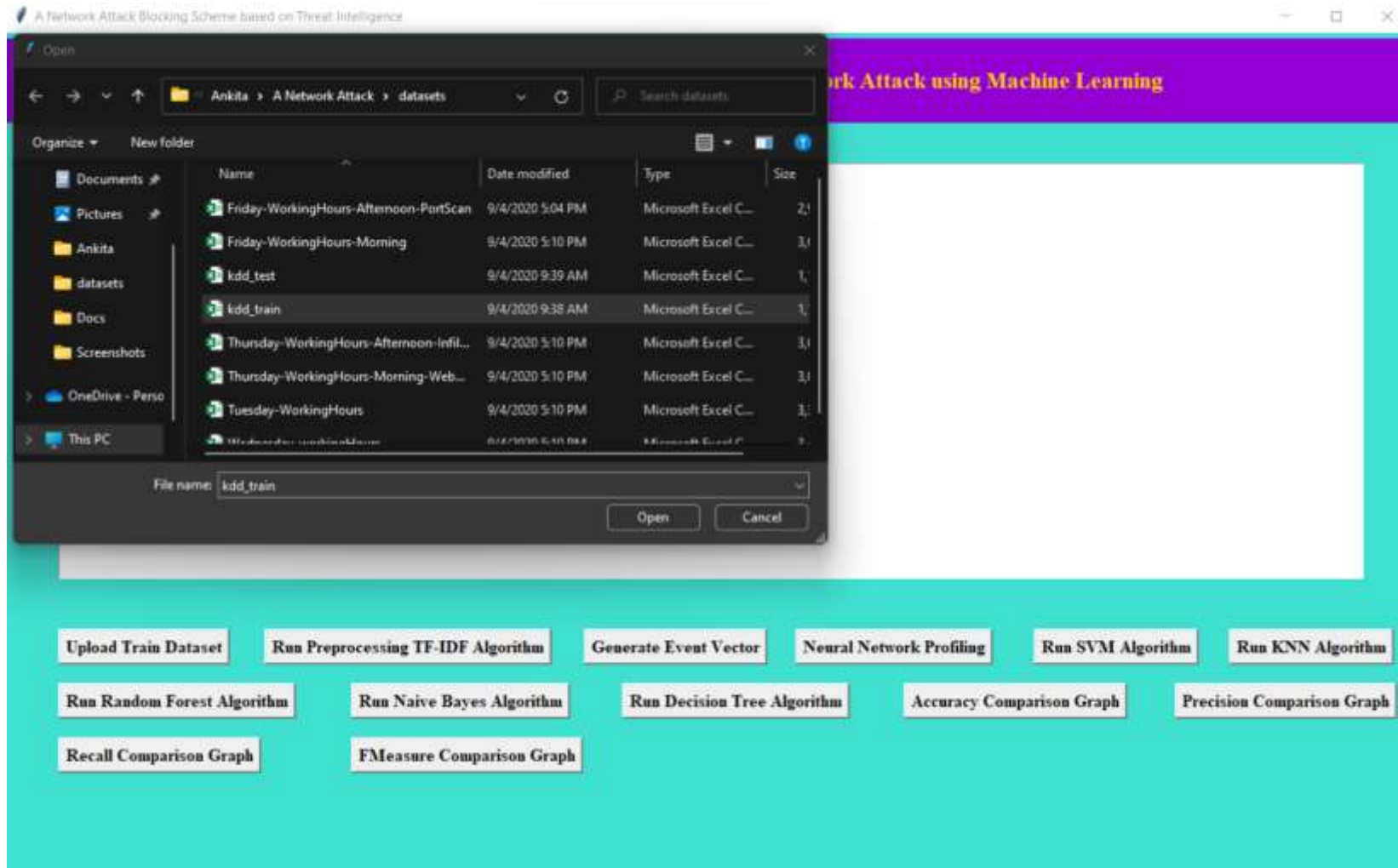


Figure 12: Uploading the dataset

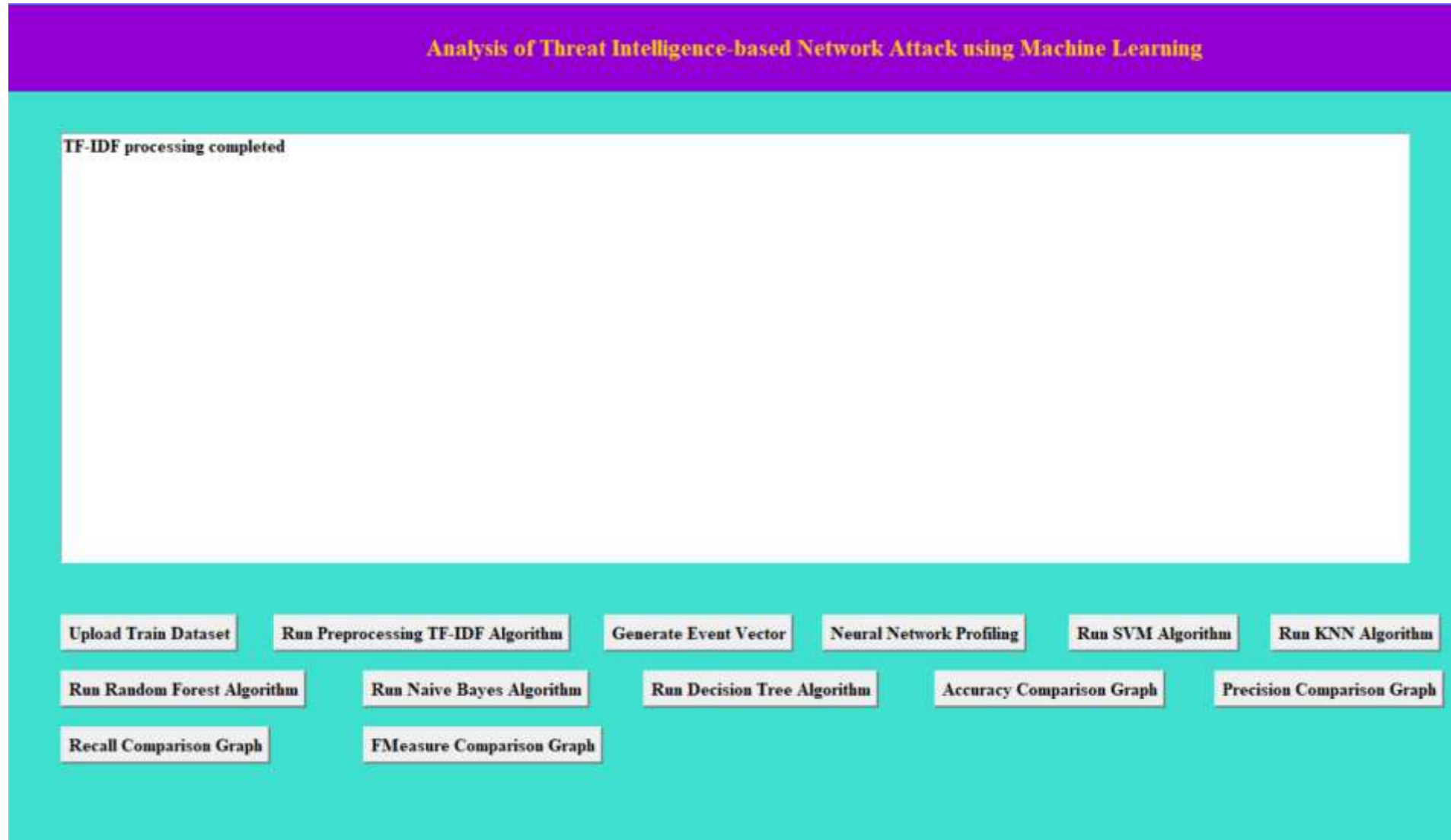


Figure 13: Running the TF-IDF Algorithm

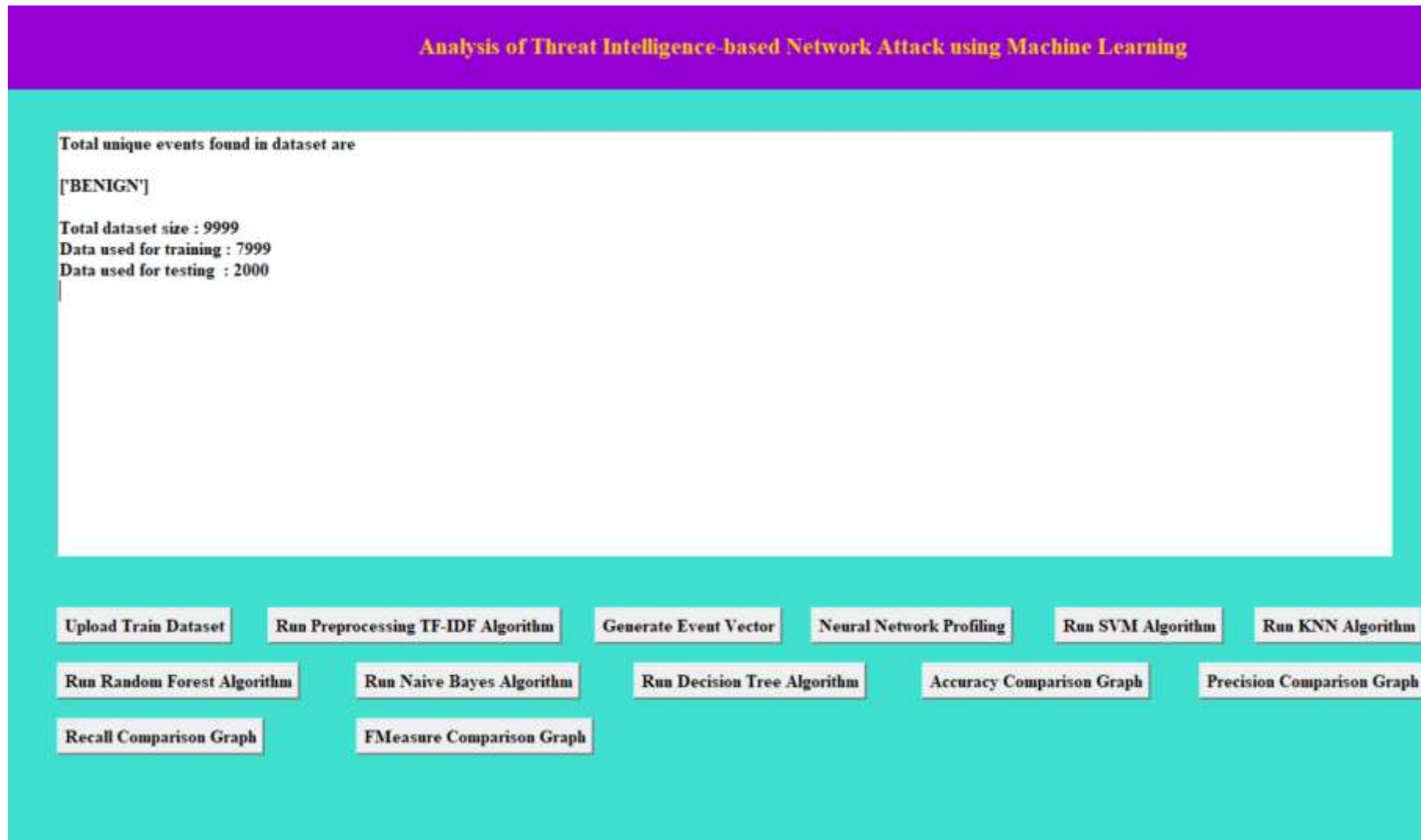


Figure 14: Generating the Event Vector

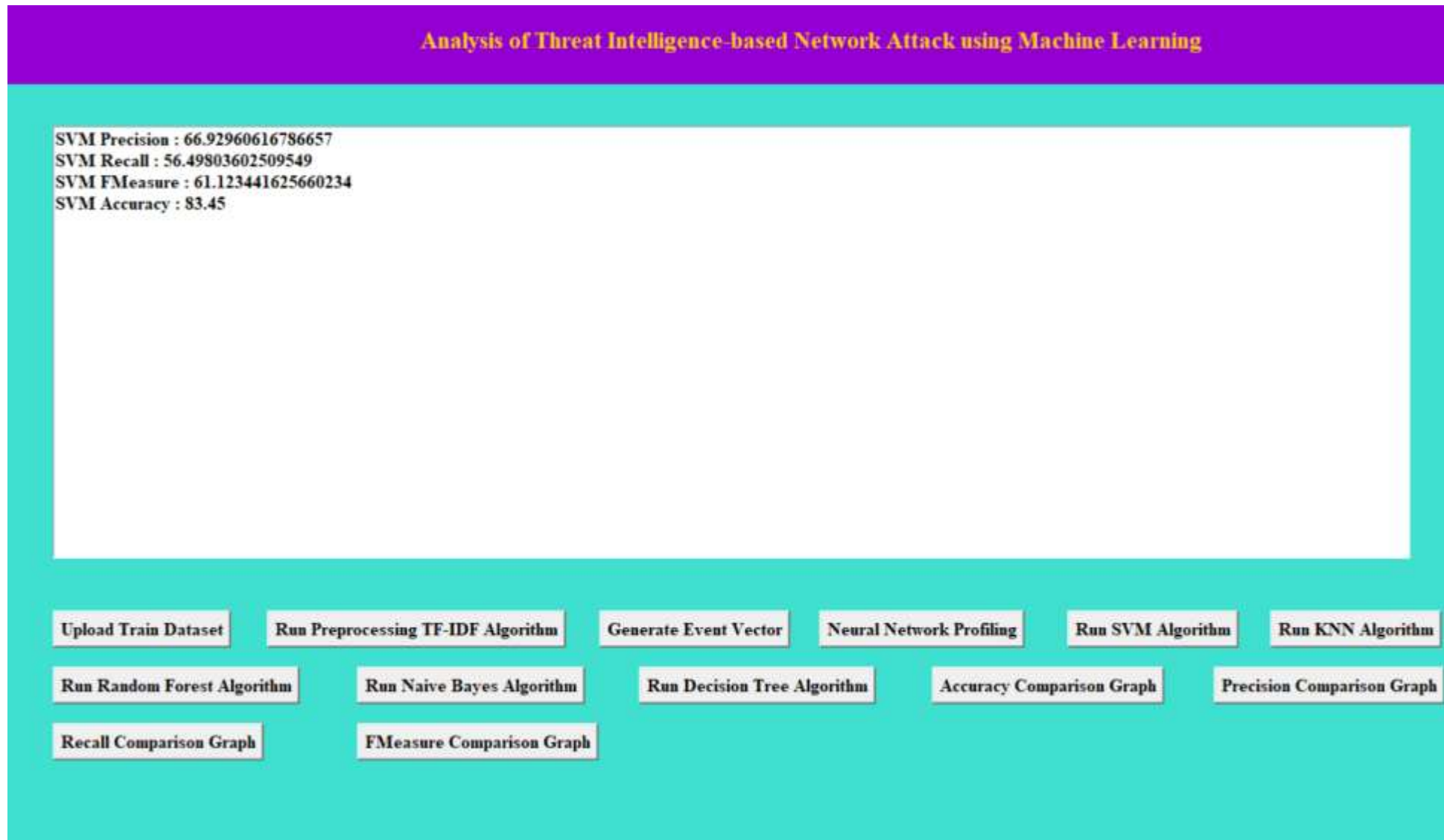


Figure 15: Support Vector Machine Algorithm Prediction Results

Analysis of Threat Intelligence-based Network Attack using Machine Learning

KNN Prediction Results

KNN Precision : 60.210288265210046
KNN Recall : 40.60731478935236
KNN FMeasure : 45.967215520505235
KNN Accuracy : 71.7

Upload Train Dataset

Run Preprocessing TF-IDF Algorithm

Generate Event Vector

Neural Network Profiling

Run SVM Algorithm

Run KNN Algorithm

Run Random Forest Algorithm

Run Naive Bayes Algorithm

Run Decision Tree Algorithm

Accuracy Comparison Graph

Precision Comparison Graph

Recall Comparison Graph

FMeasure Comparison Graph

Figure 16: KNN Algorithm Prediction Results

Analysis of Threat Intelligence-based Network Attack using Machine Learning

Random Forest Prediction Results

Random Forest Precision : 58.47078674476602

Random Forest Recall : 44.16209685318424

Random Forest FMeasure : 50.06854164180595

Random Forest Accuracy : 76.9

Upload Train Dataset

Run Preprocessing TF-IDF Algorithm

Generate Event Vector

Neural Network Profiling

Run SVM Algorithm

Run KNN Algorithm

Run Random Forest Algorithm

Run Naive Bayes Algorithm

Run Decision Tree Algorithm

Accuracy Comparison Graph

Precision Comparison Graph

Recall Comparison Graph

FMeasure Comparison Graph

Figure 17: Random Forest Algorithm Prediction Results

Analysis of Threat Intelligence-based Network Attack using Machine Learning

Naive Bayes Prediction Results

Naive Bayes Precision : 33.99400561655267
Naive Bayes Recall : 18.919307283876822
Naive Bayes FMeasure : 22.322845572845573
Naive Bayes Accuracy : 67.80000000000001

Upload Train Dataset

Run Preprocessing TF-IDF Algorithm

Generate Event Vector

Neural Network Profiling

Run SVM Algorithm

Run KNN Algorithm

Run Random Forest Algorithm

Run Naive Bayes Algorithm

Run Decision Tree Algorithm

Accuracy Comparison Graph

Precision Comparison Graph

Recall Comparison Graph

FMeasure Comparison Graph

Figure 18: Naïve Bayes Algorithm Prediction Results

Analysis of Threat Intelligence-based Network Attack using Machine Learning

Decision Tree Prediction Results

Decision Tree Precision : 3.6566666666666663

Decision Tree Recall : 6.666666666666667

Decision Tree FMeasure : 4.722850069960176

Decision Tree Accuracy : 54.85

Upload Train Dataset

Run Preprocessing TF-IDF Algorithm

Generate Event Vector

Neural Network Profiling

Run SVM Algorithm

Run KNN Algorithm

Run Random Forest Algorithm

Run Naive Bayes Algorithm

Run Decision Tree Algorithm

Accuracy Comparison Graph

Precision Comparison Graph

Recall Comparison Graph

FMeasure Comparison Graph

Figure 19: Decision Tree Algorithm Prediction Results

Analysis of Threat Intelligence-based Network Attack using Machine Learning

Deep Learning LSTM Extension Accuracy

LSTM Accuracy : 94.18455362319946

LSTM Precision : 19.615384615384617

LSTM Recall : 76.92307692307692

LSTM Fmeasure : 38.25526009826351

Deep Learning CNN Accuracy

CNN Accuracy : 99.49992299079895

CNN Precision : 90.90909090909092

CNN Recall : 5.0

CNN Fmeasure : 9.945300845350571

Upload Train Dataset

Run Preprocessing TF-IDF Algorithm

Generate Event Vector

Neural Network Profiling

Run SVM Algorithm

Run KNN Algorithm

Run Random Forest Algorithm

Run Naive Bayes Algorithm

Run Decision Tree Algorithm

Accuracy Comparison Graph

Precision Comparison Graph

Recall Comparison Graph

FMeasure Comparison Graph

Figure 20: Neural Network Profiling Prediction Results

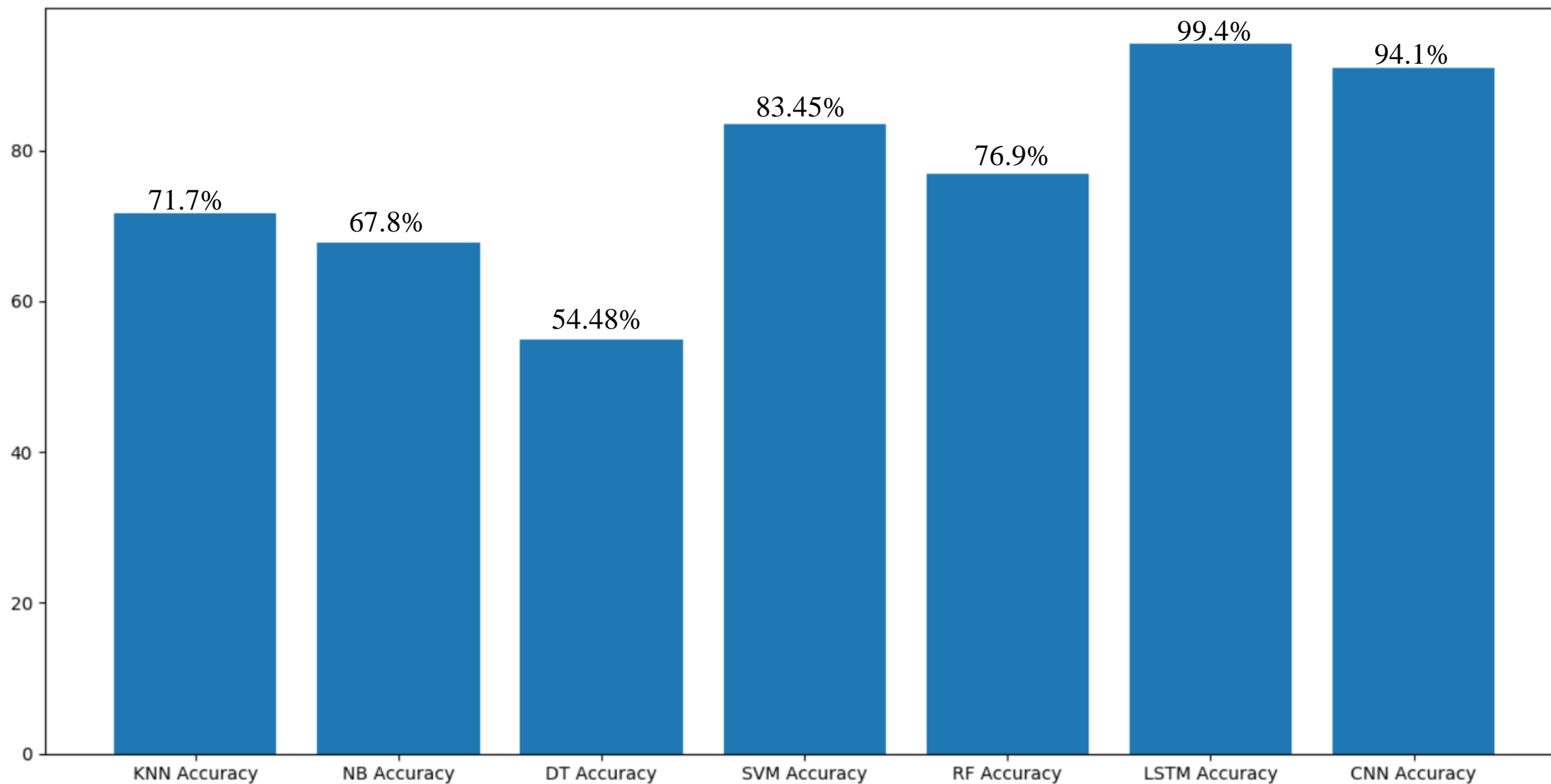


Figure 21: Comparison Graph of all Algorithm

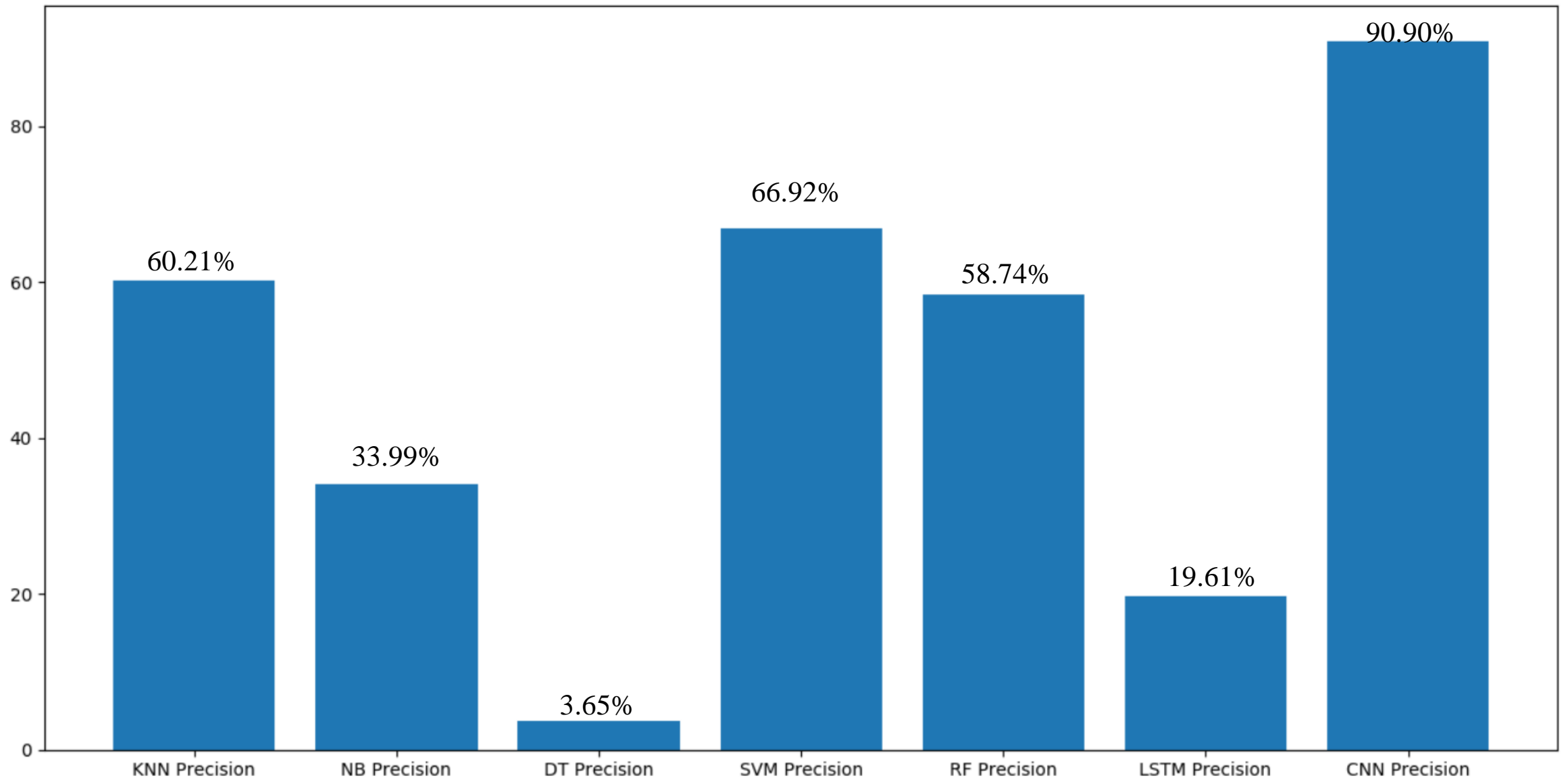


Figure 22: Precision Graph of all Algorithm

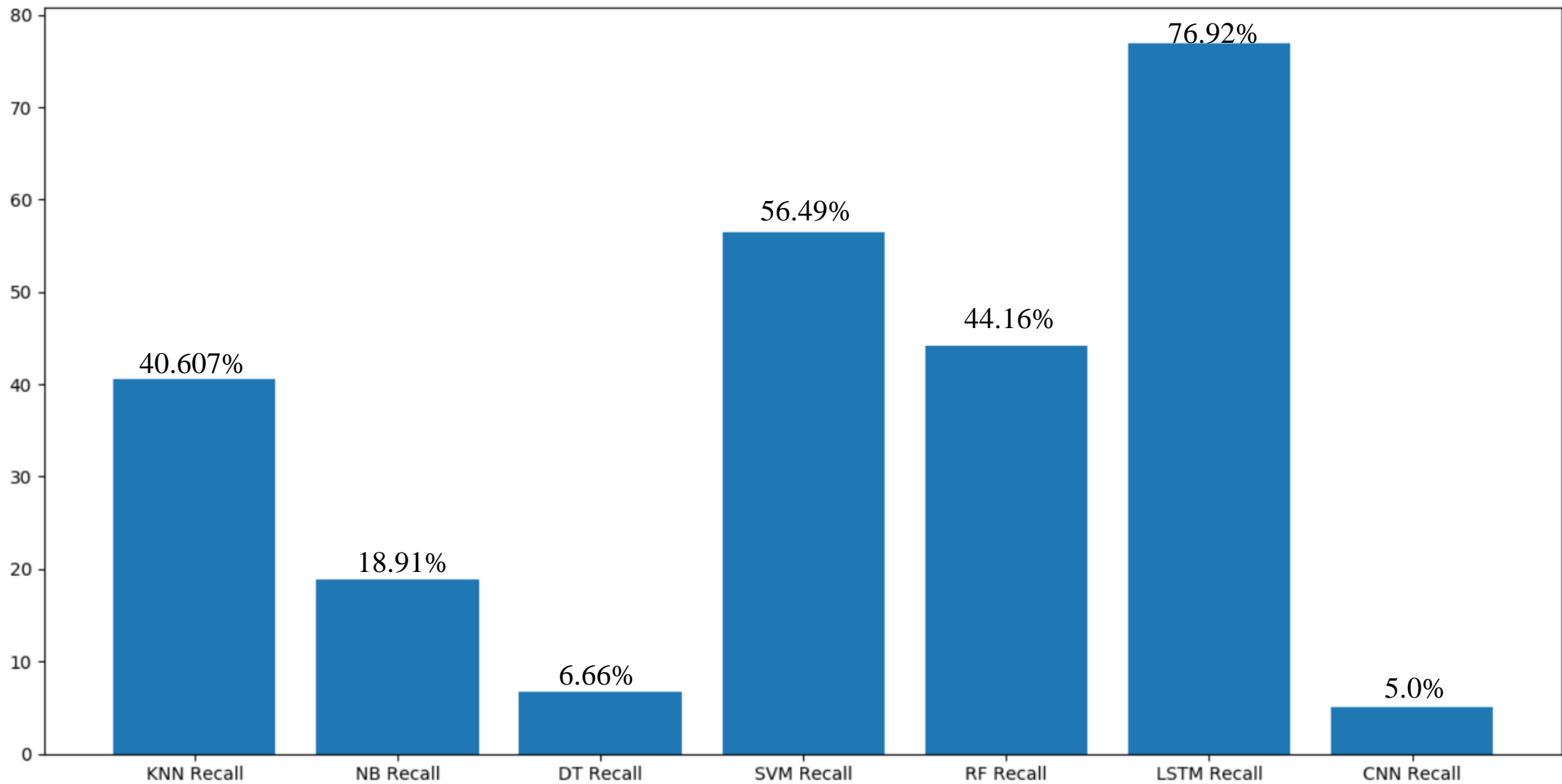


Figure 23: Recall Graph of all Algorithm

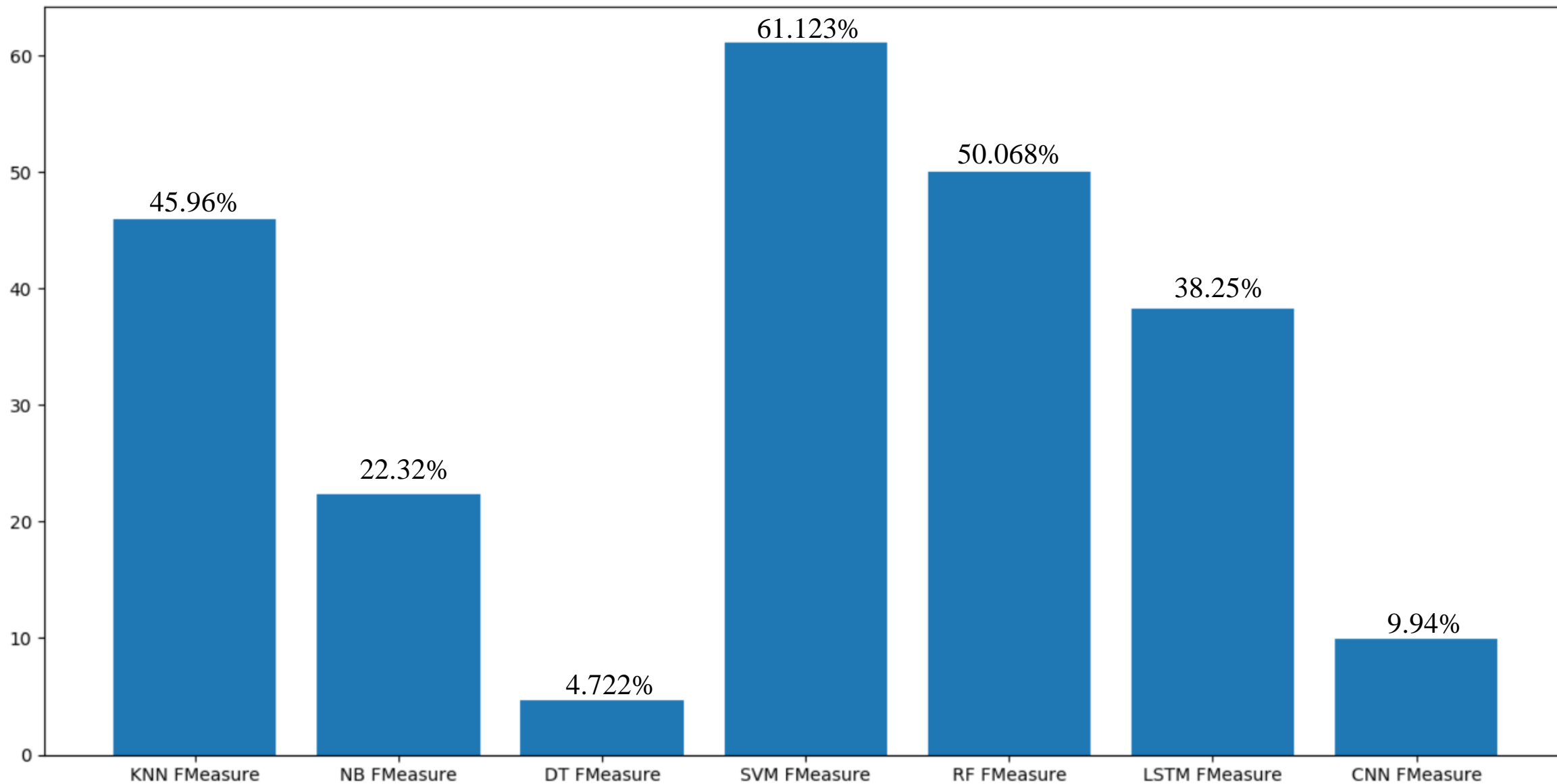


Figure 24: F1 measure Graph of all Algorithm

Conclusion

- The results are examined utilizing network attack data sets, demonstrating that it is one of the best competitor data sets for threat intelligence analysis and testing. The TI model improves both the accuracy and the detection time.
- After experimenting with various numbers of traits, it was discovered that 22 of the 76 attributes have a categorization accuracy of roughly 99 percent. And, of the 22 features, 11 are more important to total classification accuracy, while the remaining 9 are less important.
- The classification accuracy is 99.4 percent when these 13 out of 22 features are used.

Future Work

- Can reduce the number false positive rates so that the performance increases by this.
- To develop a system model with a higher precision rate while using fewer characteristics to improve overall accuracy.
- Datasets like NSL-KDD can be considered in future research.

References

- [1] Mavroeidis, Vasileios, and Siri Bromander. "Cyber threat intelligence model: an evaluation of taxonomies, sharing standards, and ontologies within cyber threat intelligence." *2017 European Intelligence and Security Informatics Conference (EISIC)*. IEEE, 2017.
- [2] Blanco-Medina, Pablo, et al. "Detecting vulnerabilities in critical infrastructures by classifying exposed industrial control systems using deep learning." *Applied Sciences* 11.1 (2021): 367.
- [3] Chen, Zhijiang, et al. "A cloud computing based network monitoring and threat detection system for critical infrastructures." *Big Data Research* 3 (2016): 10-23.
- [4] Al-Mohannadi, Hamad, Irfan Awan, and Jassim Al Hamar. "Analysis of adversary activities using cloud-based web services to enhance cyber threat intelligence." *Service Oriented Computing and Applications* 14.3 (2020): 175-187.

THANK YOU