# MACHINE LEARNING APPROACHES FOR IDENTIFYING NETWORK CYBER THREATS

**Dr. S.KISHORE VERMA, Associate Professor**

**S.SHARATH KUMAR, S.VENKATESH, M.KEERTHI , K.MANIKANTA**

SRI INDU COLLEGE OF ENGINEERING AND TECHNOLOGY, **Sheriguda (V),**

**Ibrahimpatnam (M), Rangareddy Dist –501 510**

## ABSTRACT

Contrasted with the past, improvements in PC and correspondence innovations have given broad and propelled changes. The use of new innovations give incredible advantages to people, organizations, and governments, be that as it may, messes some up against them. For instance, the protection of significant data, security of put away information stages, accessibility of information and so forth. Contingent upon these issues, digital fear based oppression is one of the most significant issues in this day and age. Digital fear, which made a great deal of issues people and establishments, has arrived at a level that could undermine open and nation security by different gatherings, for example, criminal association, proficient people and digital activists. Along these lines, Intrusion Detection Systems (IDS) has been created to maintain a strategic distance from digital assaults. Right now, learning the bolster support vector machine (SVM) calculations were utilized to recognize port sweep endeavors dependent on the new CICIDS2017 dataset with 97.80%, 69.79% precision rates were accomplished individually. Rather than SVM we can introduce some other algorithms like random forest, CNN, ANN where these algorithms can acquire accuracies like SVM – 93.29, CNN – 63.52, Random Forest – 99.93, ANN – 99.11

KEYWORDS : Digital fear,  :Intrusion Detection, bolster support vector machine (SVM), CICIDS2017 dataset

## 1.  INTRODUCTION

Contrasted with the past, improvements in PC and correspondence innovations have given broad and propelled changes. The use of new innovations give incredible advantages to people, organizations, and governments, be that as it may, messes some up against them. For instance, the protection of significant data, security of put away information stages, accessibility of information and so forth. Contingent upon these issues, digital fear based oppression is one of the most significant issues in

this day and age. Digital fear, which made a great deal of issues people and establishments, has arrived at a level that could undermine openand nation security by different gatherings, for example, criminal association, proficient people and digital activists. Along these lines, Intrusion Detection Systems (IDS) has been created to maintain a strategic distance from digital assaults. Right now, learning the bolster support vector machine (SVM) calculations were utilized to recognize port sweep endeavors dependent on the new CICIDS2017 dataset with 97.80%, 69.79% precision rates were accomplished individually. Rather than SVM we can introduce some other algorithms like random forest, CNN, ANN where these algorithms can acquire accuracies like SVM – 93.29, CNN – 63.52, Random Forest – 99.93, ANN – 99.11.

## 1.1 MOTIVATION

The use of new innovations give incredible advantages to people, organizations, and governments, be that as it may, messes some up against them. For instance, the protection of significant data, security of put away information stages, accessibility of information and so forth. Contingent upon these issues, digital fear based oppression is one of the most significant issues in this day and age. Digital fear, which made a great deal of issues people and establishments, has arrived at a level that could undermine open and nation security by different gatherings, for example, criminal association, proficient people and digital activists. Along these lines, Intrusion Detection Systems (IDS) has been created to maintain a strategic distance from digital assaults.

## 1.2 Objectives

Objective of this project is to detect cyber attacks by using machine learning algorithms like

ANN , CNN, Random forest

## 2. LITERATURE SURVEY

**R. Christopher, "Port scanning techniques and the defense against them," SANS Institute, 2001.**

Port Scanning is one of the most popular techniques attackers use to discover services that they can exploit to break into systems. All systems that are connected to a LAN or the Internet via a modem run services that listen to well-known and not so well-known ports. By port scanning, the attacker can find the following information about the targeted systems: what services are running, what users

own those services, whether anonymous logins are supported, and whether certain network services require authentication. Port scanning is accomplished by sending a message to each port, one at a time. The kind of response received indicates whether the port is used and can be probed for further weaknesses. Port scanners are important to network security technicians because they can reveal possible security vulnerabilities on the targeted system.Every publicly available system has ports that are open and available for use. The object is to limit the exposure of open ports to authorized users and to deny access to the closed ports.

**S. Staniford, J. A. Hoagland, and J. M. McAlerney, "Practical automated detection of stealthy portscans," Journal of Computer Security, vol. 10, no. 1-2, pp. 105–136, 2002.**

Portscanning is a common activity of considerable importance. It is often used by computer attackers to characterize hosts or networks which they are considering hostile activity against. Thus it is useful for system administrators and other network defenders to detect portscans as possible preliminaries to a more serious attack. It is also widely used by network defenders to understand and find vulnerabilities in their own networks. Thus it is of considerable interest to attackers to determine whether or not the defenders of a network are portscanning it regularly. However, defenders will not usually wish to hide their portscanning, while attackers will. For definiteness, in the remainder of this paper, we will speak of the attackers scanning the network, and the defenders trying to detect the scan. One concerns whether portscanning of remote networks without permission from the owners is itself a legal and ethical activity. This is presently a grey area in most jurisdictions.So we think it reasonable to consider a portscan as at least potentially hostile, and to report it to the administrators of the remote network from whence it came.

However, this paper is focussed on the technical questions of how to detect portscans, which are independent of what significance one imbues them with, or how one chooses to respond to them.

In the next section, we discuss a variety of prior work on portscan detection. Then we present the algorithms that we propose to use, and give some very preliminary data justifying our approach. Finally, we consider possibleextensions to this work, along with other applications that might be considered.The primary purpose is that of gathering information about the reachability and status of certain combinations of IP address and port (either TCP or UDP).The secondary purpose is to flood intrusion detection systems with alerts, with the intention of distracting the network defenders or preventing them from doing their jobs. We will use the term scan footprint for the set of port/IP combinations which the attacker is interested in characterizing. The most common type of portscan

footprint at present is a horizontal scan. By this, we mean that an attacker has an exploit for a particular service, and is interested in finding any hosts that expose that service.

**M. C. Raja and M. M. A. Rabbani, "Combined analysis of support vector machine and principle component analysis for ids," in IEEE International Conference on Communication and Electronics Systems, 2016, pp. 1–5.**

Compared to the past security of networked systems has become a critical universal issue that influences individuals, enterprises and governments.Based on the detection technique, intrusion detection is classified into anomaly-based and signature-based. The authors examined the performance of these features with different algorithms that included:

K-Nearest Neighbor (KNN), Adaboost, Multi-Layer Perceptron (MLP), Naïve Bayes, Random Forest (RF), Iterative Dichotomiser 3 (ID3) and Quadratic Discriminant Analysis (QDA). The highest precision value was 0.98 with RF and ID3 [4]. The execution time (time to build the model) was 74.39 s. This is while the execution time for our proposed system using Random Forest is 21.52 s with a comparable processor. Some of them were discussed here..The developers used statistical metrics such as minimum, maximum, mean and standard deviation to encapsulate the network events into a set of certain features which include: 1. The distribution of the packet size 2. The number of packets per flow 3. The size of the payload 4. The request time distribution of the protocols 5. Certain patterns in the payload Moreover, CICIDS2017 covers various attack scenarios that represent common attack families. The attacks include Brute Force Attack, Heart Bleed Attack, Botnet, DoS Attack, Distributed DoS (DDoS) Attack , Web Attack, and Infiltration Attack.Moreover SVM requires the processing of raw features for classification which increases the architecture complexity and decreases the accuracy of detecting intrusion 1.3 DEEP LEARNING Deep learning is an improved machine learning technique for feature extraction, perception and learning of machines. Deep learning algorithms performs their operations using multiple consecutive layers. There are many application areas for Deep Learning, which covers such as Image Processing, Natural Language Processing, biomedical, Customer Relationship Management automation, Vehicle autonomous systems and others

## 3. PROBLEM STATEMENT

The existing system for detecting cyber attacks in networks typically relies on traditional signature-based detection methods, which rely on pre-defined patterns of known attacks to identify new

attacks. These methods are limited in their ability to detect new and evolving types of attacks and may generate false positives or negatives.

## 3.1 LIMITATION OF SYSTEM

Strict Regulations, Difficult to work with for non-technical users, Restrictive to resources, Constantly needs Patching, Constantly being attacked

## 4. PROPOSED SYSTEM

The proposed system for detecting cyber attacks in networks using machine learning techniques aims to address the limitations of existing systems. The proposed system uses machine learning algorithms to analyze network traffic patterns and detect anomalies that may indicate a cyber attack. The system can be trained to recognize normal network behavior and identify deviations from this behavior that may indicate an attack. The proposed system can also be enhanced with additional features such as real-time monitoring, automatic response mechanisms, and integration with other security systems. Real-time monitoring allows the system to detect attacks as they occur, and automatic response mechanisms can help mitigate the damage caused by attacks. Integration with other security systems, such as firewalls and intrusion detection systems, can improve overall network security and enhance the effectiveness of the proposed system.

### 4.1 Advantages

Protection from malicious attacks on your network. Deletion and/or guaranteeing malicious elements within a preexisting network. Prevents users from unauthorized access to the network. Deny's programs from certain resources that could be infected. Securing confidential information

## 5. IMPLEMENTATION

**5.1 DATA COLLECTION**: Gathering essential data (like network traffic details) from the CICIDS2017 dataset, vital for identifying port scan attempts and potential security threats.

**5.2 DATA PRE-PROCESSING:** Cleaning, handling missing values, and organizing the data to make it compatible and optimal for machine learning algorithms to process effectively and accurately.

**5.3 FEATURE EXTRATION**: Selecting and deriving crucial attributes (e.g., packet size, protocol types) from the organized data that serve as inputs for machine learning models to identify patterns related to port scan attempts.

**5.4 EVALUATION MODEL**: Applying diverse algorithms like SVM, Random Forest, CNN, and ANN to the extracted features, training these models on a subset of data, and assessing their performance in accurately detecting port scan attempts to enhance cybersecurity measures.
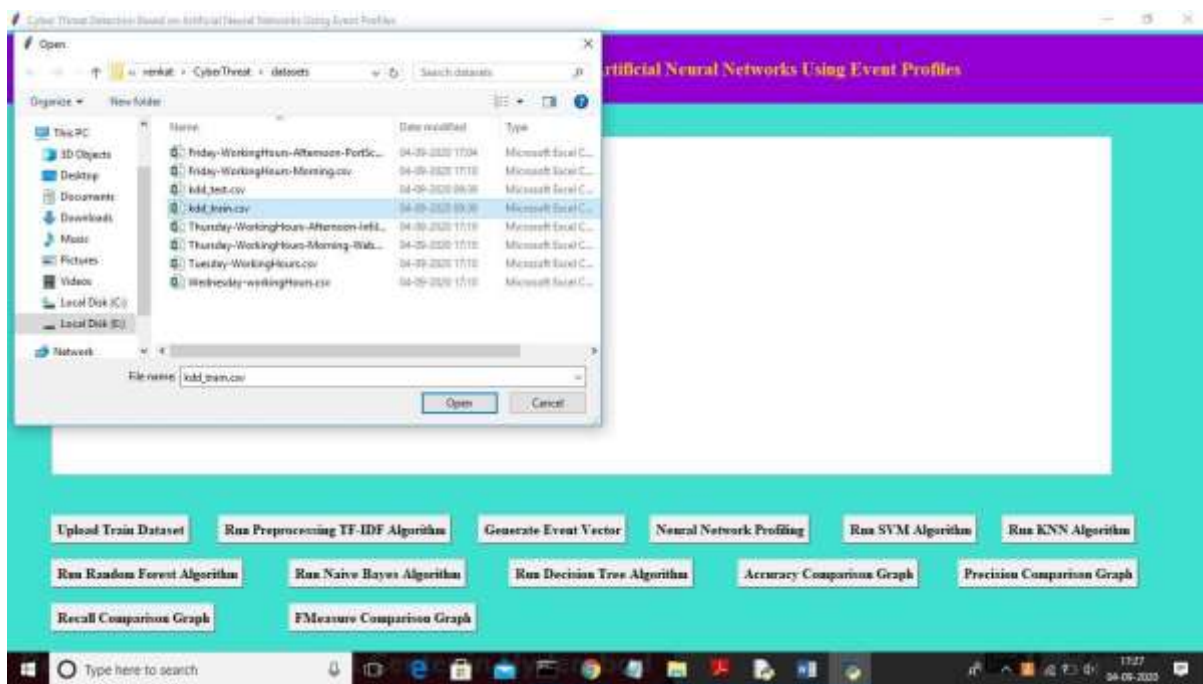
## 6. ARCHITECTURE



## 7. EXPERMENTAL RESULTS

To run project double click on 'run.bat' file to get below screen



In above screen click on 'Upload Train Dataset' button and upload dataset

In above screen uploading 'kdd_train.csv' dataset and after upload will get below screen

In above screen we can see dataset contains 9999 records and now click on 'Run Preprocessing TF-IDF Algorithm' button to convert raw dataset into TF-IDF values



In above screen TF-IDF processing completed and now click on 'Generate Event Vector' button to create vector from TF-IDF with different events

In above screen we can see total different unique events names and in below we can see dataset total size and application using 80% dataset (7999 records) for training and using 20% dataset (2000 records) for testing. Now dataset train and test events model ready and now click on 'Neural Network Profiling' button to create LSTM and CNN model

In above screen LSTM model is generated and its epoch running also started and its starting accuracy is 0.94. Running for entire dataset may take time so wait till LSTM and CNN training process completed. Here dataset contains 7999 records and LSTM will iterate all records to filter and build model.



In above selected text we can see LSTM complete all iterations and in below lines we can see CNN model also starts execution

In above screen CNN also starts first iteration with accuracy as 0.72 and after completing all iterations 10 we got filtered improved accuracy as 0.99 and multiply by 100 will give us 99% accuracy. So CNN is giving better accuracy compare to LSTM and now see below GUI screen with all details



In above screen we can see both algorithms accuracy, precision, recall and FMeasure values. Now click on 'Run SVM Algorithm' button to run existing SVM algorithm

In above screen we can see SVM algorithm output values and now click on 'Run KNN Algorithm' to run KNN algorithm

In above screen we can see KNN algorithm output values and now click on 'Run Random Forest Algorithm' to run Random Forest algorithm



In above screen we can see Random Forest algorithm output values and now click on 'Run Naïve Bayes Algorithm' to run Naïve Bayes algorithm
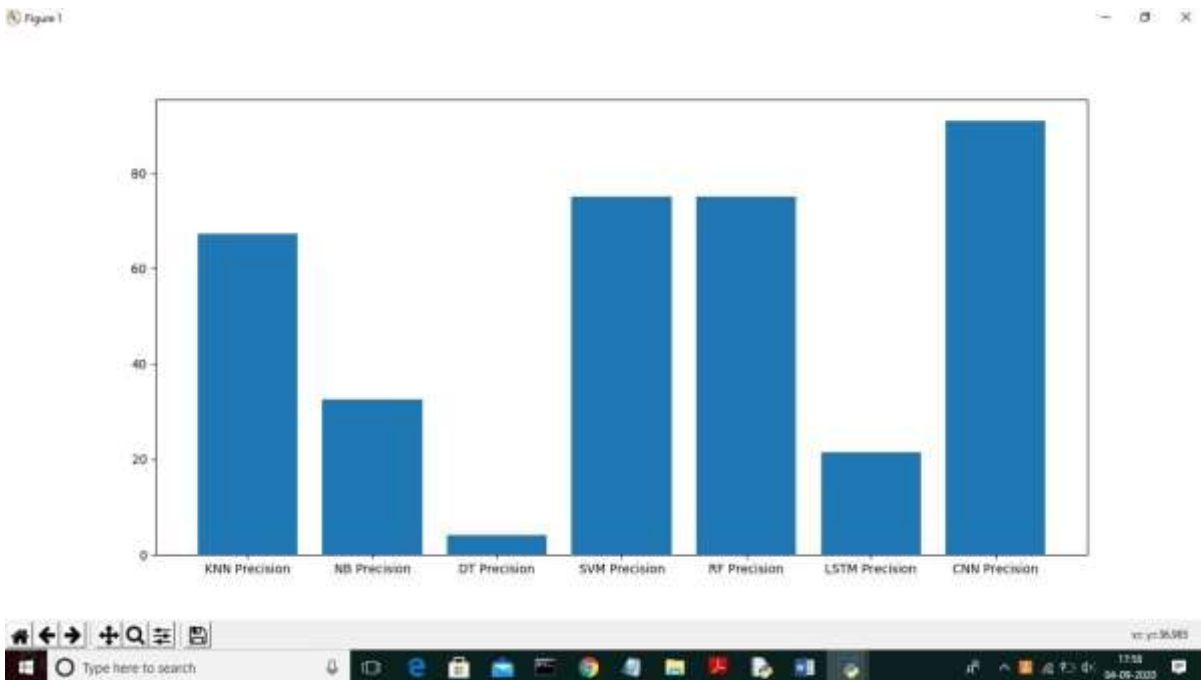
In above screen we can see Naïve Bayes algorithm output values and now click on 'Run Decision Tree Algorithm' to run Decision Tree Algorithm



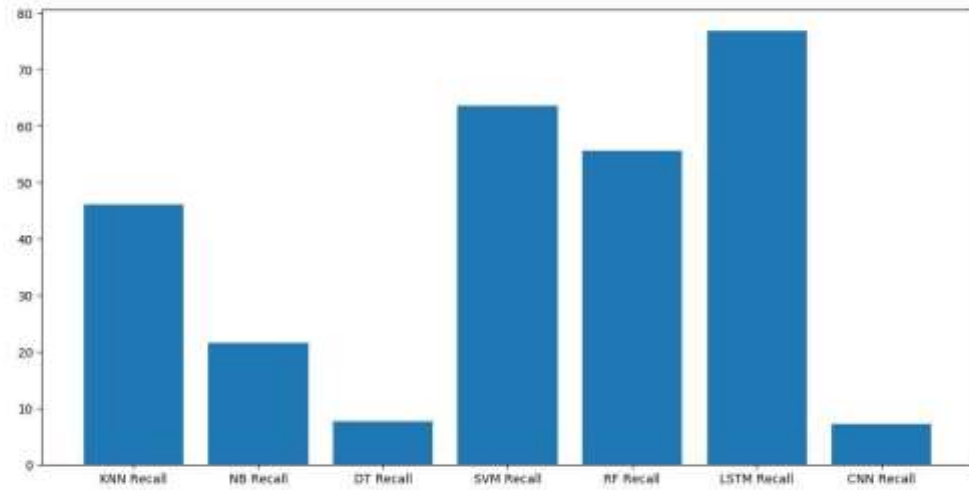Now click on 'Accuracy Comparison Graph' button to get accuracy of all algorithms

In above graph x-axis represents algorithm name and y-axis represents accuracy of those algorithms and from above graph we can conclude that LSTM and CNN perform well. Now click on Precision Comparison Graph' to get below graph



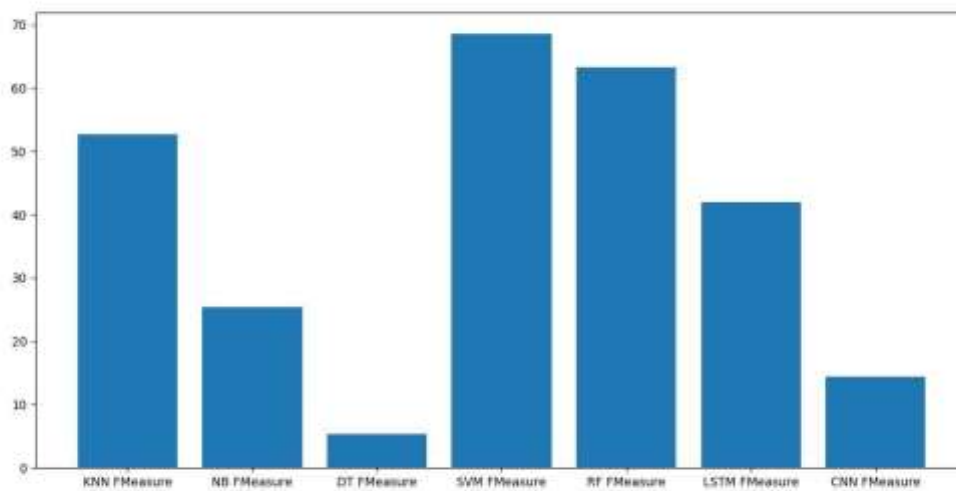In above graph CNN is performing well and now click on 'Recall Comparison Graph'

In above graph LSTM is performing well and now click on FMeasure Comparison Graph button to get below graph

From all comparison graph we can see LSTM and CNN performing well with accuracy, recall and precision.

## 7. CONCLUSIONS

Right now, estimations of help vector machine, ANN, CNN, Random Forest and profoundlearning calculations dependent on modern CICIDS2017 dataset were introduced relatively. Results show that the profound learning calculation performed fundamentally preferable outcomes over SVM, ANN, RF and CNN. We are going to utilize port sweep endeavors as well as other assault types with AI and profound learning calculations, apache Hadoop and sparkle innovations together dependent on this dataset later on. All these calculation helps us to detect the cyber attack in network. It happens in the way that when we consider long back years there may be so many attacks happened so when these attacks are recognized then the features at which values these attacks are happening will be stored in some datasets. So by using these datasets we are going to predict whether cyber attack is done or not. These predictions can be done by four algorithms like SVM, ANN, RF,

CNN this paper helps to identify which algorithm predicts the best accuracy rateswhich helps to predict best results to identify the cyber attacks happened or not.

## 8. REFERENCES

1.  K. Graves, Ceh: Official certified ethical hacker review guide: Exam 312-50. John Wiley& Sons, 2007.

2.  R. Christopher, "Port scanning techniques and the defense against them," SANS Institute, 2001.

3.  M. Baykara, R. Das¸, and I. Karado ˇgan, "Bilgi g ¨uvenli ˇgi sistemlerinde kullanilan arac¸larin incelenmesi," in 1st International Symposium on Digital Forensics and Security (ISDFS13), 2013, pp. 231–239.

4.  S. Staniford, J. A. Hoagland, and J. M. McAlerney, "Practical automated detection of stealthy portscans," Journal of Computer Security, vol. 10, no. 1-2, pp. 105–136, 2002.

5.  "SWGDE." [Online]. Available: https://www.swgde.org/. [Accessed: 30- Aug-2018].

6.  K. Ibrahimi and M. Ouaddane, "Management of intrusion detection systems based-kdd99: Analysis with lda and pca," in Wireless Networks and Mobile Communications (WINCOM), 2017 International Conference on. IEEE, 2017, pp. 1–6.

7.  N. Moustafa and J. Slay, "The significant features of the unsw-nb15 and the kdd99 datasets for network intrusion detection systems," in Building Analysis Datasets and Gathering Experience Returns for Security (BADGERS), 2015 4th International Workshop on. IEEE, 2015, pp. 25–31.

8.  "Welcome to Python.org." [Online]. Available: https://www.python.org/. [Accessed: 21-Aug-2018].

9.  S. M. Almansob and S. S. Lomte, "Addressing challenges for intrusion detection system using naive bayes and pca algorithm," in Convergence in Technology (I2CT), 2017 2nd International Conference for. IEEE, 2017, pp. 565–568.