**REPORT ON**

**"NETWORK ANOMALY DETECTION AND INSPECTION"**

**Title: AI based Anomaly Detection to detect malicious events in Network Environment**

Prepared by :

Sandip Das

TSSOT, Assam University

## **CONTENTS**

# __Introduction__

 Nowadays, there is a huge and growing concern about security in information and communication technology among the scientific community because any attack or anomaly in the network can greatly affect many domains such as national security, private data storage, social welfare, economic issues, and so on. Therefore, the anomaly detection domain is a broad research area, and many different techniques and approaches for this purpose have emerged through the years

Computer security has become a necessity due to proliferation of information technologies in everyday life. The mass usage of computerized systems has given rise to critical threats such as zero-day vulnerabilities, mobile threats, etc. Despite research in the security domain having increased significantly, are yet to be mitigated. The evolution of computer networks has greatly exacerbated computer security concerns, particularly internet security in today's networking environment and advanced computing facilities. Although Internet Protocols (IPs) were not designed to place a high priority on security issues, network administrators now have to handle a large variety of intrusion attempts by both individuals with malicious intent and large botnets. '***An anomaly is an observation which deviates so much from other observations as to arouse suspicions that it was generated by a different mechanism'***. Anomalies are considered important because they indicate significant but rare events and can prompt critical actions to be taken in a wide range of application domains; for example, an unusual traffic pattern in a network could mean that a computer has been hacked and data is transmitted to unauthorized destinations; anomalous behavior in credit card transactions could indicate fraudulent activities, and an anomaly in a MRI image may indicate the presence of a malignant tumor. Anomaly detection has been widely applied in countless application domains such as medical and public health, fraud detection, intrusion detection, industrial damage, image processing, sensor networks, robot's behavior and astronomical data.


**Types of anomalies**

Anomalies are referred to as patterns in data that do not conform to a well-defined characteristic of normal patterns. They are generated by a variety of abnormal activities, e.g., credit card fraud, mobile phone fraud, cyberattacks, etc., which are significant to data analysts. An important aspect of anomaly detection is the nature of the anomaly. An anomaly can be categorized in the following way.

- **Point anomaly:** When a particular data instance deviates from the normal pattern of the dataset, it can be considered a point anomaly. For a realistic example, if a persons' normal car fuel usage is five liters per day but if it becomes fifty liters in any random day, then it is a point anomaly.

- **Contextual anomaly:** When a data instance behaves anomalously in a particular context, it is termed a contextual or conditional anomaly; for example, expenditure on a credit card during a festive period, e.g., Christmas or New Year, is usually higher

than during the rest of the year. Although it can be high, it may not be anomalous as high expenses are contextually normal in nature. On the other hand, an equally high expenditure during a non-festive month could be deemed a contextual anomaly.

- **Collective anomaly:** When a collection of similar data instances behaves anomalously with respect to the entire dataset, the group of data instances is termed a collective anomaly. For example, in a human Electro Cardiogram (ECG) output, the existence of low values for a long period of time indicates an outlying phenomenon corresponding to an abnormal premature contraction whereas one low value by itself is not considered anomalous.

## Library used for NADI

Library used in different languages to implement network anomaly detection:

**Libraries used in python**

**TensorFlow:** It is an open source artificial intelligence library, using data flow graphs to build models. It allows developers to create large-scale neural networks with many layers. TensorFlow is mainly used for: Classification, Perception, Understanding, Discovering, Prediction and Creation.

**Pandas:** Pandas is a data manipulation library based on NumPy which provides many useful functions for accessing, indexing, merging, and grouping data easily. The main data structure (DataFrame) is close to what could be found in the R statistical package; that is, heterogeneous data tables with name indexing, time series operations, and auto-alignment of data.

**NumPy:** NumPy is a low-level library written in C (and Fortran) for high level mathematical functions. NumPy cleverly overcomes the problem of running slower algorithms on Python by using multidimensional arrays and functions that operate on arrays. Any algorithm can then be expressed as a function on arrays, allowing the algorithms to be run quickly.

NumPy is part of the SciPy project, and is released as a separate library so people who only need the basic requirements can use it without installing the rest of SciPy. NumPy is compatible with Python versions 2.4 through 2.7.2 and 3.1+.

## Dataset used for NADI

We have used NSL-KDD dataset. The NSL-KDD data set is a refined version of its predecessor KDD99 data set (which basically 21 years old). In this project the NSL-KDD data set is analyzed and used to study the effectiveness of the various classification algorithms in detecting the anomalies in the network traffic patterns.

Improvement to the KDD'99 dataset

The NSL-KDD data set has the following advantages over the original KDD data set:

- It does not include redundant records in the train set, so the classifiers will not be biased towards more frequent records.
- There are no duplicate records in the proposed test sets; therefore, the performance of the learners is not biased by the methods which have better detection rates on the frequent records.
- The number of selected records from each difficulty level group is inversely proportional to the percentage of records in the original KDD data set. As a result, the classification rates of distinct machine learning methods vary in a wider range, which makes it more efficient to have an accurate evaluation of different learning techniques.
- The number of records in the train and test sets are reasonable, which makes it affordable to run the experiments on the complete set without the need to randomly select a small portion. Consequently, evaluation results of different research works will be consistent and comparable.

# Workflow

### Task for 1st week

### Objective:

1. Learn and search about paper (IEEE transaction) on network anomaly detection and inspection technique.
2. To learn about packet capturing, for that we have used an application called **wireshark** to capture the packets going from our system to various servers. Common problems that Wireshark can help troubleshoot include dropped packets, latency issues, and malicious activity on network. It lets you put network traffic under a microscope, and provides tools to filter and drill down into that traffic, zooming in on the root cause of the problem.

### Paper used for study (IEEE Transaction paper):

1. A survey of network anomaly detection techniques Q1 Mohiuddin Ahmed, Abdun Naser Mahmood, Jiankun Hu
2. A comprehensive survey on network anomaly detection Gilberto Fernandes Jr.1 · Joel J. P. C. Rodrigues
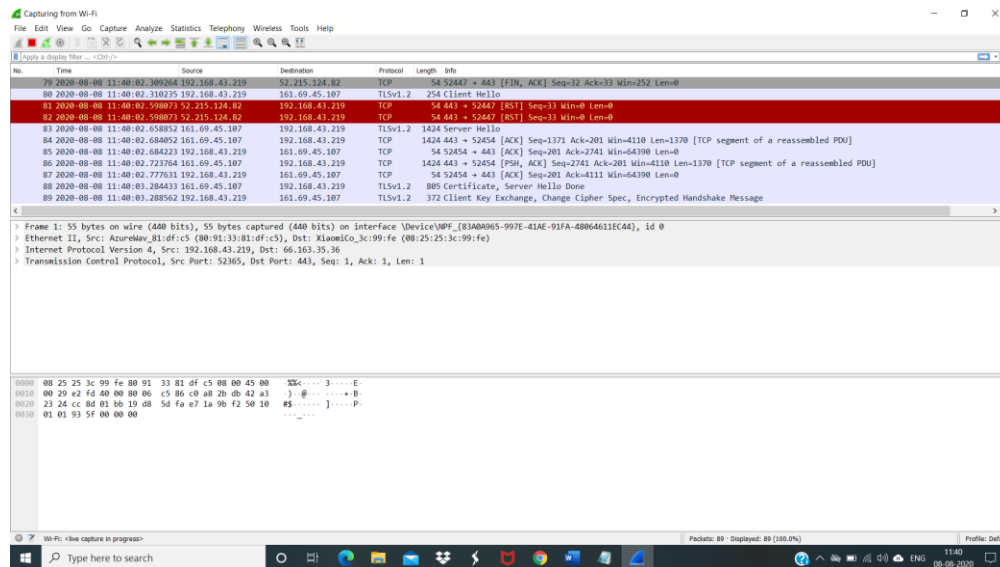3. A hybrid machine learning approach to network anomaly detection Taeshik Shon a,*, Jongsub Moon b

### Tools used:

- For packet capturing
  we have used an application called **Wireshark** to capture the packets going from our system to various servers.

### Observation:

1. We have observed various packets which are going through our system with the help of Wireshark tool in both windows and Linux platforms.

2. Each packet has information regarding its destination, source protocol, length, info attribute which basically tells about packet information. Here is snapshot of some packets.



**Task for 2<sup>nd</sup> Week**

**Objective:** Learn about the data Set And libraries used for Network anomaly detection and inspection.

**Library used for NADI**: TensorFlow, pandas, sklearn. numpy

**Application**: Anaconda Jupyter notebook.

**Observation:**

1. We have learned about pandas, TensorFlow library of python.
2. We have installed TensorFlow module in Anaconda Environment.
3. We Have learned about different function of these libraries, which are used for file read, data cleaning, model building and used for anomaly detection.
   Some import functions Are:

| read_csv() | Pandas function to read Csv file. |
|---|---|
| read_table() | Pandas function to read txt file as table |
| read_excel() | Pandas function to read excel file |
| Shape | Pandas function to read total number of columns |
| Isnull() | Pandas function to check NAN value |
| Isnull().sum() | Pandas function to find total NAN values |
| Groupby() | To count different group of attribute |
| Mean() | To find mean value |

| | |
|---|---|
| head() | To read first five rows of datafrome |
| tail() | To read last column of dataframe |
| Pd.Copy() | To copy one dataframe to other dataframe |
| Get_dummies | To create dummy variable |
| train_test_split() | To spilt data into training and testing data set |

## Task for 3rd and 4th week

**Objective:** to learn about the dataset and data cleaning

**Dataset used:** NSL-KDD

**Library used:** pandas

**Observation**

1. Dataset consist of various columns such protocol, and each row of contain information regarding type of attack

2. NSL-KDD dataset has no NAN values in each Attribute



3. We have observed that it consists of 23 attacks as show in image

- back            956
- buffer_overflow     30
- ftp_write           8
- guess_passwd        53
- imap            11
- ipsweep         3599
- land            18
- loadmodule          9
- multihop            7
- neptune         41214
- nmap            1493
- normal          67343
- perl            3
- phf             4
- pod             201
- portsweep        2931
- rootkit         10
- satan           3633
- smurf           2646
- spy              2
- teardrop        892
- warezclient         890
- warezmaster         20

## Task for 5th and 6th week

**Objective:** build a machine learning model for network anomaly detection.

**Dataset used:** NSL-KDD

**Technique:**  Neural Network, Convolutional Neural Network, Autoencoders, z_score

**Accuracy:** we have split the dataset into 75% training data And 25% testing dataset, then we check for accuracy and got 95% accuracy

```
Validation score: 0.950180986854639
```



9

## Task for 7th week

**Objective:** build a machine learning model to detect anomaly (abnormalities) using different method/ algorithm

**Dataset used:** NSL-KDD

**Technique used:** Random forest classifiers

**Accuracy:** this time we trained the machine learning model with Random forest classifier technique and we were able to detect the abnormalities with accuracy of 99%. At first, we have split the data in training and testing data into 75% and 25% data respectively. Then we feed the data to random forest classifier with n_estimator=20 i.e. number of decision tree for evaluation is 20 decision trees. Then we check for accuracy and we got an accuracy of 99%.

```
In [77]: #implementing random forest to detect anomaly abnormalities
         from sklearn.ensemble import RandomForestClassifier
         model1 = RandomForestClassifier(n_estimators=20)
         model1.fit(x_train, y_train)

Out[77]: RandomForestClassifier(n_estimators=20)

In [90]: model1.score(x_test, y_test)

Out[90]: 0.9976503460976694
```

**Summarization and conclusion:**

In this internship program we learn about network anomaly detection technique and about various others topic such as packet capturing and dataset like KDD99, NSL-KDD according to weeks as mention above of report. Then We have also built the machine learning model using simple neural network, and random forest classifiers to detect anomaly (abnormalities) with accuracy of 95% and 99% respectively. Now further we want to work with this project and modify the technique to learn about how to predict the attack type with anomaly detection. And in our project, the data set is spilt into training dataset with 75% of main dataset and testing dataset of 25%.and we have simple neural network and random forest techniques to analysis the NSL_KDD for anomaly detection. The analysis shows that NSL-KDD dataset is very ideal for comparing different intrusion detection models. In future we can try to improve the Random Forest algorithm and neural network algorithm to build an efficient intrusion detection system with prediction of attack type by identifying all 23 attack.

**Reference:**

1. A survey of network anomaly detection techniques Q1 Mohiuddin Ahmed, Abdun Naser Mahmood, Jiankun Hu
2. A comprehensive survey on network anomaly detection Gilberto Fernandes Jr.1 · Joel J. P. C. Rodrigues
3. A hybrid machine learning approach to network anomaly detection Taeshik Shon a,*, Jongsub Moon b
4. A Detailed Analysis on NSL-KDD Dataset Using Various Machine Learning Techniques for Intrusion Detection by S. Revathi and Dr. A. Malathi
5. Vipin Kumar, Himadri Chauhan, Dheeraj Panwar, "K-Means Clustering Approach to Analyze NSL-KDD Intrusion Detection Dataset", International Journal of Soft Computing and Engineering (IJSCE) ISSN: 2231-2307, Volume3, Issue-4, September 2013.