

Anomaly Intrusion Detection System Using (NAB) Machine Learning Technique

¹Sheela Patel, ²Shivendra Dubey, ³Mukesh Dixit

¹M. Tech Scholar, ²Research Guide, ³Head of Department

¹²³Department of Computer Science Engineering, REC, Bhopal

Abstract- The aim of this analysis is the creation and associate build up the system to forestall an organism against each well-known and new attacks, and functions as an adaptive distributed defense system or adaptive artificial system. Artificial Immune Systems abstract the structure of immune systems to include memory, fault detection and adaptive learning. We tend to propose associate system primarily based real time intrusion detection system exploitation supervised learning algorithmic rule. The formula is tested on the KDD 99 information, wherever it achieves an occasional warning rate whereas maintaining a high detection rate. This can be true even just in case of novel attacks, that might an important improvement over alternative algorithms.

KEYWORDS: ANN, KDD-99, Supervised Learning, AIS, NAB, Confidence Factor.

I INTRODUCTION

Now a day's development of any country or origination is depending upon its information technology system and all the information whether it's confidential, personal or public is shared through internet or network. So any country or organization needs to develop their information sharing network throughout the world with rapid speed. There is a rapid development in making such types of networks which available worldwide and have confidential information. But some time the intruder can attack over network where network based or client based firewall not capable enough to provide complete security against these types of threads.

Computer security is a very important issue to any or all users of computer systems. The rise of the web, computer attacks are increasing and may simply cause numerous dollar harm to a corporation. Detection of those attacks is a very important issue of pc security. Intrusion Detection Systems (IDS) technology is a good approach in addressing the issues of network security. The main objective of Intrusion Detection System is to observe unauthorized use, misuse and abuse of pc systems and laptop by each systems insiders and external intruders. There are many strategies following implement intrusion detection like statistical analysis knowledgeable systems, and state transition

approaches etc., and these many approaches is based on the system were planned in recent years. In order to provide complete security against these word wide thread IDS system play a key role. IDS system identifies the unauthorized activity that compromise the integrity, confidentiality and availability of confidential information.

Conventional IDS is based on continuous monitoring of well know attack by their extensive knowledge of signature to detect intrusion. This method based on pattern recognitions of various audit streams and detect intrusion by comparing their pattern provide by human expert. The pattern has been manually revised for a new type of intrusion whenever discover. The basic limitation of this pattern based Method is cannot detect emerging cyber thread.

Artificial Immune System is an emerging technology in order to fine the intruders or making the IDS. Recently AIS is a new bio-inspired model, which is applied for solving various security problems in the field of information security, genetic algorithms, neural networks, evolutionary algorithms and swarm intelligence. As one of the solutions to intrusion detection problems, AIS have shown their advantages. To improve the correlation factor and minimizing the false alarm generation we used the concept of AIS and NAB Algorithm to identify the intrusion in the system.

Detection and prevention of anomaly over the internet in real time scenario is a big challenge. The versatile feature and dynamic nature of anomaly attack emerges the issue of discovery and avoidance of attack. The abnormality assault is umbrella of different assault such DOS, Probe, U2R, R2L and numerous mixes of assaults. For the discovery of abnormality assault utilized firewall, interruption identification framework, antivirus and numerous more application programming are used [1,2,3]. The handling of identification is ease back because of vast number of interruption characteristic, now different creators utilized machine learning method and highlight diminishment prepare for the characterization and detection in intrusion detection system. Big data is the accumulation of extensive informational collections and it winds up plainly hard handling utilizing reasonable customary information

preparing applications or database administration instruments. The difficulties incorporate procuring, putting away, seeking, sharing, exchanging, breaking down and imagining.

The pattern of big data is due to the other useful information that can be derived from analysis of large set of related data, allowing correlations to be made to spot business trends, prevent diseases, determine quality of research, link legal citations, contest crime, and find current roadway traffic conditions. In the circumstance of the data explosion phenomenon, existing performance models for Map-Reduce are applicable for specific production workloads, but are to reveal the real capabilities of the processing system under heavy workloads that process tens of terabytes of data. While processing a query in big data, speed is a significant demand. The combination of suitable index for big data and current preprocessing technology will be a desirable solution when we encounter this kind of problems. Big Data is used in many real-world areas such as telecommunications, health care, pharmaceutical or financial businesses. Machine learning offers various algorithm for classification, clustering and combination of clustering and classification. The clustering techniques provide various algorithms such as k-means, k-mod, FCM and many more algorithms. Instead of clustering the classification algorithm gives more accuracy in terms of detection.

II SYSTEM METHODS

In this section we present the conclusion results of existing intrusion detection techniques for detection DOS attacks. Intrusion detection system in a very popular and computationally expensive task.

Network intrusion detection is an important component for network management and it is defense mechanism for network security. A real-time network intrusion detection system has been presented in this work. This proposed Support Vector Machine based network intrusion detection system is evaluated with KDD 99 dataset. The proposed system is developed with the consideration of big streaming data. The experimental results show that the proposed system is feasible for stream processing of network traffic data for detection of network intrusion with high accuracy [4].

Intrusion detection systems play an important role in network security. Feature selection is the major challenging issue in IDS in order to reduce the useless and redundant features among the attributes. In this report, a hybrid learning approach through combination of K - Means clustering and SVM classifier are proposed. In hybrid IDS we have used RBF kernel function of SVM for classification purpose. They took the help of K- Means clustering

technique to reduce large heterogeneous dataset to a number of small homogeneous subsets. The proposed approach is compared and evaluated using KDD CUP 99 dataset. Only selected attributes are used for cluster formation as well as for classification purpose. As a result complexities get reduced and consequently performances get increased. Simulation results prove that accuracy rate and detection rate of DOS, PROBE, U2R and R2L increases by using this proposed method. Furthermore, for reduced feature attribute dataset performance of these hybrid IDS slightly increases as compare to all 41 attribute set. As well as the false alarm rate also decreases of proposed technique [5].

They have presented an immune system inspired unsupervised intrusion detection system. Unlike other methods it assumes the coexistence of viruses and self, and therefore allows adaptability without supervision. It does not require labeled data, but learns the signatures of normal connections and intrusions from the dataset itself. It consists of two units: the T cell and the B-cell. The T-cells are formulated using a hidden Markov model, while additionally incorporating information from past observations to improve the model adaptively.

The B-cells combine inputs from T-cells with broader-level feature information to weed out false alarms. Its two – layer structure, similar to biological immune systems, achieves a high detection rate and helps keep false alarm rate in check; while other models show either very low detection rates or unacceptably high false alarm rates. The algorithm can operate in real time. Since viruses mutate, the ability of an IDS to detect viruses previously not encountered is of paramount importance. Our model is able to do that, and represents a significant improvement over other models in this respect. This is evident from results for the ipsweep dataset, where our model successfully detected a type of attack it had not been exposed to in the training stage [6].

III PROPOSED METHODOLOGY

The planned design contains varied modules every outlined with a particular purpose and connected along to spot the precise unauthorized person within the given system.

Intruder Data: it's a knowledge set accessible online so as to perform analysis work. It's a data on that the planned formula can work. In our planned work use KDD-99 cup dataset.

1 T-Cell theory over unauthorized activity information: This a part of the planned model takes the unauthorized activity data as associate input and apply the T-Cell formula and send the result to B-cell. Before causing it to completely different stages it'll found the traditional and abnormal feature.

2 B-cell Function: This step is employed to calculate the degree of belief of the chosen information set. It helps to gather the proof.

3 Determination of probability of an attack: during this a part of our work the classification and optimisation has performed victimisation support vector machine.

The detail rationalization of operating steps concerned in planned methodology:

Step1: With the assistance of nerve fibre Cell formula we tend to classified information, whether the info is traditional or affected with anomaly or we are able to say, abnormal.

The formula operates in two steps: Firstly it identifies whether or not anomalies occurred within the past supported the computer file, Secondly it correlates the known anomalies with the potential causes, generating associate anomaly scene per suspect. After applying the formula, we tend to reason the info into five classes specifically traditional, Denial of service attack (DOS), user a pair of root attack, remote a pair of native attack and inquiring.

Step2: B-cell Function theory is employed to figure the chance of evidences that indicate support the attack or traditional category.

The use of B-cell Function steady spreads out, principally as a result of its wont to address massive amounts of uncertainties that are inherent of natural environments. This new approach considers sets of propositions and assigns to every of them an interval [Belief, Plausibility].

Step3. After the classification we calculate the entropy of the attack treated as signal. For the calculation of entropy let us consider set having possible event. Each of which we assumed, occurs some numbers of times. Thus if there are n distinct possible event X_1, X_2, \dots, X_n , and the event occurred at frequency N_1, N_2, \dots, N_n .

Algorithm for Detection Phase:

The algorithm which is using in our proposed methodology for detect the anomaly behavior of test data set discussed below. (Numenta Anomaly Benchmark [NAB]) algorithm used and elaborate in the following points.

Step 1: Start

- Select the test excel data sheet (KDD-Cup-99 Dataset)
- Define the variable loop value to read the data sheet length
- Format test data into string
- Create new two arrays (final_res and final_act) of Zeros
- Read the test excel data sheet
- Step 2:
 - Load trained data set for matching
 - Define main category
 - Taking particular case i.e. case-14, case-23, case-otherwise

Pri_col = 14; (For case = 14)
Pri_col = 23; (For case = 23)
Pri_col = 44; (For case = otherwise)

Step 3: Define some variables,

Wrong = []; cmp = []; actual = []; search = [1: Upto]; wn = []; wa = []

Step 4: Length Searching

for l = 1:length(search)
 r_no = search(l)
 read_data =

data1(r_no,1:42);

Step 5: For desire category finding,

Comparing the two string [Test data (d) and trained match data]

Step 6: Finding Confidence Factor

If decision ==High, then (it is case of anomaly detection)

If decision ==Low, then (it is case of normal detection)

Step 7: Final Matching,

If match (length) then final_act == final_res;
If final_act != final_res, then wrong predicted;

Step 8: For Normal, the value of final_res = final_act = 1

If value of final_act == final_res =2, (Detected as DOS)

If value of final_act == final_res =3, (Detected as R2L_U2L)

If value of final_act == final_res =4, (Detected as Probe)

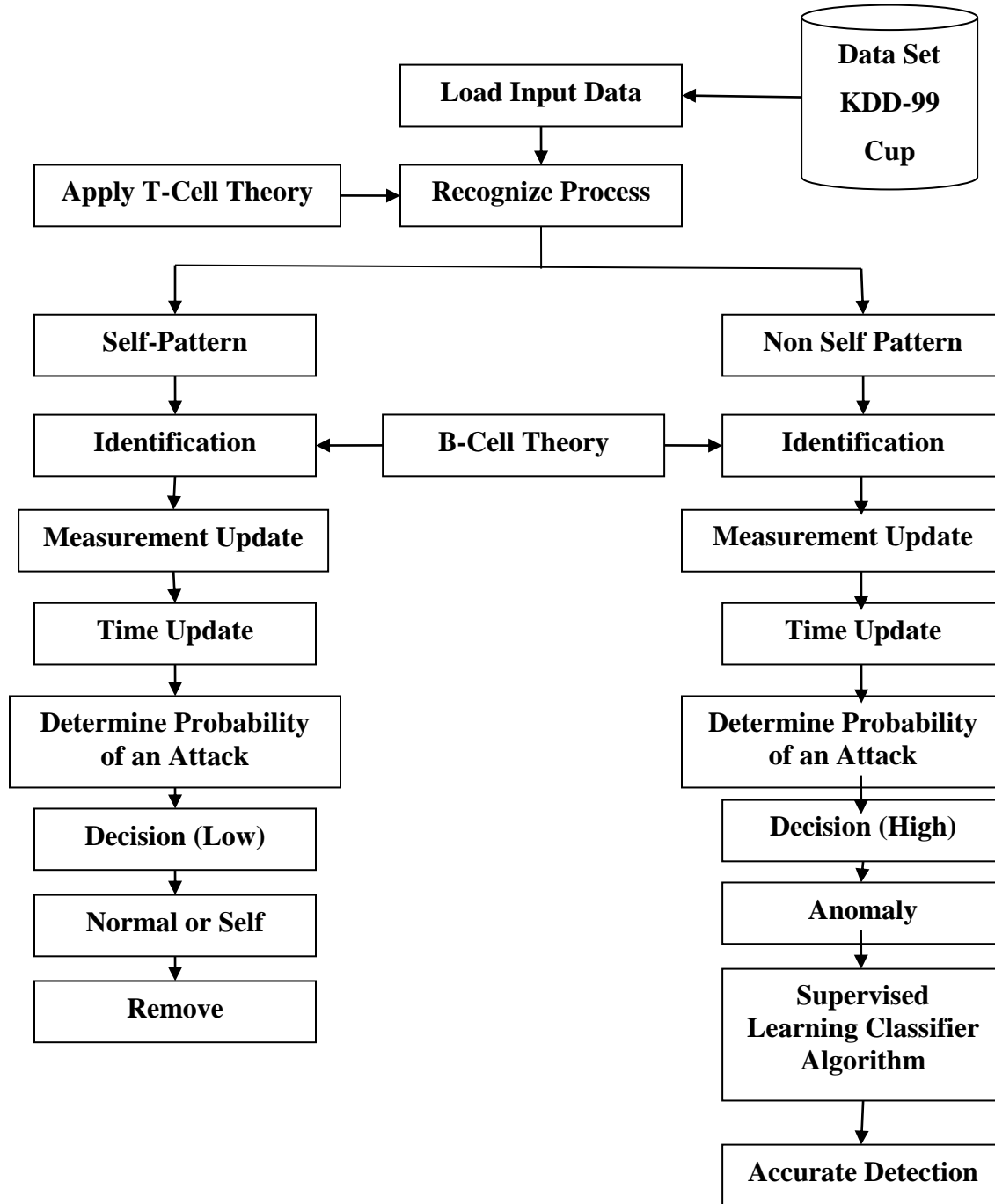


Figure 1 Model Flow Design Chart of proposed Methodology

Figure 1 shows the design for the planned new methodology for intrusion detection that's supported one in all the formula of artificial system referred to as the T-Cell and B-Cell theory. The nerve fibre T-Cell formula facilitate us to resolve the matter of correlation and B-Cell theory resolve the matter of unknown and known evolving harmful attacks and also describe the process which are involving in the

process of NAB supervised classifier and confidence factor with the help of decision making output (low or high), that also define the accurate detection of unauthorized activities.

IV SIMULATION RESULTS

The data set employed in the experiments is "KDD Cup 1999 Data", that may be a subversion of Defense Advanced (Defense

Advanced Research Projects Agency) 1998 dataset. The KDD cup ninety nine dataset Includes a group of forty one options derived from every association and a label that specifies the standing of association records as either traditional or specific attack sort. These options had all types of continuous and symbolic with extensively variable ranges falling in four categories:

A variety of attacks incorporated within the dataset make up following four major classes:

Denial of Service Attacks: A denial of service attack is associate attack wherever the assailant constructs some computing or memory resource absolutely occupied or unobtainable to manage legitimate necessities, or reject legitimate users right to use a machine.

User to Root Attacks: User to Root exploits are a class of exploits wherever the assailant initiate by accessing a standard user account on the system (possibly achieved by trailing down the passwords, a dictionary attack, or social engineering) and profit of some condition to realize root access to the system.

Remote to User Attacks: a far off to User attack takes place once associate assailant who has the potential to send packets to a machine over a network however doesn't have associate degree account on it machine, makes use of some vulnerability to realize native access as a user of that machine.

Probes: searching may be a class of attacks wherever associate degree assailant examines a network to gather info or discover well-known vulnerabilities. These network investigations square measure moderately valuable for associate assailant who is staging an attack in future. Associate assailant who has a record, of that machines and services are accessible on a given network, will build use of this info to appear for fragile points.

Figure 2, shows the main window of proposed IDS system. In this main window of proposed IDS system has been implemented in MATLAB 2009b framework.

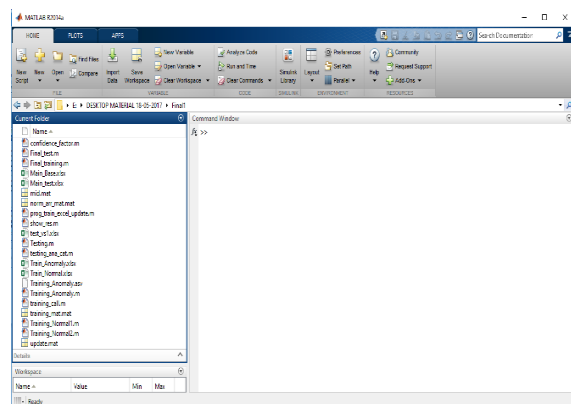


Figure 2 Main MATLAB Window Environments

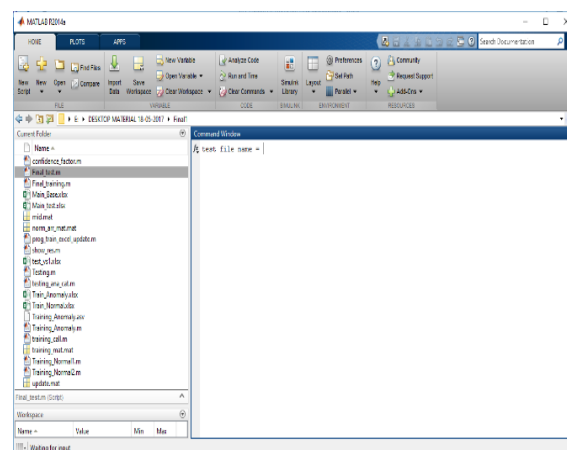


Figure 3 Training Phase of Proposed Algorithm in MATLAB Environment

Figure 3 show the training phase of proposed algorithm in MATLAB environment. The training phase of the test file follows some tasks. The sequence of steps is as follows.

1. Start the software module environment.
2. Select the data set which has to be trained by the training module from the main data base folder.
3. Fetch each and every attribute class and count. These three points include the training phase of our projected module.
4. Then construct the new data for new collected information.
5. Last but not the least of our project method module is updating this information into the data set and gets trained data set.

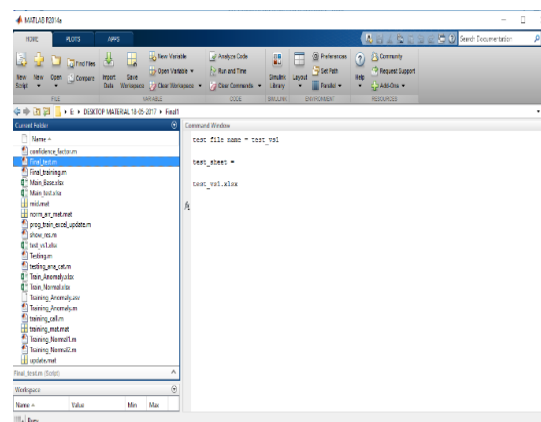


Figure 4 Testing Phase of Proposed Methodology for Test File

Figure 4 shows that the testing phase of test data files of proposed methodology. In this process basically find the abnormalities about the test data set. These finding basically distinguished between normal and abnormal behavior of our data set. The sequences that our proposed algorithm follows: The searching according to length and location, Matching process, and determine the confidence factor and then find out the detection rate and accuracy of the system.

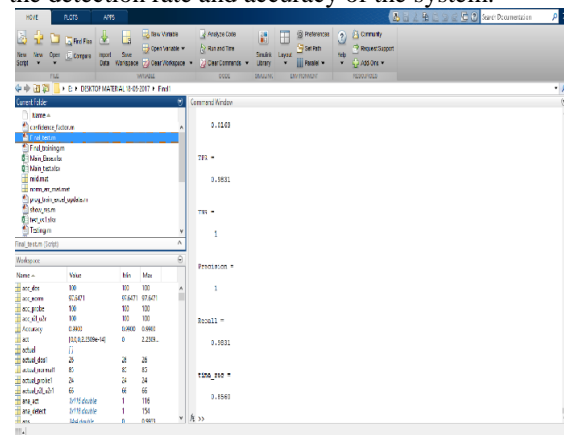


Figure 5 Results of Our Proposed Methodology

Figure 5 show that the results of our proposed methodology. In this figure shows the accuracy results for all attributes class data set.

V COMPARISON WITH EXISTING MECHANISM (ANOMALY)

This section shows that the comparison in term of accuracy and detection rate with the existing mechanism. In this section as existing mechanism we shows the results of the system using numerous algorithm like K-means clustering, support vector machine and k-means with RBF kernels of SVM.

5.1 Accuracy Results (In Percentage) for All Attribute Set

METH OD	KMSV M	KM	SV M	(PROPOSE D)
DOS	93.33	86.67	40	98.50
PROBE	100	87.50	75	100
U2L	93.75	68.75	62.5 0	93.78
R2L	87.50	75	68.7 5	90.13
ALL	80.28	73.24	49.3 0	90.07

Table 1 Accuracy Results (In Percentage) for All Attribute Set

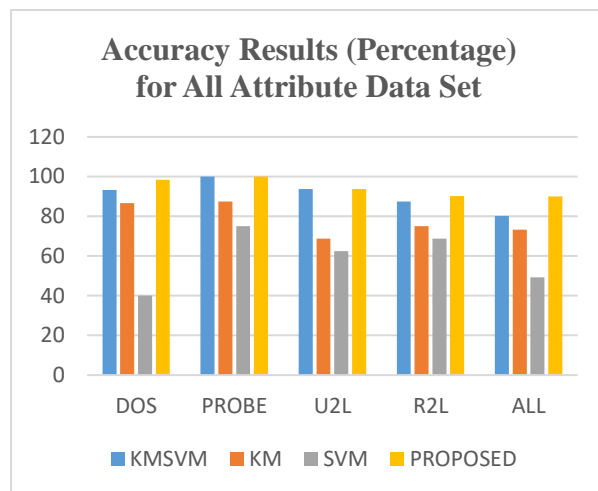


Figure 6 Accuracy Results (In Percentage) for All Attribute Data Set

In table 1 shows the accuracy (percentage) value of different algorithm used for IDS, in this table we describe the algorithm used system accuracy k-means support vector machine (KMSVM), K-Means Clustering (KM) and Support Vector Machine (SVM) and proposed hybrid model. Also we plot these value of accuracy of different algorithm used IDS system for all attribute class data set in figure 6.

5.2 Accuracy Results (In Percentage) for Reduced Attribute Set

METHOD DATA- SET	KMSV M	KM	SVM	(PROPOSED)
DOS	100	100	57.9 8	100
PROBE	91.30	91.3 0	73.9 1	100
U2L	100	77.7 8	88.8 9	100
R2L	94.12	82.3 5	70.5 9	98.50
ALL	87.01	76.6 2	62.3 4	93.34

Table2 Accuracy Results (In Percentage) for Reduced Attribute Set

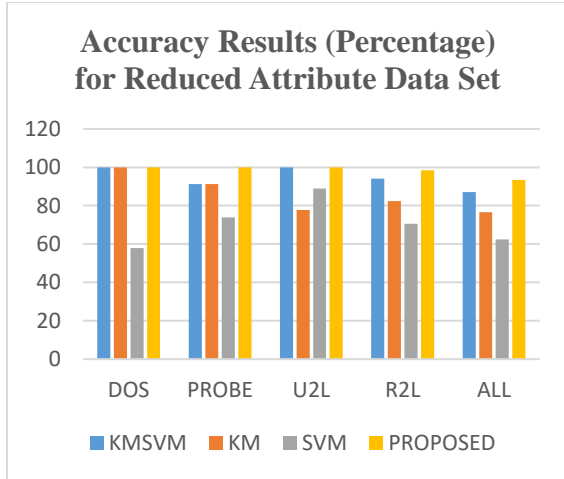


Figure 7 Accuracy Results (In Percentage) for Reduced Attribute Data Set

In table 2 shows the accuracy (percentage) value of different algorithm used for IDS, in this table we describe the algorithm used system accuracy k-means support vector machine (KMSVM), K-Means Clustering (KM) and Support Vector Machine (SVM) and proposed hybrid model. Also we plot these value of accuracy of different algorithm used IDS system for all attribute class data set in figure 7. Also Figure 6 and Figure 7 shows comparison of the simulation result. It gives the comparison of the degree of Accuracy rate of IDS system by using traditional classification method namely KMSVM, SVM and KM with our proposed method Hybrid model. Hybrid modal increases the accuracy rate by encapsulating B and T-Cell theory with NAB algorithm. As shows in figure 6 KMSVM, KM and SVM classification algorithm alone having accuracy rate for all attributes data set never reaches even 93.75% whereas hybrid model having accuracy rate up to 94.50%. As shows in figure 7, KMSVM, KM and SVM classification algorithm alone having accuracy rate for reduced attributes data set reaches even average accuracy 97% whereas hybrid model having accuracy rate up to 98.50%.

5.3 Detection Rate Results (In Percentage) for All Attribute Set

METHOD \ DATA-SET	KMSVM	KM	SVM	(PROPOSED)
DOS	91.67	85.71	100	99.30
PROBE	100	89.47	75	98.70
U2L	90.91	66.67	62	95.80
R2L	84.62	73.33	100	100
ALL	79.37	72.86	100	99.80

Table 3 Detection Rate Results (In Percentage) for All Attribute Set

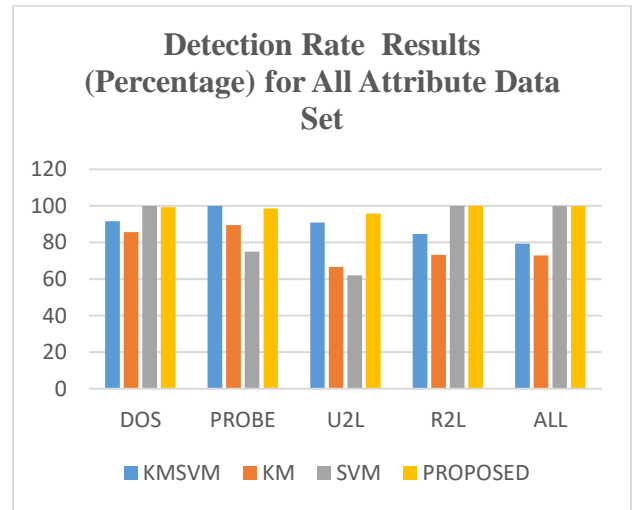


Figure 8 Detection Results (In Percentage) for All Attribute Data

In table 3 shows the detection rate (percentage) value of different algorithm used for IDS, in this table we describe the algorithm used system accuracy k-means support vector machine (KMSVM), K-Means Clustering (KM) and Support Vector Machine (SVM) and proposed hybrid model. Also we plot these value of detection rate of different algorithm used IDS system for all attribute class data set in figure 8.

VI CONCLUSION

In order to overcome all these deficiency from IDS, system over network, we propose a novel dual detection of IDS based on AIS that integrating the NAB algorithm. The training phase helps us to solve the problem of correlation and NAB theory resolves the problem of unknown and rapidly evolving

harmful attacks. The simulation results shows that the proposed method has improved the accuracy rates, minimizing false +ve and false -ve alarm generation and to increase the efficiency and accuracy of the IDS system.

REFERENCES

- 1) Baojiang Cui and Shanshan He “Anomaly detection model based on Hadoop platform and Weka interface”, Innovative Mobile and Internet Services in Ubiquitous Computing, 2016, Pp 84-89.
- 2) Sherenaz Al-Haj Baddar, Alessio Merlo and Mauro Migliardi “Anomaly Detection in Computer Networks: A State-of-the-Art Review”, Wireless Mobile Networks, Ubiquitous Computing, and Dependable Applications, 2015, Pp 29-64.
- 3) Zhengbing Hu, SergiyGnatyuk, Oksana Koval, Viktor Gnatyuk and SerhiiBondarovets “Anomaly Detection System in Secure Cloud Computing Environment”, I. J. Computer Network and Information Security, 2017, Pp 10-21.
- 4) Muhammad Asif Manzoor, Yasser Morgan, “Real-time Support Vector Machine Based Network Intrusion Detection System Using Apache Storm”, IEEE 2016.
- 5) Prof. UjwalaRavale, Prof. NileshMarathe, Prof. Puja Padiya, “Feature Selection Based Hybrid Anomaly Intrusion Detection System Using K Means and RBF Kernel Function”, Elsevier 2015.
- 6) ManjariJha, Raj Acharya, “An Immune inspired Unsupervised Intrusion Detection System for Detection of Novel Attacks”, IEEE 2016.