# An Optimal Approach of Classifying in Machine Learning happening in Cyber Security

[1] Dr R B Kulkarni
[1] Associate Professoe,IT Department, Government College of Engineering,Karad

*Abstract*— **The Digital flexibility is a quickly developing perspective that is accomplishing acknowledgment. Negative Cyber-assaults are those that oppositely impact the accessibility, uprightness or secrecy of IT arrange frameworks and related administrations and data. Earlier research works have carried on information control by an adversary as a worry, yet their works neglected to sum up the experiments. Many focused on formulating assault vectors inverse to explicit AI calculations and applications, for example, the Support Vector Machine (SVM) classifier. In our proposed work, an autonomous methodology on flexibility assessment and the development of enemy versatile classifiers utilizing Cluster Tree Map (CTM) Algorithm is finished. All information types in the area of Cyber Network information investigation are focused. The goal is to make a familiarity with any such strategy able to do effectively demonstrating the innovativeness and expertise of digital assailants and in this manner creating solo learning model. Better expected precision is achieved by utilizing Scalable Resilience Machine Learning Classifiers (SR-MLC).**

**Keywords: Resilience, Cluster, Classifiers, Cyber attackers Resilience, Cluster, Cloud, Cyber Security, Artificial Intelligence, Machine Learning, Network, Network Security, Analytics.**

## 1. INTRODUCTION

AI, as indicated by its, is a field of programming building that created from looking at structure affirmation and computational learning speculation in man-made brainpower. It is the learning and working of calculations that can pick up from and make desires on educational lists. These frameworks work by advancement of a model from point of reference commitments to demand to choose data driven estimates or choices instead of sticking to firm static program guidelines. [1] Machine learning techniques have been associated in various zones of science due to their uncommon properties like adaptability, flexibility, and potential to rapidly adjust to new and darken challenges. Digital security is a rapidly creating field mentioning a great deal of thought by virtue of amazing advances in informal organizations, web advances, cloud and portable condition, web based banking, shrewd matrix, etc. Different AI techniques have been adequately made to address such expansive issues in digital security. With the fast evolvement of web and adaptable headways, assault strategies are moreover winding up progressively refined in a few systems and staying away from real mark based techniques. AI techniques offer potential arrangements that can be used for settling such troublesome and complex conditions as a result of their ability to change quickly to new and cloud conditions. Distinctive AI strategies have been viably tended to unlimited issues in PC and information security.[2]

AI calculations are used to clarify a normally growing extent of collection issues in acknowledgment, vision, examination and inducing over the entire scope of processing stages [3]. AI calculations work in two phases: preparing and testing. In preparing, choice models are created subject to a stamped preparing enlightening file. In testing, the model is associated with order new info cases. [4] The use of AI calculations for digital security purposes offers climb to request of opposing quality, to be explicit: Can we assess the effort expected of an ominous to control a structure that relies upon AI frameworks? Could the adversarial adaptability of such structures be officially exhibited and evaluated? Would we have the option to quantify this quality with the ultimate objective that particular structures can be taken a gander at using empiric estimations? Past works have demonstrated how a horrible one can control a system subject to AI techniques by changing a part of its data sources.

Grouping is an astoundingly significant data exploratory AI instrument that empowers us to appreciate heterogeneous data by social occasion data with near qualities subject to a couple of criteria. Popular chart apportioning strategies, for instance, the Girvan–Newman

algorithm[5], sparsest-cuts[6], ghastly partitioning[7], and general conductance based systems (related to phantom methods by methods for Cheeger's inequality)[8] may be viewed as dealing with an edge-based flexibility issue concerning a diagram while simultaneously yielding the segments coming about on account of the ejection of the basic edge set as the game plan of bunches.

In our proposed work, Cluster Tree Map (CTM) calculation using organized information from flexibility information is made to focus on all information types in the zone of Cyber Network data examination which thus is used to gathering and characterize the strength information in the effective manner. The explanation behind the proposed work is to make an unaided learning model by making a free strategy on the contribution of substance specialists to assess include control costs. A free accessibility on versatility assessment and the development of foe strong classifiers is performed using CTM Algorithm. Therefore a familiarity with any such strategy is made able to do effectively displaying the inventiveness and expertise of digital aggressors. The goals and our commitments in the proposed work are

☐     To make an independent method on the involvement of content experts to evaluate feature manipulation costs.

☐     To make an independent presence on resilience estimation and the construction of unfavourable resilient classifiers.

☐     To make numeric techniques depend on the opinions of experts might seem preferable.

☐     To make an awareness of any such method capable of correctly modelling the innovativeness and cyber attackers' skill.

☐     To utilize unsupervised learning model and all these are done using CTM algorithm.

## 2. REVIEW OF RELATED LITERATURES

Anna L. Buczak [9] depicted information mining (DM) and AI (ML) techniques for digital investigation identified with assault location. Some well known digital informational collections used in ML/DM Special accentuation was made on the usage of different ML and DM procedures in the digital area, both for recognition of abuse and peculiarity. Thus, it was difficult to make one proposal for every technique, in view of the sort of assault the framework should identify.

Nguyen et al. [10] clarified ML systems for arrangement of traffic in Network. The strategies clarified didn't rely upon effectively understood port numbers however on factual traffic highlights. The procedures are separated and checked on according to their alternatives of ML systems. The promising results of ML-based IP traffic order opened a few new roads for related research spaces, similar to the use of ML in interruption identifications, inconsistency location in client information and control, directing traffic, and building system profiles for proactive system continuous checking and the board.

Teodoro et al. [11] focused on irregularity interruption procedures in arrange. Measurable, information based, and AI approaches were exhibited, however their examination doesn't present a full arrangement of AI procedures.

Wu et al. [12] focused on Computational Intelligence procedures and their applications to aggressor location. Procedures to be specific Swarm Intelligence Artificial Neural Networks (ANNs), Evolutionary Computation, Fuzzy Systems, and Artificial Immune Systems are clarified in detail. Since just Computational Intelligence strategies are clarified, a few ML/DM techniques like choice trees, bunching, and rule mining have not been consolidated. Its attributes, similar to adjustment, high operational speed, adaptation to non-critical failure, and blunder strength in the region of boisterous data, fit the prerequisite of building a decent arrangement of distinguishing interruption.

Revathi and Malathi [13] focused on AI interruption strategies. The creators introduced an explanative arrangement of AI calculations on the NSL-KDD interruption location dataset, however their examination just fused an abuse discovery setting. Interestingly, this work clarified abuse discovery as well as peculiarity location. The results delineated that NSL-KDD dataset is a lot of perfect for examination of different sort of interruption identification models.

Buczak and Guven [14] focussed on AI methods and their use in interruption recognition. Calculations like Neural Networks, Genetic Algorithms, Support Vector Machine, Bayesian Networks, Fuzzy Logics, and Decision Tree were depicted in detail. Exactness, time for ordering an obscure occurrence with a prepared model, and unpredictability, understandability of the last arrangement

(characterization) of every ML or DM strategy gave better results.

Sahoo et al. [15] presented the conventional system of Malicious Detection of URL as an AI procedure and isolated and reviewed their commitments that address various elements of the issues like element portrayal and calculation structure. Nonetheless, they didn't clarify the specialized subtleties of the calculation.

Pervez and Farid [16] proposed a sifting calculation dependent on SVM classifier to choose different interruption characterization exhibitions on the NSL-KDD interruption identification dataset. The technique kept up the order precision of the SVM classifier yet it utilizes a diminished arrangement of information highlights from preparing information.

## 3. METHODOLOGY

Dynamic investigation report database is taken from the outset. Dynamic examination of malware fuses the executing procedure of malware, screens its qualities, and creates a profile.
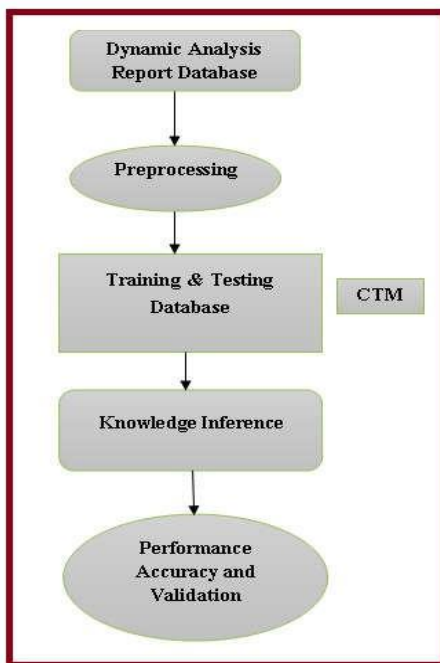


**Figure 1. Flow chart of Proposed Methodology**

It identifies the obscure malware by registering its comparative known malware profile [17]. After location, the following stage is pre-preparing pursued via preparing and testing the considered dataset. Testing process is conveyed by CTM (Cluster Tree Map) calculation. CTM is utilized to group and arrange the flexibility information in the powerful way and giving the better outcome. The following stage is Knowledge surmising. Induction is a database framework strategy used to assault databases in which malignant clients reason touchy data from complex databases at a higher stage. At last, execution investigation is done and the outcomes are approved. The stream outline of the proposed strategy is given in figure 1.

### 3.1 Attacks exploiting Machine Learning Systems

Situational mindfulness for Cyber Defenses consolidates 7 principal perspectives: "Checking a. current condition, impact of assault, how conditions create, entertainer lead, why and how the current situation is caused, quality and how possible fates of the current situation." Based on this definition, it might be said that situational mindfulness gives the client the closer view and brief evaluation of framework. Likewise, perceiving factors won't be anything but difficult to the point that even a couple of affiliations don't have mindfulness about their digital missions. Additionally, they have neither any sufficient after parts in their frameworks nor sensors which will lead them for assaults. Under these confinements, it is very difficult to have a ground-breaking situational mindfulness. Best procedure is perceiving the fundamental missions which the affiliation has; at that point executing the following instruments, in conclusion portraying proactive answers for the security of framework. A couple of assessments have been cultivated in this field. While a part of the examinations [18] [19] all things considered portray the term of situational mindfulness and standard systems and applications which have been used in order to take care of it; others revolve around increasingly express methods, for instance, continuous multistage assault mindfulness and mission-driven digital situational mindfulness [20], [21].

Assault trees expect a colossal procedure in showing structure security and framework to the extent frailty and danger conspicuous confirmation [22]. They can be mapped in various structures. For the most part, while hubs conquer assaults, the root hub is the worldwide target of the aggressor which can in like manner depicted

as an occasion. Youngster hubs are the improvements of this goal and branches are the manner by which aggressor can't be refined any more. Each way in assault tree indicates exceptional assault. Furthermore, assault trees can in like manner be mastermind abstractly as opposed to graphically. In literary structure, the 'AND' and additionally deteriorations are used and the results of achieve sub objectives were shown by them. Assault models can be described to fabricate the sensibility of assault trees age and reuse [23].

Assault configuration is the mapping of different sorts of assaults that fuses a. the target of the foreordained assault, b. the preconditions for use, c. the methods for practicing assault, d. post conditions which are substantial if the assault is made adequately [23]. The preconditions contain doubts which are associated with the typical acts of aggressor and the characteristics of the assault. The capacities, resources, access and data can be given for example to preconditions [23].

The foundation of disseminated reproduction and assault tree gives the plausibility of utilizing vulnerabilities originating from neighborhood weakness database. In assault tree, the representation depends on vulnerabilities. An assault tree hub incorporates a. the gadget b. helplessness which is the wellspring of an assault to this gadget (counting pre/post conditions), c. picked up benefit because of this assault on this gadget. A while later, by utilizing assault trees, sway scores of the vulnerabilities and criticality levels, hazard examination is made and chance scores are determined for every way in an assault tree. In the informative supplement, the example of assault tree can be seen. [24]

An explanative scientific categorization of different assaults whose goal is misusing AI frameworks are: (a) Causative assaults changing the preparation procedure; (b) Attacks on uprightness and accessibility, making bogus positives as a rupture into a framework; (c) Exploratory assaults abusing the predominant vulnerabilities; (d) Targeted assaults towards a specific info; (e) Indiscriminate assaults in which information sources won't work out. First type is a safeguard against exploratory assaults, where an aggressor can make an assessment dispersion that the student predicts ineffectively. For shielding against this assault, the safeguard can restrict the entrance to the preparation method and information, making it harder for an assailant

to utilize figuring out. Additionally, the more troublesome theory space is, the harder for an aggressor to derive the scholarly speculation. Also, a safeguard can constrain the input gave to an assailant with the goal that it gets harder to break into the framework. Second type is a protection against causative assaults, where an aggressor can advance both assessment and preparing disseminated assignments.

**3.2 Cluster Tree Mapping**
So as to think about the information and needing bunch information focuses to comprehend their aggregate conduct, bunching is one of the go-to methods. Clustering systems can amass characteristics into a couple of comparative fragments where information inside each gathering is like one another and particular crosswise over gatherings.

**Cluster Tree Map Algorithm**

```
FeaturesSelect (Dataset)
Begin
//load dataset
source ?  new DataSource(Dataset);
dataset ?  source.getDataSet();
//create FeaturesSelection object
filter ?  new FeaturesSelection();
//create evaluator and search algorithm objects
eval ?  new CfsSubsetEval();
search ?  new GreedyStepwise();
//set the algorithm to search backward
search.setSearchBackwards(true);
//set the filter to use the evaluator and search algorithm
filter.setEvaluator(eval);
filter.setSearch(search);
//specify the dataset
filter.setInputFormat(dataset);
//apply
newData ?  Filter.useFilter(dataset, filter);
//save
f?  new File(fp);
ff?  f.getName();
driveLetter ?  ff.split("\\")[0];
fn?  "attrs_"+driveLetter+".arff";
fns?  "E:/input/"+fn;
saver ?  new ArffSaver();
saver.setInstances(newData);
saver.setFile(new File(fns));
saver.writeBatch();
FilteredFeaturesDataset? fns;
End
```

**J48 Algorithm**
Classification is the process of building a model of classes from a set of records that contain class labels. Decision Tree Algorithm is to find out the way the attributes-vector behaves for a number of instances.

Also on the bases of the training instances the classes for the newly generated instances are being found. This algorithm generates the rules for the prediction of the target variable. With the help of J48 algorithm the critical distribution of the data is easily understandable

## 3.3 CTM_Test (Filtered Features Dataset)

```
Begin
datafile ?  readDataFile(fpp);
validation? null;
data ?  newInstances(datafile);
data.setClassIndex(data.numFeaturess() - 1);
// Do 10-split cross validation
split[][] ?  crossValidationSplit(data, 10);
// Separate split into training and testing arrays
trainingSplits[] ?  split[0];
testingSplits[] ?  split[1];
// Use a set of classifiers
ctm?  new       j48();
// Run for each model
// Collect every group of predictions for current model in a FastVector
predictions ?  newFastVector();
// For each training-testing split pair, train and test the classifier
for  i < trainingSplits.length
Begin
validation ?  classify(mlprbf, trainingSplits[i], testingSplits[i]);
predictions.appendElements(validation.predictions());
End
// Calculate overall accuracy of current classifier on all splits
accuracy ?  calculateAccuracy(predictions);
// Print current classifier's name and accuracy in a complicated, but nice-looking way.
validation.toClassDetailsString();
validation.toMatrixString();
validation.toSummaryString();
new GenerateROC(fpp);
visualizetree(fpp);
End
```

## 4. RESULTS AND DISCUSSION

In our proposed work, an independent method on the involvement of content experts to evaluate feature manipulation costs has been applied and followed by resilience estimation. Then, unfavourable resilient classifiers are constructed. Numeric techniques depending on the prefer ability of experts has been found to make them aware of any such method which can correctly model the innovativeness and cyber attackers' skill. This unsupervised learning model is done using Cluster Tree Map (CTM) and J48 Algorithm. The CTM algorithm is applied to know the behaviour of group data points and to select and filter the appropriate features. J48 Algorithm has been formulated for classification and to find the behaviour of attributes.

The true positive and false positive rate, precision, recall F-Measure, ROC area and class values of cluster 0 and cluster 1 are estimated. Form table 1, it is observed that

the weighted average of the true positive and false positive rate, precision, recall F-Measure, ROC are 0.998, 0.002, 0.998, 0.998, 0.998 and 0.999 respectively.

**Table 1. Weighted Average of calculated results**

| | TP rate | FP rate | Precisio-n | Recall | F-Measure | ROC area | Class |
|---|---|---|---|---|---|---|---|
| | 0.998 | 0.001 | 0.999 | 0.998 | 0.999 | 0.999 | cluster0 |
| | 0.999 | 0.002 | 0.998 | 0.999 | 0.998 | 0.999 | cluster1 |
| Weighted average | | | | | | | |
| | 0.998 | 0.002 | 0.998 | 0.998 | 0.998 | 0.999 | |

**Table 2. Proposed Calculation results**

| | |
|---|---|
| Correctly Classified Instances ( 5204 ) | 99.85% |
| Incorrectly Classified Instances ( 8 ) | 0.15% |
| Kappa statistic | 0.9969 |
| Mean absolute error | 0.0017 |
| Root mean squared error | 0.0357 |
| Relative absolute error | 0.33% |
| Root relative squared error | 7.12% |
| Total number of instances | 5212 |
| Algorithm Processing Time(Milliseconds): | 4061 |

From the calculation results of table 2, it is observed that the Correctly Classified Instances (5204) is 99.85 %, Incorrectly Classified Instances (8) is 0.15%, Kappa Statistics is 99.69%, Mean Absolute error is 0.17%, Root mean squared error is 3.57%, Relative absolute error is 0.33%, Root relative squared error is 7.12%. The total number of instances is about 5212. The time required for the algorithm to process is 4061 milliseconds.

In table 3, Sensitivity and specificity for both the existing and proposed system are calculated from the TP, TN, FP, and FN values. Sensitivity values of existing and proposed systems are 0.981 and 0.999 respectively. It is observed that sensitivity of proposed system is greater than the existing system. The obtained Specificity values for existing and proposed system are 0.983 and 0.998 respectively. Proposed system shows greater specificity than existing system.

**Table 3. Sensitivity and Specificity of Proposed and Existing System**

| | TP | TN | FP | FN | Sensitivity | Specificity |
|---|---|---|---|---|---|---|
| Existing | 2434 | 2675 | 54 | 47 | 0.981 | 0.983 |
| Proposed | 2736 | 2468 | 5 | 3 | 0.999 | 0.998 |

From the graph shown in figure 3, it is observed that proposed methodology shows better probability than the existing method.

ISSN (Online) 2394-2320

**International Journal of Engineering Research in Computer Science and Engineering (IJERCSE)**
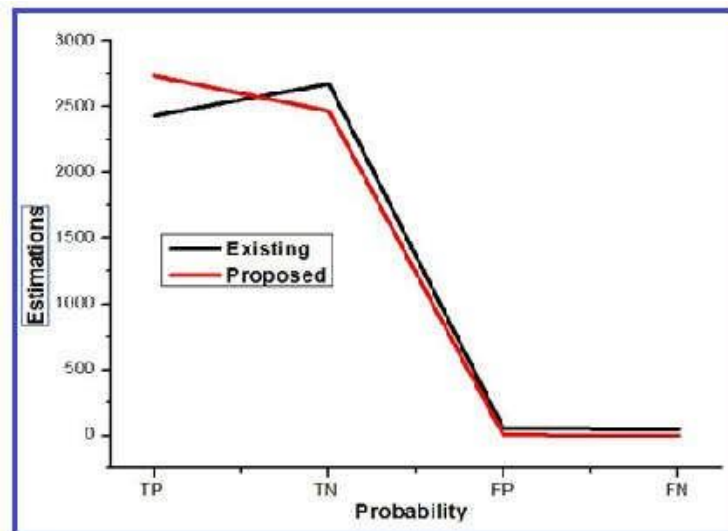**Vol 5, Issue 9, September 2018**

**Figure 3. Graph for probability comparison**

The parameters like TP rate, FP rate, precision, recall, F-Measure, ROC area, Class are estimated and compared with the existing method. The weighted average values of label 1 and label 2 for all the above mentioned parameters were found to be better compared to the existing methodology and it is shown in table 4.

**Table 4. Validation table**

| Param eters | TP Rate | | FP Rate | | Precision | | Recall | | F-Measure | | ROC Area | | Class | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Algorithms | Exi stin g | Pro pose d | Exi stin g | Pro pose d | Exi stin g | Pro pose d | Exi stin g | Pro po se d | Exi stin g | Pro pose d | Exi stin g | Pro pose d | Existi ng | Pro pose d |
| Label 1 | 0.97 8 | 0.99 8 | 0.01 7 | 0.00 1 | 0.98 1 | 0.99 9 | 0.97 8 | 0.99 8 | 0.98 | 0.99 9 | 0.98 4 | 0.99 9 | Benig n | clust er0 |
| Label 2 | 0.98 3 | 0.99 9 | 0.02 2 | 0.00 2 | 0.98 | 0.99 8 | 0.98 3 | 0.99 9 | 0.98 1 | 0.99 8 | 0.98 4 | 0.99 9 | Malig nant | clust er1 |
| Weight ed Avg. | 0.98 1 | 0.99 8 | 0.02 | 0.00 2 | 0.98 1 | 0.99 8 | 0.98 1 | 0.99 8 | 0.98 1 | 0.99 8 | 0.98 4 | 0.99 9 | | |

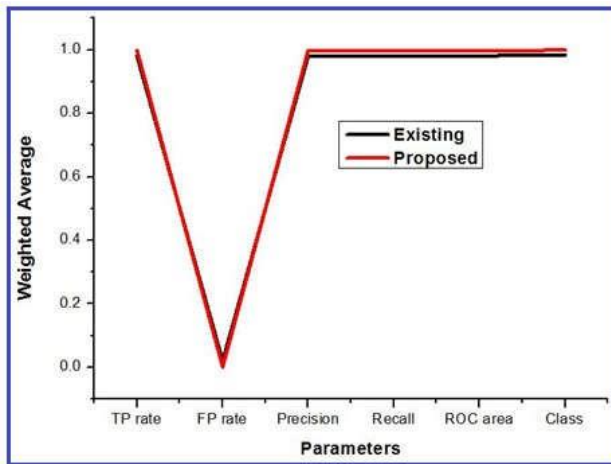The graph shown in figure 4 depicts the comparison of weighted average of both the existing and proposed methodology.

**Figure 4. Graph for parameters comparison**

**Table 5. Comparison of parameters**

| Parameters | | Existing | | Proposed |
|---|---|---|---|---|
| Correctly Classified Instances | 5109 | 0.980614 | 5204 | 0.99847 |
| Incorrectly Classified Instances | 101 | 0.019386 | 8 | 0.00154 |
| Kappa statistic | | 0.9611 | | 0.9969 |
| Mean absolute error | | 0.0264 | | 0.0017 |
| Root mean squared error | | 0.1354 | | 0.0357 |
| Relative absolute error | | 0.052693 | | 0.00331 |
| Root relative squared error | | 0.269474 | | 0.07116 |
| Total Number of Instances | | 5210 | | 5212 |
| Ignored Class Unknown Instances | 2 | | 0 | |
| Algorithm Processing Time(Milliseconds): | | 36926 | | 4061 |

From the table 5, it can be observed that correctly classified instances (0.99847), incorrectly classified instances (0.00154), Kappa statistic (0.9969), mean absolute error (0.0017), root mean squared error (0.0357), relative absolute error (0.00331), root relative squared error (0.07116), total number of instances (5212), ignored class unknown instances (0), algorithm processing time (4061) estimations are much better in proposed methodology than existing methodology.

## 5. CONCLUSION

CTM (Cluster Tree Map) has been used to cluster and classify the resilience data in the effective manner. The proposed methodology provided better results when all data types in the area of Cyber Network data analytics are focussed. Better accuracy is achieved by CTM algorithm.

Feature manipulation costs are estimated. To make an independent approach on the involvement of content experts to estimate feature manipulation costs using CTM Algorithm.Thus, an independent availability on resilience evaluation is made and adversary resilient classifiers are constructed. Unsupervised learning model is utilized to create an awareness of any such method capable of correctly modelling the creativeness and skill of cyber attackers. In future, Resilience of attacks can be done by advanced classifiers rather than Machine Learning Classifiers for achieving better results.

## REFERENCES

[1]     Simon, A., & Singh, M. (2015). An Overview of M Learning and its Ap. International Journal of Electrical Sciences Electrical Sciences & Engineering (IJESE), 22.

[2]     Ford, V., & Siraj, A. (2014, October). Applications of Machine Learning in Cyber Security. In Proceedings of the 27th International Conference on Computer Applications in Industry and Engineering.

[3]     P. Dubey. Recognition, mining and synthesis moves computers to the era of tera. Intel Tech. Magazine, 9(2):1–10, Feb. 2005

[4]     Venkataramani, S., Raghunathan, A., Liu, J., & Shoaib, M. (2015, June). Scalable-effort classifiers for energy-efficient machine learning. In Proceedings of the 52nd Annual Design Automation Conference (p. 67). ACM.

[5]     Girvan, M.; Newman, M.E. Community structure in social and biological networks. Proc. Natl. Acad. Sci. USA 2002, 99, 7821–7826. [CrossRef] [PubMed]

[6]     Shi, J.; Malik, J. Normalized cuts and image segmentation. IEEE Trans. Pattern Anal. Mach. Intell. 2000, 22, 888–905.

[7]     Alpert, C.J.; Kahng, A.B.; Yao, S.Z. Spectral partitioning with multiple eigenvectors. Discret. Appl. Math. 1999, 90, 3–26. [CrossRef]

[8]     Chung, F. Spectral Graph Theory; American Mathematical Society: Providence, RI, USA, 1997.

[9]     Buczak, A. L., & Guven, E. (2016). A survey of data mining and machine learning methods for cyber security intrusion detection. IEEE Communications Surveys & Tutorials, 18(2), 1153-1176.

[10]     T. T. T. Nguyen and G. Armitage, "A survey of techniques for internet traffic classification using machine learning," IEEE Commun. Surv. Tuts., vol. 10, no. 4, pp. 56–76, Fourth Quart. 2008.

[11]     P. Garcia-Teodoro, J. Diaz-Verdejo, G. Maciá-Fernández, and E. Vázquez, "Anomaly-based network intrusion detection: Techniques, systems and challenges," Comput. Secur., vol. 28, no. 1, pp. 18–28, 2009.

[12]     S. X. Wu and W. Banzhaf, "The use of computational intelligence in intrusion detection systems: A review," Appl. Soft Comput., vol. 10, no. 1, pp. 1–35, 2010.

[13]     S. Revathi and A. Malathi, ''A detailed analysis on NSL-KDD dataset using various machine learning techniques for intrusion detection,'' in Proc. Int. J. Eng. Res. Technol., 2013, pp. 1848–1853.

[14]     L. Buczak and E. Guven, ''A survey of data mining and machine learning methods for cyber security intrusion detection,'' IEEE Commun. Surveys Tuts., vol. 18, no. 2, pp. 1153– 1176, 2nd Quart., 2016.

[15]     D. Sahoo, C. Liu, and S. C. H. Hoi. (2017). ''Malicious URL detection using machine learning: A survey.'' [Online]. Available: https://arxiv.org/abs/1701.07179

[16]     M. S. Pervez and D. M. Farid, ''Feature selection and intrusion classification in NSL-KDD CUP 99 dataset employing SVMs,'' in Proc. 8th Int. Conf. Softw., Knowl., Inf. Manage. Appl. (SKIMA), 2014, pp. 1–6.

[17]     Moshiri, E., Abdullah, A. B., Mahmood, R. A. B. R., & Muda, Z. (2017). Malware Classification Framework for Dynamic Analysis using Information Theory. Indian Journal of Science and Technology, 10(21).

[18]     G. P. Tadda, J. S. Salerno, "Overview of Cyber Situation Awareness," in Cyber Situational Awareness Issues and Research, vol. 46, Springer, 2010, pp. 15-35.

[19]     L. D. Cumiford, "Situation Awareness for Cyber Defense," Information for the Defense Community, 2006.

[20]     S. Mathew, D. Britt, R. Giomundo, and S. Upadhyaya, "Real-Time Multistage Attack Awareness Through Enhanced Intrusion Alert Clustering," in Proc. Military Communications Conference, 2005, pp. 1801-1806.

[21]     S. Jajodia, S. Noel, P. Kalapa, and M. Albanese, "Cauldron missioncentric cyber situational awareness with defense in depth," in Proc. Military Communications Conference, 2011, pp. 1339-1344.

[22]     I. Ray and N. Poolsapassit, ''Using Attack Trees to Identify Malicious Attacks from Authorized Insiders,'' in Proc. Computer Security - ESORICS 2005 Lecture Notes in Computer Science, v. 3679, 2005, pp. 231-246.

[23]     A. P. Moore, R. J. Ellison, and R. C. Linger, "Attack Modeling for Information Security and Survivability," Carnegie Mellon Software Engineering Institute, 2001.

[24]     Goodall, J. R. (2008). Introduction to visualization for computer security. In VizSEC 2007 (pp. 1-17). Springer, Berlin, Heidelberg.

[25]     M. Barreno et al., "The security of machine learning", Journal Machine Learning, Vol. 81, Issue 2, pp. 121- 148, November 2010.