# A Review of Intrusion Detection System using Machine Learning Approach

**[1]SH Kok, [2]Azween Abdullah, [3]NZ Jhanjhi, [4]Mahadevan Supramaniam**

*[1,2,3]School of Computer and IT (SoCIT), Taylor's University, Malaysia.*

*[4]Research & Innovation Management Centre, SEGI University, Malaysia.*

*ORCIDs: [1](0000-0001-9477-8988), [2](0000-0003-4425-8604), [3](0000-0001-8116-4733), [4](0000-0002-3734-0899)*

## Abstract

Intrusion Detection System (IDS) is an important tool use in cyber security to monitor and determine intrusion attacks This study aims to analyse recent researches in IDS using Machine Learning (ML) approach; with specific interest in dataset, ML algorithms and metric. Dataset selection is very important to ensure model build is suitable for IDS use. In addition, dataset structure can affect effectiveness of ML algorithm. Thus, ML algorithm selection is dependent on the structure of the selected dataset. After that, metric will provide a quantitative evaluation of ML algorithms towards specific dataset. This study found that soft computing techniques are getting considerable attention, as many have applied it here. In addition, many researchers are focusing on the classification of IDS, which is beneficial in determining known intrusion attacks. However, it may pose a problem in detecting anomalous intrusion, which may include new or modified intrusion attacks. For dataset, many researchers were still using KDDCup99 and its variant NSL-KDD, although they are almost 20 years old. This continuous trend could result in static progress in IDS, while intrusion attacks continue to evolve together with new technologies and user behaviours. Ultimately, this situation will result in the obsolete use of IDS as part of a cyber security tool. Three most used metrices for performance evaluation for IDS are accuracy, True Positive Rate (TPR) and False Positive Rate (FPR). This is expected, because these metrices provide important indications that are very relevant to IDS functionality.

**Keywords** - Computation Intelligence, Dataset, Intrusion Detection System, Machine Learning, Soft Computing.

## I. INTRODUCTION

Intrusion Detection System (IDS) is an important tool use in cyber security to monitor and determine intrusion attack. There are three types of IDS; network IDS, host IDS, and Application IDS. Network IDS monitors network packet to detect intrusion attack. While host IDS monitors a single host (server or computer). Lastly, application IDS monitors several known high risk applications.

To determine whether an intrusion attack has occurred or not, IDS depends on few approaches. First is signature-based approach, where known intrusion attack signature is stored in the IDS database to match with current system data. When the IDS finds a match, it will recognise it as an intrusion. This approach provides a fast and accurate detection. However, the drawback of this is to have periodic update of the signature database. In addition, the system could be compromised before the newest intrusion attack can be updated.

The second approach is anomaly-based, or behaviour-based, where IDS will determines an attack when the system operates out of the norm. This approach can detect both known and unknown attacks. However, the drawback of this approach is low accuracy with high false alarm rate.

Lastly, hybrid-based approach uses both signature-based and anomaly-based approaches. This approach uses signature-based approach to detect known attacks, and anomaly-based approach to detect unknown attacks. Combining both approaches can ensure a more effective detection, but may increase computational cost.

Machine Learning (ML) uses statistical modeling approach to learn past data pattern, and then predicts the most likely outcome using new data. Therefore, ML algorithm has been applied to IDS using anomaly-based approach. As stated above, the challenge here is to build a model that can give high accuracy with low false alarm rate.

Therefore, this study aims to analyse recent researches in IDS using ML approach; with specific interest in dataset, ML algorithms and metric. Dataset selection is very important to ensure model build is suitable for IDS use. In addition, dataset structure can affect effectiveness of ML algorithm. Thus, ML algorithm selection is dependent on the structure of the selected dataset. After that, metric will provide a quantitative evaluation of ML algorithms towards specific dataset.

## II. APPROACH

In order to ensure we review researches of interest only, we preset some important criteria. Firstly, the article must be published in year 2015 and later. This is to ensure we get only the most recent researches, so that our study is relevant and not outdated.

Secondly, the article must be published in scientific journal or conference. This is to ensure the validity of the content, which have been peer reviewed and approved.

Thirdly, the article must use ML for IDS. This is our objective for this study, so we must work within the scope of our study.

## III. MACHINE LEARNING

ML algorithm can be categorized into 11 categories. This is shown in Fig. 1. Bayesian category uses Bayes Theorem of probability, which determines the probability of specific outcome to come true. The most popular algorithm in this category is Naïve Bayes.

Decision tree has a tree like structure that starts from root nodes, which is the best predictor. Then progresses through its branches until it reach a leave node. This is the decision outcome.

**Fig. 1** Category of ML algorithms adopted [1]

Dimensional reduction is to find features that are important to the outcome. This will removes irrelevant and redundant features. It is mostly performed during the pre-processing phase. The most popular algorithm is Principal Component Analysis (PCA)

Instance-based is also known as memory-based learning. This category of algorithm finds the most similar instances, or training data, that matches the new data to make prediction. The most popular algorithm in this category is k-Nearest Neighbour (kNN).

Clustering is grouping of data points that are close together to form its own group. This category of algorithm works well in unsupervised learning approach, which do not require labelled data. The most popular algorithm in this category is k-Means.

Regression algorithm try to build model that can represent the relationship between variables. It is derived from statistical analysis. The most popular algorithm in this category is Logistic Regression.

Neural network is inspired by the brain cell called neuron that forms the biological neural network. This category finds patterns from the data to make its prediction. Normally it would require large amount of data to produce a good prediction. The most popular algorithm in this category is Perceptron.

Ensemble is a method of combining the result of several algorithms before producing the final outcome. There are typically 2 methods, bagging and boosting.

Table 1 is the summary of researches article found in this study. Information extracted and summarized in this table are dataset, method (or algorithm) and accuracy metric being used in their researches.

**Table 1** List of recent researches in IDS from 2015 to 2018

| Reference | Dataset | Method | Accuracy (%) |
|---|---|---|---|
| [2] | TRAbID (Probe, DoS) | Decision Tree (DT) and Naïve Bayes (NB) | Probe; DT (98.42), NB (97.29) DoS; DT (99.90), NB (99.66) |
| [3] | CIDDS-001 | k-Nearest Neighbour (kNN) and k-means | kNN (99.0), k-means (99.7) |
| [4] | ISCX 2012 | Recursive Feature Addition (RFA) with SVM | 92.90 |
| [5] | 10% KDD KDD DoS NSL-KDD UNSW-NB15 | Constrained-optimization-based Extreme Learning Machines (cELM) | Binary (98.90), multi-class (99.90) |
| [6] | Real network traffic | Principal Component Analysis (PCA) and Ant Colony Optimization | 96.00 |
| [7] | Real network traffic | Fuzzy Logic | 96.50 |
| [8] | KDDCup99 | Multi-level hybrid Support Vector Machine (SVM) and ELM | 95.80 |
| [9] | NSL-KDD | SVM-Radial Basis Function (RBF) | 98.10 |
| [10] | NSL-KDD | Single hidden layer feed-forward neural network (SLFN) | 84.10 |
| [11] | KDDCup99 | Hybrid k-means and SVM-RBF | 88.70 |
| [12] | NSL-KDD UNSW-NB15 | Hybrid Artificial Bee Colony (ABC) and Artificial Fish Swarm (AFS) | NSL-KDD (99.00) UNSW-NB15 (98.90) |
| [13] | KDDCup99 UNSW-NB15 | Genetic Algorithm (GA) as search and Logistic Regression (LR) as learning algorithm | KDDCup99 (99.90) UNSW-NB15 (81.40) |
| [14] | NSL-KDD | Hypergraph based Genetic Algorithm (HG-GA) | 97.10 |
| ]15\ | Self-generated SCADA network | Hierarchical Neuron Architecture based Neural Network (HNA-NN) | 93.10 |
| [16] | NSL-KDD | Time-varying chaos particle swarm optimization (TVCPSO) | 97.20 |
| [17] | NSL-KDD | Marginal density ratio | 99.20 |
| [18] | NSL-KDD | Clustering ELM (Clus-ELM) | 77.00 |

| Reference | Dataset | Method | Accuracy (%) |
|---|---|---|---|
| [19] | KDDCup99 Kyoto University Benchmark Dataset (KUBD) | Anomaly-detection method based on the change of cluster centres (ADBCC), then k-NN | KDDCup99 (93.30) KUBD (95.80) |
| [20] | NSL-KDD (exclude U2R) | Discrete wavelet transform (DWT) | 96.70 |
| [21] | NSL-KDD | Weighted one-against-rest SVM (WOAR-SVM) | 80.70 |
| [22] | NSL-KDD | Two-layer classification, Genetic Algorithm for Detectors Generartions (GADG), Random Forest Tree | 98.600 |
| [23] | NSL-KDD ISCX 2012 | Hybrid Artificial Bee Colony (ABC) and AdaBoost | 98.90 |
| [24] | Gure-KDD KDDCup99 | Improved many-objective optimization (I-NSGA-III) | Gure-KDD (99.60) KDDCup99 (99.40) |
| [25] | KDDCup99 | Cluster center and nearest neighbor (CANN) | 99.9 |
| [26] | NSL-KDD | Hybrid J48, Meta Pagging, Random Tree, REPTree, AdaBoostM1, Decision Stump, Naïve Bayes | Binary (99.80), Multiclass (98.60) |
| [27] | KDDCup99 NSL-KDD UNSW-NB15 | Dendron (DT and GA) | KDDCup99 (98.90), NSl-KDD (97.60), UNSQ-NB15 (84.30) |

Further analysis of this study found that 65% of recent researches focus on classification, utilizing supervised machine learning techniques (as shown in Fig. 2). Therefore, only labelled datasets were used in these researches. However, it may pose a problem in detecting anomalous intrusion, which includes new or modified intrusion attacks.



**Fig. 2** Research focus area of IDS from 2015-2018

In addition, this study also found that 44% of recent researches used the soft-computing (ensemble and hybrid) approach to tackle IDS problem, as shown in Fig. 3. This proves that soft-computing techniques are getting considerable attention from researchers in IDS.



**Fig. 3** Approach used in IDS research from 2015-2018

## IV. DATASET

Dataset is the key component to train machine learning to detect anomaly threats. However, the analysis from this study shows that many researchers are still relying on an outdated dataset, KDDCup99 and NSL-KDD (a variant of KDD00 dataset), which have been criticized by many as outdated and not relevant in current network infrastructure. This dataset was produced in 1999, which is almost 20 years old. Rapid development and changes in Information Technology such as cloud computing, social media and Internet of Things are changing the landscape of network infrastructure. These changes have the driving force in changing threat attack itself. Therefore, many research results that demonstrate high accuracy is being viewed as overstated, because the dataset

being used does not represent the current threat or infrastructure.

The KDDCup99 dataset is a popular dataset and has been used for the Third International Knowledge Discovery and Data Mining Tools Competition. Each connection instance is described by 41 attributes (38 continuous or discrete numerical attributes and 3 symbolic attributes). Each instance is labelled as either normal or a specific type of attack. These attacks fall under one of the four categories: Probe, DoS, U2R, and R2L [9], as described below.

Probing: This type of attack collects information of target system prior to initiating an actual attack.

Denial of Service (DoS): This type of attack results in unavailability of network resources to legitimate requests by exhausting the bandwidth or by overloading computational resources.

User to Root (U2R): In this case, an attacker starts out with access to a normal user account on the system and is able to exploit the system's vulnerabilities to gain root access to the system.

Remote to Local (R2L): In this case, an attacker who does not have an account on a remote machine sends a packet to that machine over a network and exploits some vulnerabilities to gain local access as a user of that machine.



**Fig. 4** Dataset being used in IDS research from 2015-2018

The NSL-KDD dataset was developed in 2009, but it is actually an improved version of the KDDCup99 dataset. NSL-KDD tries to improve KDDCup99 dataset by removing redundant records, including the imbalanced number of instances and the variety of attack classes [2]. However, it still inherited the fundamental limitation of the dataset.

KDDCup99 has many drawbacks. Firstly, this dataset was developed in 1999 using a Solaris-based operating system to collect a wide range of data due to its easy deployment. However, there are significant differences in today's operating systems which barely resemble Solaris. In this age of Ubuntu, Windows and MAC, Solaris has almost no market share.

Secondly, the traffic collector used in KDD datasets, TCPdump, is very likely to become overloaded and drop

packets from a heavy traffic load. More importantly, there is some confusion about the attack distributions of these datasets. According to an attack analysis, Probe is not an attack unless the number of iterations exceeds a specific threshold, while label inconsistency has been reported [26].

Thirdly, the emergence of new technologies such as cloud computing, social media and the Internet of Things has changed the network infrastructure drastically. These changes will also result in new types of threat.

The other two popular datasets are ISCX 2012 and UNSW-NB15. ISCX 2012 is a dataset created by Information Security Centre of Excellence (ISCX) at University of New Brunswick in 2012. This dataset consists of seven days of data with labelling of normal (one) or attack (two). The dataset has no classification of the types of attack, thus it will only provide binary classification. However, this dataset is no longer available. This is because the centre has created a new dataset, called CICIDS2017 [28]. The centre has also changed its name to Canadian Institute for Cybersecurity (CIC). Unfortunately, no article was found using this new dataset at the time of this study.

Another popular dataset is UNSW-NB15, this dataset was created by Australia Centre for Cyber Security (ACCS) using IXIA PerfectStorm to generate nine types of attack. These nine types of attack are namely fuzzers, analysis, backdoors, DoS, exploits, generic, reconnaissance, shellcode, and worms. The dataset has a total of 47 features with two labels. First is named as 'Label', where zero indicates normal and one indicates an attack. Second label is named as 'attack_cat', which provides the type of attack [29].

## V. METRIC

Metric is the quantitative evaluation of ML algorithm performance towards specific dataset. It provides a way for comparison, to determine which model performance better and by how much. Most metrices can be derived from a confusion matrix table, as shown in Table 2 below.

Accuracy is the most often used metric. This metric provides the ratio of correctly predicted outcome compared to total observed outcome [15]. Therefore it is being used as the primary metric for comparison in this study. The formula is shown in equation 1:

$$\frac{TN + TP}{TN + TP + FP + FN} \qquad (1)$$

True Positive Rate (TPR) has three other names, but all used the same formula. These names are recall, sensitivity, and detection rate. This metric is the ratio of correctly predicted positive outcome compared to actually positive observation [15]. The formula is shown in equation 2 below:

$$\frac{TP}{TP + FN} \qquad (2)$$

**Table 2.** Confusion matrix table

| | | Predicted Class | |
|---|---|---|---|
| | | Negative (Normal) | Positive (Attack) |
| Actual Class | Negative (Normal) | True Negative (TN) | False Positive(FP) |
| | Positive (Attack) | False Negative (FN) | True Positive(TP) |

False Positive Rate (FPR) is also called false alarm rate (FAR) or fall-out. This metric is the ratio of wrongly predicted positive outcome compared to actual negative observation [15]. The formula is shown in equation 3 below:

$$\frac{FP}{FP + TN} \qquad (3)$$

True Negative Rate (TNR) is also called specificity. This metric is the ratio of correctly predicted negative outcome compared to actually negative observation [15]. The formula is shown in equation 4 below:

$$\frac{TN}{TN + FP} \qquad (4)$$

False Negative Rate (FNR) is also called miss rate. This metric is the ratio of wrongly predicted negative outcome compared to actually positive observation [21]. The formula is shown in equation 5 below:

$$\frac{FN}{FN + TP} \qquad (5)$$

Precision is the ratio of correctly predicted positive outcome compared to positive prediction [15]. The formula is shown in equation 6 below:

$$\frac{TP}{TP + FP} \qquad (6)$$

F-measure is also called F-score. This metric provide performance evaluation based on precision and recall [22]. The formula is shown in equation 7 below:

$$\frac{2TP}{2TP + FP + FN} \qquad (7)$$

Time is the measurement of efficiency. Two phases of measurement can be performed. One measurement during training phase, and the second during testing phase. There are other metrices found in this study, but are less common, thus those will not be discussed here.

This study found that two metrices were used in more than 70% of researches. These are accuracy and TPR. Accuracy provides good indication of how well the algorithm can predict the correct outcome. This is important, because it shows how much the result can be trusted to be correct.

TPR, or better known as detection rate, provide an indication of how well the algorithm can detect and intrusion attack. The purpose of IDS is to detect an attack, thus this metric is important.

Another metric that was used in more than 50% of researches is FPR. Another name for this metric is False Alarm Rate (FAR). This metric provides indication whether the algorithm will produce many false alarms. This is important, because it shows how much more work is needed to further filter out these false alarms observation, after the IDS. This is most probably performed by a human expert.



**Fig. 5** Percentage of metric being used in IDS research from 2015-2018

## VI. CONCLUSION

Soft computing techniques are getting considerable attention from researchers in IDS. This is because this technique is easy to apply and often produce better result compared to single algorithm. Proper combination of multiple algorithms is the way forward. Most researchers are focusing on the classification of IDS, which is beneficial in determining known intrusion attacks. However, it may pose a problem in detecting anomalous intrusion, which may include new or modified intrusion attacks. Therefore to produce a more robust IDS, clustering algorithm should be considered for future development. KDDCup99 and its variant NSL-KDD datasets are the two most widely used datasets, although they are almost 20 years old. This continuous trend could result in static progress in IDS, while intrusion attacks continue to evolve together with new technologies and user behaviours. Ultimately, this situation will result in the obsolete use of IDS as part of a cyber security tool. Therefore new dataset that represent current environment setup, both software and hardware, is important. The latest publicly available dataset is CICIDS2017, should be explored.

Three most used metrices for performance evaluation for IDS are accuracy, TPR and FPR. This is expected, because these metrices provide important indications that are very relevant to IDS functionality. In order to simply the evaluation process, it is possible to develop a metric that can combine all three metrices.

## REFERENCES

[1] J. Brownlee, "A Tour of Machine Learning Algorithms", https://machinelearningmastery.com/a-tour-of-machine-learning-algorithms/ 2013

[2] E. K. Viegas, A. O. Santin, and L. S. Oliveira, "Toward a reliable anomaly-based intrusion detection in real-world environments," Comput. Networks, vol. 127, pp. 200–216, 2017.

[3] A. Verma and V. Ranga, "Statistical analysis of CIDDS-001 dataset for Network Intrusion Detection Systems using Distance-based Machine Learning," Procedia Comput. Sci., vol. 125, pp. 709–716, 2018.

[4] T. Hamed, R. Dara, and S. C. Kremer, "Network intrusion detection system based on recursive feature addition and bigram technique," Comput. Secur., vol. 73, pp. 137–155, 2018.

[5] C. R. Wang, R. F. Xu, S. J. Lee, and C. H. Lee, "Network intrusion detection using equality constrained-optimization-based extreme learning machines," Knowledge-Based Syst., vol. 147, pp. 68–80, 2018.

[6] G. Fernandes, L. F. Carvalho, J. J. P. C. Rodrigues, and M. L. Proença, "Network anomaly detection using IP flows with Principal Component Analysis and Ant Colony Optimization," J. Netw. Comput. Appl., vol. 64, pp. 1–11, 2016.

[7] A. H. Hamamoto, L. F. Carvalho, L. D. H. Sampaio, T. Abrão, and M. L. Proença, "Network Anomaly Detection System using Genetic Algorithm and Fuzzy Logic," Expert Syst. Appl., vol. 92, pp. 390–402, 2018.

[8] W. L. Al-Yaseen, Z. A. Othman, and M. Z. A. Nazri, "Multi-level hybrid support vector machine and extreme learning machine based on modified K-means for intrusion detection system," Expert Syst. Appl., vol. 67, pp. 296–303, 2017.

[9] I. Sumaiya Thaseen and C. Aswani Kumar, "Intrusion detection model using fusion of chi-square feature selection and multi class SVM," J. King Saud Univ. - Comput. Inf. Sci., vol. 29, no. 4, pp. 462–472, 2017.

[10] R. A. R. Ashfaq, X. Z. Wang, J. Z. Huang, H. Abbas, and Y. L. He, "Fuzziness based semi-supervised learning approach for intrusion detection system," Inf. Sci. (Ny)., vol. 378, pp. 484–497, 2017.

[11] U. Ravale, N. Marathe, and P. Padiya, "Feature selection based hybrid anomaly intrusion detection system using K Means and RBF kernel function," Procedia Comput. Sci., vol. 45, no. C, pp. 428–435, 2015.

[12] V. Hajisalem and S. Babaie, "A hybrid intrusion detection system based on ABC-AFS algorithm for misuse and anomaly detection," Comput. Networks, vol. 136, pp. 37–50, 2018.

[13] C. Khammassi and S. Krichen, "A GA-LR wrapper approach for feature selection in network intrusion detection," Comput. Secur., vol. 70, pp. 255–277, 2017.

[14] M. R. Gauthama Raman, N. Somu, K. Kirthivasan, R. Liscano, and V. S. Shankar Sriram, "An efficient intrusion detection system based on hypergraph - Genetic algorithm for parameter optimization and feature selection in support vector machine," Knowledge-Based Syst., vol. 134, pp. 1–12, 2017.

[15] S. Shitharth and D. Prince Winston, "An enhanced optimization based algorithm for intrusion detection in SCADA network," Comput. Secur., vol. 70, pp. 16–26, 2017.

[16] S. M. Hosseini Bamakan, H. Wang, T. Yingjie, and Y. Shi, "An effective intrusion detection framework based on MCLP/SVM optimized by time-varying chaos particle swarm optimization," Neurocomputing, vol. 199, pp. 90–102, 2016.

[17] H. Wang, J. Gu, and S. Wang, "An effective intrusion detection framework based on SVM with feature augmentation," Knowledge-Based Syst., vol. 136, pp. 130–139, 2017.

[18] S. Roshan, Y. Miche, A. Akusok, and A. Lendasse, "Adaptive and online network intrusion detection system using clustering and Extreme Learning Machines," J. Franklin Inst., vol. 355, no. 4, pp. 1752–1779, 2018.

[19] C. Guo, Y. Ping, N. Liu, and S. S. Luo, "A two-level hybrid approach for intrusion detection," Neurocomputing, vol. 214, pp. 391–400, 2016.

[20] S. Y. Ji, B. K. Jeong, S. Choi, and D. H. Jeong, "A multi-level intrusion detection method for abnormal network behaviors," J. Netw. Comput. Appl., vol. 62, pp. 9–17, 2016.

[21] A. A. Aburomman and M. Bin Ibne Reaz, "A novel weighted support vector machines multiclass classifier based on differential evolution for intrusion detection systems," Inf. Sci. (Ny)., vol. 414, pp. 225–246, 2017.

[22] A. S. Amira, S. E. O. Hanafi, and A. E. Hassanien, "Comparison of classification techniques applied for network intrusion detection and classification," J. Appl. Log., vol. 24, pp. 109–118, 2017.

[23] M. Mazini, B. Shirazi, and I. Mahdavi, "Anomaly network-based intrusion detection system using a reliable hybrid artificial bee colony and AdaBoost algorithms," J. King Saud Univ. - Comput. Inf. Sci., 2018.

[24] Y. Zhu, J. Liang, J. Chen, and Z. Ming, "An improved NSGA-III algorithm for feature selection used in intrusion detection," Knowledge-Based Syst., vol. 116, pp. 74–85, 2017.

[25] W. C. Lin, S. W. Ke, and C. F. Tsai, "CANN: An intrusion detection system based on combining cluster centers and nearest neighbors," Knowledge-Based Syst., vol. 78, no. 1, pp. 13–21, 2015.

[26] S. Aljawarneh, M. Aldwairi, and M. B. Yassein, "Anomaly-based intrusion detection system through feature selection analysis and building hybrid efficient model," J. Comput. Sci., vol. 25, pp. 152–160, 2018.

[27] D. Papamartzivanos, F. Gómez Mármol, and G. Kambourakis, "Dendron: Genetic trees driven rule induction for network intrusion detection systems," Futur. Gener. Comput. Syst., vol. 79, pp. 558–574, 2018.

[28] A. H. L. and A. A. G. Iman Sharafaldin, "Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization," Proc. 4th Int. Conf. Inf. Syst. Secur. Priv., no. Cic, pp. 108–116, 2018.

[29] N. Moustafa and J. Slay, "UNSW-NB15: a comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set)," 2015 Mil. Commun. Inf. Syst. Conf., no. November, pp. 1–6, 2015.

[30] A. A. Aburomman and M. B. I. Reaz, "A survey of intrusion detection systems based on ensemble and hybrid classifiers," Comput. Secur., vol. 65, pp. 135–152, 2017.