

Anomaly Classification Using Genetic Algorithm-Based Random Forest Model for Network Attack Detection

Adel Assiri*

Management Information Systems Department, College of Business, King Khalid University, Abha, 61421, Saudi Arabia

*Corresponding Author: Adel Assiri. Email: adaseri@kku.edu.sa

Received: 22 August 2020; Accepted: 14 September 2020

Abstract: Anomaly classification based on network traffic features is an important task to monitor and detect network intrusion attacks. Network-based intrusion detection systems (NIDSs) using machine learning (ML) methods are effective tools for protecting network infrastructures and services from unpredictable and unseen attacks. Among several ML methods, random forest (RF) is a robust method that can be used in ML-based network intrusion detection solutions. However, the minimum number of instances for each split and the number of trees in the forest are two key parameters of RF that can affect classification accuracy. Therefore, optimal parameter selection is a real problem in RF-based anomaly classification of intrusion detection systems. In this paper, we propose to use the genetic algorithm (GA) for selecting the appropriate values of these two parameters, optimizing the RF classifier and improving the classification accuracy of normal and abnormal network traffics. To validate the proposed GA-based RF model, a number of experiments is conducted on two public datasets and evaluated using a set of performance evaluation measures. In these experiments, the accuracy result is compared with the accuracies of baseline ML classifiers in the recent works. Experimental results reveal that the proposed model can avert the uncertainty in selection the values of RF's parameters, improving the accuracy of anomaly classification in NIDSs without incurring excessive time.

Keywords: Network-based intrusion detection system (NIDS); random forest classifier; genetic algorithm; KDD99; UNSW-NB15

1 Introduction

Network-based intrusion detection system (NIDS) is a network security tool that works together with popular data encryption algorithms and firewalls to protect network resources and services [1]. The work to develop an effective NIDS to detect malicious activities and network intrusion attacks is still the motivation for developers and researchers. In the literature of intrusion detection system (IDS), a number of methods and models have been proposed to prevent the networks from malicious threats and attacks. For instance, Song et al. [2], Gong et al. [3], and Murugesan et al. [4] offered several techniques to trace back the IP address. Nguyen et al. [5], Crotti et al. [6], and Callado et al. [7] introduced a number of methods to classify IP traffics of the networks. Dharmapurikar et al. [8], Zhou et al. [9], Chen et al. [10],



This work is licensed under a Creative Commons Attribution 4.0 International License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Hu et al. [11], Das et al. [12], and Mabu et al. [13] developed many techniques for intrusion detection in the networks and explained the performance, the advantages, and disadvantages of these developed techniques. Hadlington [14] presented a study summaries the human causes that leads to some cyber security violation issues. Alternatively, the machine learning (ML) methods can be a suitable and adaptive approach for detecting abnormal network traffics due to the continuous changes of attacks patterns.

Recently, ML methods have been used for solving many problems in different applications [15,16]. The ML-based data analysis is utilized as a tool for automatic classification [17–19], decision making [20], and prediction [21,22]. For network attacks detections, threats, and malicious executables, supervised and unsupervised learning technique has been achieved a promising result [23–26]. ML learning approach is capable to give the networks devices the ability to learn attacks patterns from past data traffics and detect the unknown and new attacks [27–29].

In previous works, numerous studies have been proposed using ML methods for network intrusion detection. The work introduced by Solanki et al. [30] has computed the accuracy of decision tree (C4.5) and support vector machine (SVM) methods to detect intrusion attacks. The two methods were tested on a public dataset contains four attacks. The authors reported that the accuracy of SVM was less than the accuracy of C4.5 method. The authors in [31] offered a work based on a number of ML classifiers to detect the common attacks in the networks and determining the best. They determined the suitable classifier for each of attack and reported that the most classifiers have achieved a high accuracy result to detect the denial of service attacks.

Gao et al. [32] developed a technique to analyze the normal and abnormal data traffics using hidden Markov model. The authors did a set of experiments and achieved 63.2% accuracy result. Gomez et al. [33] proposed a study for network-based IDS based on fuzzy logic. The authors in [9] introduced an approach for classifying periodic patterns of network traffics and detecting normal or abnormal behaviors using Fourier transform method.

An audit technique based on the frequency happened in the data traffics of the networks has proposed by Ye et al. [34]. However, the data used for testing in this study was simple, pure and did not reflect the real states of network traffics. Additionally, a Chi-square test is used by Goonatilake et al. [35] for detecting abnormal network traffics in IDS. The study in [36] has compared and analyzed the performance of five architectures of artificial neural networks for IDS. The study shows that the quasi-Newton and conjugate gradient descent attained improved accuracy results. The works in [37,38] have developed some models for detecting network intrusion attacks using ML methods with swarm intelligence algorithm.

Some comparative studies on ML methods have been proposed for network defense [39] and botnet attack detection [40,41]. However, the performance of ML methods for anomaly classification in NIDSs still needs more improvements in terms of time cost and accuracy. Recently, Khan et al. [42] introduced a comparative study on ML methods for NIDS. They have mentioned that the accuracy results of random forest (RF) is better than the other ML methods.

In this study, a genetic-based RF model is proposed and compared with the baseline ML methods for network intrusion detection in the state-of-the-arts. The experiment is conducted on two available public datasets, namely, KDD99 [43] and UNSW-NB15 [44]. The main contribution of the study is to apply the genetic algorithm (GA) to select the appropriate values of RF classifier for improving its accuracy result for network intrusion detection. Moreover, another contribution of the work is to present a comparative study on ML methods for anomaly classification in NIDS using a set of evaluation measures.

The rest of the paper is structured as follows: Section 2 describes the research methods and the main steps of the proposed IDS model. The experiments and discussion on the used datasets are given in Section 3. Finally, Section 4 summarizes the conclusion of the work.

2 Research Methods

2.1 GA

GA was presented initially by Holland [45]. It is a form of inductive learning strategy to provide another method to conventional optimization methods based on adaptive search techniques. GA can find the near-optimal solution for problems that need complex optimization. It is a stochastic method depends on some natural phenomena based on natural selection and genetic inheritance. GA is the most common class of EA [46]. GA works on a population of individuals or chromosomes that represent the candidate solutions for a given problem. Each individual compete with others to reproduce based on Darwin's principle (survival of the fittest) in each generation of evolution.

All the individuals are evaluated by a fitness function that expresses the importance of the individual as a solution. Then select the best parent individuals and apply the crossover and mutation operator to produce the new individuals (offspring) for the next generation. Crossover operator combines the features of two selected parents to create two offspring. Mutation operator changes one or more components of the selected individual in order to prevent any stagnation that may occur during the search process. After a number of generations in evolution when the stopping criterion is met, the individuals that survived in the population are considered the optimal solutions [29]. Algorithm 1 summarizes the main steps of GA.

Algorithm 1: Genetic algorithm (GA)

```

1: procedure GENETIC ALGORITHM
2:   Generate randomly the initial population of solutions
3:   Evaluate the initial population by fitness function
4:   repeat until (a stop condition is satisfied) do
5:     Select a pair of parents based on fitness
6:     Create two offspring using Crossover
7:     Apply Mutation to the offspring
8:     Evaluate the new offspring by fitness function
9:     Evaluate the new offspring by fitness function
10:  end do
11:  return best parameters' values
12: end procedure

```

2.2 RF Classifier

The RF classifier is a powerful ML tool that can be used for solving classification and regression problems. RF is one of the ensemble learning methods that can build a number of decision trees [47]. For building trained RF model, two steps of randomness are used:

- Individually and randomly, each decision tree is constructed using different samples of the training dataset.
- During the construction of each tree, a part of m samples is randomly selected from the training dataset. The split point of these m samples is used as best split. In a case of new sample c , the RF can classify or predict c by aggregated decision trees. For RF that has n decision trees, the output is the probability of the class label y for the sample c given a feature vector x . The equation of RF ensemble learning can be computed as follows:

$$P(y/x) = \frac{1}{n} \sum_{i=1}^n P_i(y/x) \quad (1)$$

In other words, the RF can average the probability of decision trees obtained using different random samples of the original dataset [47]. Fig. 1 visualizes the construction process of RF according to ensemble learning concept.

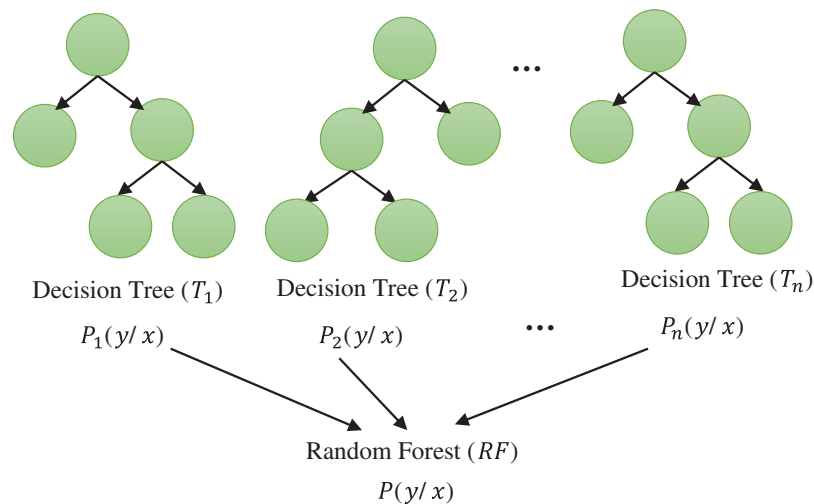


Figure 1: Construction process of random forest (RF) according to ensemble learning concept

The RF classifier has been used in a wide range of applications, such as image classification [48], network intrusion detection [49], and neuroimaging [50]. Algorithm 2 defines the RF steps.

Algorithm 2: Random forest (RF) pseudo-code

- 1: **procedure** RANDOM FOREST ALGORITHM
 - 2: **for** $i = 1$ to T trees **do**
 - 3: Pick n data points ($D_{i=1..n}$) with replacement from training dataset (D)
 - 4: Build full decision tree on $D_{i=1..n}$
 - 5: for each split, consider only k features that are picked uniformly at random new features for every split
 - 6: Prune tree to minimize out-of-bag erro
 - 7: **end for**
 - 8: Average all T trees
 - 9: **end procedure**
-

In this research, we explore the application of GA-based RF for detecting intrusion attack throughout the features of network data traffic.

2.3 GA-Based RF Model for IDS

The idea behind the GA-based RF model is to optimize the RF classifier by selecting the appropriate parameters' values and improve the detection rate of NIDS by using the optimized RF. The GA can generate random values for the specific parameters of RF and build a new decision boundary that has a highest value of GA fitness function. In detail, the datasets for training and testing the GA-based RF model are prepared from the network data traffics. The decision boundary of GA-based RF model is trained using training set and GA. After that, the trained GA-based RF model with the appropriate parameters' values is tested to detect normal and abnormal class label of samples in the testing set. Fig. 2 illustrates the main steps to build GA-based RF model for IDS.

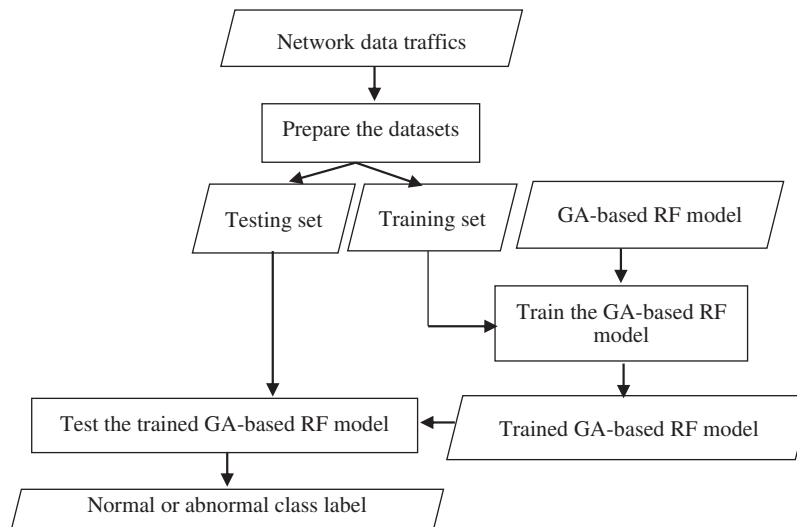


Figure 2: The main steps to build genetic algorithm (GA)-based random forest (RF) model for intrusion detection system (IDS)

3 Experiments and Discussion

The study experiments are conducted on a laptop has a CPU processor Intel Core i7-4510U with 2.0 GHz, 8 GB RAM, and a 64 bit Windows 10 operating system. Python programming language is used to implement the experiments. Two public datasets, namely, KDD99 and UNSW-NB15 are employed to evaluate and compare the proposed model.

3.1 Datasets Description

As mentioned above, the datasets used in the experiments are KDD99 [43] and UNSW-NB15 [44] datasets. The KDD99 dataset is divided into two sets: A training set contains 145,586 samples and testing set includes 73,269 samples. The UNSW-NB15 dataset is also separated into two sets: a training set consists of 175,341 samples and a testing set has 82,332 samples. These datasets are processed and normalized to be suitable for training and testing the models. Figs. 3 and 4 display the distribution of samples in the training and testing sets according to normal and abnormal network traffics.

To evaluate the proposed GA-based RF and other baseline classifiers, the training samples of two sets are used first to train these classifiers and build trained models; then, these trained models are tested on the two testing sets.

3.2 Performance Evaluation Measures

The results of experiments are assessed based on three measures. These measures are accuracy, sensitivity, and precision, computed as follows:

$$\text{Accuracy} = \frac{(TP + TN)}{(TP + TN + FP + FN)} \quad (2)$$

$$\text{Recall (Sensitivity)} = \frac{TP}{(TP + FN)} \quad (3)$$

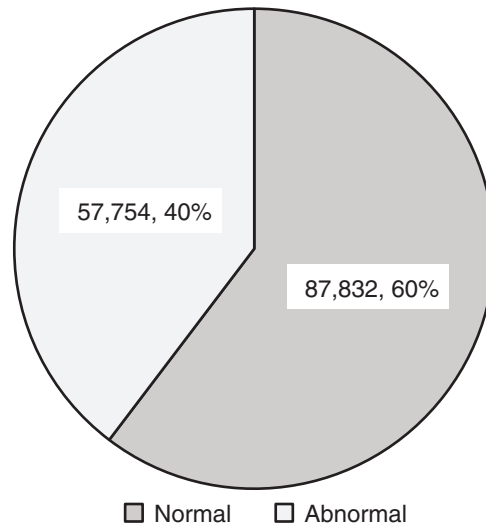


Figure 3: The number of normal or abnormal samples in the KDD99 training set

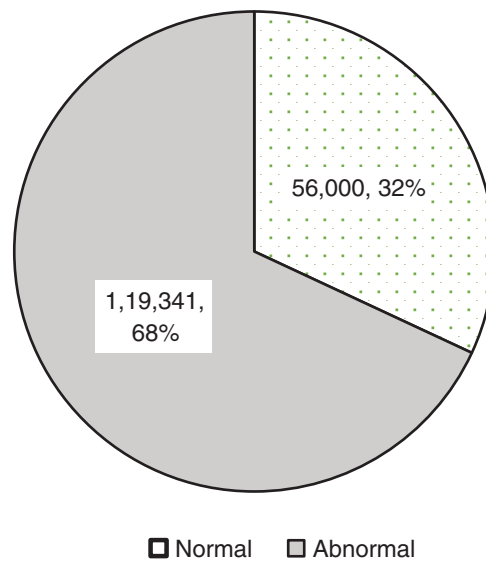


Figure 4: The number of normal and abnormal samples in the UNSW-NB15 training set

$$\text{Precision} = \frac{TP}{(TP + FP)} \quad (4)$$

$$\text{F1-Score} = 2 * \frac{(\text{Precision} * \text{Recall})}{(\text{Precision} + \text{Recall})} \quad (5)$$

FP and FN are the number of false positives and negatives. TP and TN are the number of true positives and negatives.

3.3 Results and Comparisons

In this section, the results of the experiments are presented and compared with the results of recent related work. After building the GA-based RF model using the KDD99 training set, the best values of the minimum number of instances for each split and the number of trees in the RF are selected to be 17 and 2, respectively. For the UNSW-NB15 training set, the value of the minimum number of instances for each split is 4 and the value of the number of trees in the forest is also 2. The other parameters of RF are fixed to have the default values. [Tabs. 1](#) and [2](#) demonstrate the results of confusion matrices for testing the model on the KDD99 and UNSW-NB15 testing sets.

Table 1: Results of confusion matrix for normal and abnormal classification of the KDD99 testing set

	Normal traffic	Abnormal traffic
Normal traffic	47630	283
Abnormal traffic	1808	23548

Table 2: Results of confusion matrix for normal and abnormal classification of the UNSW-NB15 testing set

	Normal traffic	Abnormal traffic
Normal traffic	28267	8733
Abnormal traffic	2228	43104

From the results of confusion matrices, the performance evaluation measures are computed and shown in [Tabs. 3](#) and [4](#). As seen in the [Tab. 3](#), the GA-based RF achieves 97.2% of the accuracy and 97.0% for the weighted average of precision and recall on the KDD99 testing set. In addition, it obtains 86.7% of the accuracy and 87.0% for the weighted average of precision and recall on the UNSW-NB15 testing, which is noisy and more complex.

Table 3: The performance evaluation results for anomaly classification of the KDD99 testing set

Class Na	Evaluation measure		
	Precision	Recall	F1-score
Normal traffic	96.0%	99.0%	98.0%
Abnormal traffic	99.0%	93.0%	96.0%
Accuracy	97.2%		
Macro avg.	98.0%	96.0%	97.0%
Weighted avg.	97.0%	97.0%	97.0%

To compare the accuracy results of optimized RF classifier to classify anomalies with the traditional RF and other baseline classifiers in the recent work [42], [Tabs. 5](#) and [6](#) show the accuracy results on the same testing sets of KDD99 and UNSW-NB15.

Table 4: The performance evaluation results for anomaly classification of the UNSW-NB15 testing set

Class Na	Evaluation measure		
	Precision	Recall	F1-score
Normal traffic	93.0%	76.0%	84.0%
Abnormal traffic	83.0%	95.0%	89.0%
Accuracy	86.7%		
Macro avg.	88.0%	86.0%	86.0%
Weighted avg.	87.0%	87.0%	86.0%

Table 5: The accuracy results of the GA-based RF model compared with the baseline classifiers in recent work [42] using KDD99 testing set

Work/authors [ref.]	Classifier name	Accuracy	Weighted average of precision	Weighted average of recall
Khan et al. [42]	NB	94.68%	95%	95%
	KNN	96.01%	96%	96%
	SVM-Poly	94.04%	94%	94%
	NB-KE	94.43%	95%	94%
	SMO	95.11%	95%	95%
	SVM-RBF	94.95%	95%	95%
	DS	93.98%	94%	94%
	DT	96.22%	96%	96%
	RF	96.79%	97%	97%
	HT	92.66%	93%	92%
This work	GA-based RF	97.20	97%	97%

Note. SVM, support vector machine; RF, random forest; GA, genetic algorithm.

Table 6: The accuracy results of the GA-based RF model compared with the baseline classifiers in recent work [42] using UNSW-NB15 testing set

Work/authors [ref.]	Classifier name	Accuracy	Weighted average of precision	Weighted average of recall
Khan et al. [42]	NB	76.39%	78%	76%
	KNN	84.49%	86%	85%
	SVM-Poly	68.34%	69%	68%
	NB-KE	76.22%	77%	76%
	SMO	85.34%	86%	85%
	SVM-RBF	83.22%	84%	83%

Table 6 (continued).

Work/authors [ref.]	Classifier name	Accuracy	Weighted average of precision	Weighted average of recall
	DS	76.63%	84%	77%
	DT	84.55%	86%	85%
	RF	83.63%	87%	84%
	HT	59.44%	76%	59%
This work	GA-based RF	86.70%	87%	87%

Note. SVM, support vector machine; RF, random forest; GA, genetic algorithm.

As shown in the [Tabs. 5 and 6](#), the accuracy results highlighted in the boldface font clarify that the GA-based RF improves the accuracy of the RF due to selecting the best values of its parameters and outperforms the other ML baseline classifiers.

4 Conclusion and Future Work

In this paper, a GA-based RF model is proposed to classify normal and abnormal networks traffics for IDS. The GA is used for selecting the appropriate values for two parameters of RF. These parameters are the minimum number of instances for each split and the number of trees in the forest, optimizing the RF classifier and improving the accuracy of anomaly classification and intrusion detection. A set of experiments were conducted on two public dataset and evaluated using a set of performance evaluation measures. The experimental results revealed that the selection of suitable values of RF classifier has improved the accuracy of network anomaly classification compared to the RF with default values. Moreover, the proposed GA-based RF model outperforms the ML models with high detection rates of 97.20% for KDD99 test set and 86.70% for UNSW-NB15 test set. In the future work, the proposed model will be used with feature selection methods to detect the types of attacks in the abnormal network traffic and enhance the network-based IDS.

Acknowledgement: The author would like to express his gratitude to King Khalid University, Saudi Arabia for providing administrative and technical support.

Funding Statement: The authors received no specific funding for this study.

Conflicts of Interest: The author declares that he have no conflicts of interest to report regarding the present study.

References

- [1] G. Li, Z. Yan, Y. Fu and H. Chen, "Data fusion for network intrusion detection: A review," *Security and Communication Networks*, vol. 2018, 8210614, 2018.
- [2] D. X. Song and A. Perrig, "Advanced and authenticated marking schemes for IP traceback," in *Proc. IEEE INFOCOM 2001. Conf. on Computer Communications. Twentieth Annual Joint Conf. of the IEEE Computer and Communications Society (Cat No. 01CH37213)*, Anchorage, AK, pp. 878–886, 2001.
- [3] C. Gong and K. Sarac, "A more practical approach for single-packet IP traceback using packet logging and marking," *IEEE Transactions on Parallel and Distributed Systems*, vol. 19, no. 10, pp. 1310–1324, 2008.
- [4] V. Murugesan, M. Shalinie and N. Neethimani, "A brief survey of IP traceback methodologies," *Acta Polytechnica Hungarica*, vol. 11, no. 9, pp. 197–216, 2004.

- [5] T. T. Nguyen and G. Armitage, "A survey of techniques for internet traffic classification using machine learning," *IEEE Communications Surveys & Tutorials*, vol. 10, no. 1–4, pp. 56–76, 2008.
- [6] M. Crotti, F. Gringoli, P. Pelosato and L. Salgarelli, "A statistical approach to IP-level classification of network traffic," in *2006 IEEE Int. Conf. on Communications*, Istanbul, Turkey, pp. 170–176, 2006.
- [7] A. Callado, C. Kamienski, G. Szabo, B. P. Gero, J. Kelner *et al.*, "A survey on internet traffic identification," *IEEE Communications Surveys & Tutorials*, vol. 11, no. 3, pp. 37–52, 2009.
- [8] S. Dharmapurikar and J. W. Lockwood, "Fast and scalable pattern matching for network intrusion detection systems," *IEEE Journal on Selected Areas in Communications*, vol. 24, no. 10, pp. 1781–1792, 2006.
- [9] M. Zhou and S. D. Lang, "Mining frequency content of network traffic for intrusion detection," in *Proc. of the IASTED Int. Conf. on Communication, Network, and Information Security*, pp. 101–107, 2003.
- [10] L. Chen and J. Leneutre, "A game theoretical framework on intrusion detection in heterogeneous networks," *IEEE Transactions on Information Forensics and Security*, vol. 4, no. 2, pp. 165–178, 2009.
- [11] W. Hu, W. Hu and S. Maybank, "Adaboost-based algorithm for network intrusion detection," *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 38, no. 2, pp. 577–583, 2008.
- [12] A. Das, D. Nguyen, J. Zambreno, G. Memik and A. Choudhary, "An FPGA-based network intrusion detection architecture," *IEEE Transactions on Information Forensics and Security*, vol. 3, no. 1, pp. 118–132, 2008.
- [13] S. Mabu, C. Chen, N. Lu, K. Shimada and K. Hirasawa, "An intrusion-detection model based on fuzzy class-association-rule mining using genetic network programming," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 41, no. 1, pp. 130–139, 2011.
- [14] L. Hadlington, "Human factors in cybersecurity; examining the link between Internet addiction, impulsivity, attitudes towards cybersecurity, and risky cybersecurity behaviours," *Heliyon*, vol. 3, no. 7, e00346, 2017.
- [15] A. Gumaei, M. M. Hassan, A. Alelaiwi and H. Alsalman, "A hybrid deep learning model for human activity recognition using multimodal body sensing data," *IEEE Access*, vol. 7, pp. 99152–99160, 2019.
- [16] A. Gumaei, M. M. Hassan, M. R. Hassan, A. Alelaiwi and G. Fortino, "A hybrid feature extraction method with regularized extreme learning machine for brain tumor classification," *IEEE Access*, vol. 7, pp. 36266–36273, 2019.
- [17] A. Gumaei, R. Sammouda, A. M. S. Al-Salman and A. Alsanad, "Anti-spoofing cloud-based multi-spectral biometric identification system for enterprise security and privacy-preservation," *Journal of Parallel and Distributed Computing*, vol. 124, pp. 27–40, 2019.
- [18] A. Gumaei, R. Sammouda, A. M. S. Al-Salman and A. Alsanad, "An improved multispectral palmprint recognition system using autoencoder with regularized extreme learning machine," *Computational Intelligence and Neuroscience*, vol. 124, pp. 27–40, 2018.
- [19] A. Gumaei, R. Sammouda, A. M. Al-Salman and A. Alsanad, "An effective palmprint recognition approach for visible and multispectral sensor images," *Sensors*, vol. 18, no. 5, pp. 1575, 2018.
- [20] S. K. Pal and A. Skowron, *Rough-Fuzzy Hybridization: A New Trend in Decision Making*. Berlin: Springer, 1999.
- [21] S. M. Weiss and C. A. Kulikowski, *Computer Systems That Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems*. Burlington: Morgan Kaufmann Publishers Inc., 1991.
- [22] M. Al-Rakhami, A. Gumaei, A. Alsanad, A. Alamri and M. M. Hassan, "An ensemble learning approach for accurate energy load prediction in residential buildings," *IEEE Access*, vol. 7, pp. 48328–48338, 2019.
- [23] J. Z. Kolter and M. A. Maloof, "Learning to detect and classify malicious executables in the wild," *Journal of Machine Learning Research*, vol. 7, pp. 2721–2744, 2006.
- [24] S. Siddiqui, M. S. Khan, K. Ferens and W. Kinsner, "Detecting advanced persistent threats using fractal dimension based machine learning classification," in *Proc. of the 2016 ACM on Int. Workshop on Security and Privacy Analytics*, pp. 64–69, 2016.
- [25] F. A. Khan, A. Gumaei, A. Derhab and A. Hussain, "A novel two-stage deep learning model for efficient network intrusion detection," *IEEE Access*, vol. 7, pp. 30373–30385, 2019.
- [26] A. Derhab, M. Guerroumi, A. Gumaei, L. Maglaras, M. A. Ferrag *et al.*, "Blockchain and random subspace learning-based ids for sdn-enabled industrial IoT security," *Sensors*, vol. 19, no. 14, pp. 3119, 2019.

- [27] S. Mukkamala, G. Janoski and A. Sung, "Intrusion detection using neural networks and support vector machines, in *Proc. of the 2002 International Joint Conf. on Neural Networks IJCNN'02 (Cat. No. 02CH37290)*, Honolulu, HI, 1702–1707, 2002.
- [28] M. M. Hassan, A. Gumaei, A. Alsanad, M. Alrubaian and G. Fortino, "A hybrid deep learning model for efficient intrusion detection in big data environment," *Information Sciences*, vol. 513, pp. 386–396, 2020.
- [29] M. Alqahtani, A. Gumaei, H. Mathkour and M. M. B. Ismail, "A genetic-based extreme gradient boosting model for detecting intrusions in wireless sensor networks," *Sensors*, vol. 19, no. 20, pp. 4383, 2019.
- [30] M. Solanki and V. Dhamdhere, "Intrusion detection system using means of data mining by using C 4.5 algorithm," *International Journal of Application or Innovation in Engineering & Management*, vol. 4, no. 5, pp. 2319–4847, 2015.
- [31] H. A. Nguyen and D. Choi, "Application of data mining to network intrusion detection: classifier selection model," in *Asia-Pacific Network Operations and Management Sym.*, Berlin, 2008.
- [32] B. Gao, H. Y. Ma and Y. H. Yang, "Hmms (hidden markov models) based on anomaly intrusion detection method," in *Proc. Int. Conf. on Machine Learning and Cybernetics*, Beijing, China, pp. 381–385, 2002.
- [33] J. Gomez and D. Dasgupta, "Evolving fuzzy classifiers for intrusion detection," in *Proc. of the IEEE Workshop on Information Assurance*, West Point, NY, pp. 321–323, 2002.
- [34] N. Ye, X. Li, Q. Chen, S. M. Emran and M. Xu, "Probabilistic techniques for intrusion detection based on computer audit data," *IEEE Transactions on Systems, Man, and Cybernetics—Part A: Systems and Humans*, vol. 31, no. 4, pp. 266–274, 2001.
- [35] R. Goonatilake, A. Herath, S. Herath, S. Herath and J. Herath, "Intrusion detection using the chi-square goodness-of-fit test for information assurance, network, forensics and software security," *Journal of Computing Sciences in Colleges*, vol. 23, no. 1, pp. 255–263, 2007.
- [36] D. Vu and V. R. Vemuri, "Computer network intrusion detection: A comparison of neural networks methods," *Differential Equations and Dynamical Systems*, vol. 10, no. 1, pp. 2, 2002.
- [37] T. S. Kala and A. Christy, "An intrusion detection system using opposition based particle swarm optimization algorithm and PNN," in *2019 Int. Conf. on Machine Learning, Big Data, Cloud and Parallel Computing*, Faridabad, India, pp. 184–188, 2019.
- [38] M. Mahmood and B. Al-Khateeb, "Review of neural networks and particle swarm optimization contribution in intrusion detection," *Periodicals of Engineering and Natural Sciences*, vol. 7, no. 3, pp. 1067–1073, 2019.
- [39] A. Ali, Y. H. Hu, C. C. G. Hsieh and M. S. Khan, "A comparative study on machine learning algorithms for network defense," *Virginia Journal of Science*, vol. 68, no. 3, pp. 1, 2017.
- [40] S. Ryu and B. Yang, "A comparative study of machine learning algorithms and their ensembles for botnet detection," *Journal of Computer and Communications*, vol. 6, no. 05, pp. 119–129, 2018.
- [41] A. Bansal and S. Mahapatra, "A comparative analysis of machine learning techniques for botnet detection," in *Proc. of the 10th Int. Conf. on Security of Information and Networks*, pp. 91–98, 2017.
- [42] F. A. Khan and A. Gumaei, "A comparative study of machine learning classifiers for network intrusion detection," in *Int. Conf. on Artificial Intelligence and Security*, pp. 75–86, 2019.
- [43] KDD Cup 1999 Data, "*Kdd.ics.uci.edu*," 2018. [Online]. Available: <https://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>.
- [44] UNSW Canberra at the Australian Defense Force Academy, UNSW-NB15 Dataset. Canberra, Australia, 2015. [Online]. Available: <https://www.unsw.adfa.edu.au/australian-centre-for-cybersecurity/cybersecurity/ADFA-NB15-Datasets/>.
- [45] J. H. Holland, *Adaptation in Natural and Artificial Systems*. Ann Arbor: University of Michigan Press, 1975.
- [46] M. Ozdemir, "Evolutionary computing for feature selection and predictive data mining," M.S. thesis. Rensselaer Polytechnic Institute, New York, 2002.
- [47] V. F. Rodriguez-Galiano, B. Ghimire, J. Rogan, M. Chica-Olmo and J. P. Rigol-Sanchez, "An assessment of the effectiveness of a random forest classifier for land-cover classification," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 67, pp. 93–104, 2012.

- [48] P. Du, A. Samat, B. Waske, S. Liu and Z. Li, "Random forest and rotation forest for fully polarized SAR image classification using polarimetric and spatial features," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 105, pp. 38–53, 2015.
- [49] J. Zhang, M. Zulkernine and A. Haque, "Random-forests-based network intrusion detection systems," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 38, no. 5, pp. 649–659, 2008.
- [50] K. R. Gray, P. Aljabar, R. A. Heckemann, A. Hammers and D. Rueckert *et al.*, "Random forest-based similarity measures for multi-modal classification of Alzheimer's disease," *Random Forest-Based Similarity Measures for Multi-Modal Classification of Alzheimer's Disease*, vol. 65, pp. 167–175, 2013.