



A Novel Features Selection Method for Misuse Intrusion Detection System based on RNA Encoding and Raita Algorithm

Dunia Alawi Jarwan¹, Omar Fitian Rashid^{2,*}, M. Jasim Mohammed¹, Shaymaa E. Sarhan¹,
Hind Moutaz Al-Dabbas³, Maythem K. Abbas⁴

¹Department of Mathematics, College of Science, University of Anbar, Ramadi, Iraq

²Department of Geology, College of Science, University of Baghdad, Baghdad, Iraq

³Department of Computer Science, College of Education for Pure Science/Ibn Al-Haitham, University of Baghdad, Baghdad, Iraq

⁴Asia Pacific University of Technology and Innovation Technology Park Malaysia, Bukit Jalil, 57000 Kuala Lumpur, Malaysia

Emails: dunia.alawi@uoanbar.edu.iq; omar.f@sc.uobaghdad.edu.iq; mohadmath87@uoanbar.edu.iq;
Shaymaa.e.alqaisi@uoanbar.edu.iq; hind.moutaz@ihcoedu.uobaghdad.edu.iq; maythem.abbas@apu.edu.my

Abstract

The significance of the Intrusion Detection System (IDS) is due to its capability in detecting attacks over the network. The current paper proposes a new feature selection method for misuse intrusion detection systems based on RNA encoding, where the proposed method includes five steps. Firstly, the KDD-Cup99 dataset is used and then select random records are used for both training and testing. Secondly, RNA encoding to encode each possible value in the dataset into RNA characters. Thirdly, the keys and their locations are extracted by dividing the achieved RNA sequences from previous steps into blocks with different sizes, then finding the most repeated blocks, choosing them as keys, and storing their location. The next step is the proposed feature selection method based on the extracted keys and their locations, depending on the place of the key within the feature number. Finally, the Raita algorithm for matching to search for keys before and after the applied features selection method. In terms of IDS performance evaluation, experimental outcomes of the proposed feature selection method show the capability of optimizing the time complexity and metrics.

Keywords: Features Selection; Intrusion Detection; Misuse; RNA encoding; Matching algorithm

1. Introduction

The Internet dominates our modern world and provides major convenience to clients in different areas, such as banking, e-commerce, education, and social media [1]. The high use of computer networks led to the "Information security" problem. Information security aims to control access and protect information security, integrity, and personal information. The main objectives of attackers are compromising computer systems, stealing personal information, and destroying part of the system. Companies apply different security techniques to overcome data breaches as well as to protect private information, consequently guaranteeing system services availability [2,3]; where these techniques include encryption, authentication, firewalls, etc., that have been applied. However, the attackers try to penetrate the network systems to find weak spots that may lead to stealing sensitive information. As a result, these attacks pose a threat to the network system. In general, IDS can identify malicious patterns [4].

An IDS is an important part of cybersecurity that looks for and reacts to unauthorized or harmful activities on a network. Because cyber threats are growing more advanced and common, IDSs serve to detect them and notify administrators. These tools scan network activity, log data and other digital signs to detect things like infections by malware, theft of confidential data or break of security policies. The primary types of IDSs are NIDS, which watch network traffic, and HIDS, which watch over individual devices or hosts. An IDS is effective if it is very

accurate at telling whether a pattern is normal or abnormal. Modern IDSs depend on both known attack patterns and the use of machine learning to find any unusual behaviour. Signature-based systems focus on known dangers, but it is the anomaly-based systems that are best at discovering new threats. Nevertheless, there are issues with each approach, including the chance of false positives and the high cost of running the analysis in real time. Since new cyber threats keep appearing, it is important to make IDS solutions more flexible and smarter to protect digital infrastructure [5]. The main contributions of the proposed method are:

- A new feature selection method is proposed based on RNA encoding and the Raita algorithm to make feature selection more effective for misuse-based intrusion detection systems.
- RNA encoding is used as a new technique to represent network traffic because it mimics the way biological sequences are used for better data analysis.
- The data was simplified, so now the IDS can train quickly and use less computing power.
- Proved that the method is capable of filtering out features that do not add value and in fact lower the IDS's performance.
- A flexible and extendable framework was provided that can be built upon or merged with advanced detection systems.

2. Related Work

In order to evolve the IDS performance, the feature selection of network data is currently a modern research subject [6]. Various feature selection methods for IDS have been proposed in different publications. Mohy-Eddie et al. [7] developed a new IDS using the K-Nearest Neighbours (K-NN) algorithm. They used different feature selection techniques such as principal component analysis, genetic algorithm, and univariate statistical tests. A new IDS is proposed by [8] based on the machine learning method, and then enhances the performance of the proposed method by using modified Singular value decomposition to extract features. A new IDS is suggested by [9] by using Fuzzy Rough sets that are applied for attribute selection and Allen's interval algebra, then applied features selection method that led to a decrease in the false alarm rate results and, at the same time, increased the accuracy result. A new features selection technique is proposed by [10] using DNA encoding and DNA key positions, where the outcomes of the suggested method showed that the proposed method has a fast detection time compared with other systems. A novel feature selection technique proposed in [11] merges statistical importance, deploying standard deviation besides mean and median differences. Zorarpaci [12] proposed a swift wrapper feature selection method and fast ensemble classifier for IDS; the proposed technique led to enhanced runtime performance. A data-driven IDS was proposed by [13] using feature selection and deep learning, where the standard deviation besides the association rule of data mining is applied to remove the redundant features, decrease the computational load, and enhance the accuracy. Fang et al. [14] built a feature selection method for industrial control systems by using a genetic algorithm, and the proposed method decreased the classification complexity and enhanced the accuracy rate.

Al-Bakaa and Al-Musawi [15] proposed a novel anomaly IDS method by using a recurrence quantification analysis to recognize abnormal behaviour in an individual feature extracted, and this led to finding the minimum number of features. [16] suggest a novel feature selection procedure, where the proposed method is based on a genetic algorithm used to determine the optimal feature subsets. A hybrid classification was accomplished using both logistic regression and decision trees. A new IDS technique for the Internet of Things is proposed by [17] based on Linear SVMs, where four feature selection techniques are applied to find the best features; these four methods are Importance Coefficient-, Forward- and Backward-Sequential-, and Correlation Coefficient.

A new DNA encoding idea is proposed by [18] for misuse IDS using the UNSW-NB15 dataset, where the achieved results showed that the suggested DNA encoding led to high DR values for all nine attack types. [19] propose a novel IDS, where this system is generated by the combination of genetic algorithm-based feature selection and multiple support vector machine classifiers, where this system can select features of each attack category rather than for all the attacks. An IDS technique is presented by [20] by using convolutional neural networks and the new UNSW-NB15 dataset, then using a synthetic minority oversampling algorithm and Bayesian Gaussian mixture model to decrease bias toward the majority dataset class. Wei et al. [21] proposition improvement for feature selection in IDS based on a multi-objective immune algorithm then used the neural network method to train the classification model based on the selected feature, and the classification results considered as the target fitness value for each individual. The DNA sequence is used for IDS [22]; the proposed misuse IDS has three stages, which are DNA encoding, shortest tandem repeat extraction, and finally, a matching algorithm based on the Horspool algorithm.

3. Materials and methods

A new feature selection method for the misuse of IDS is proposed. The proposed method consists of five steps: dataset selection, RNA encoding, keys and their locations extraction, features selection, and matching algorithm. Figure 1 depicts the steps of the proposed method.

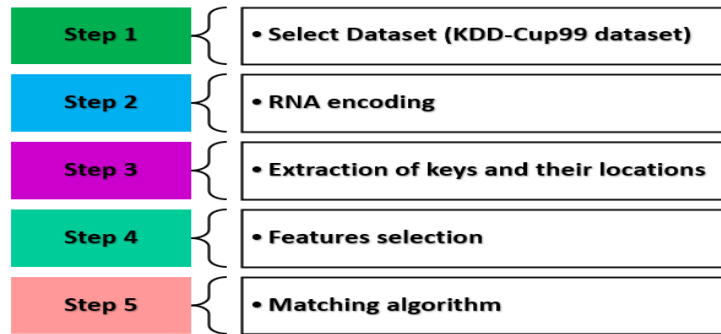


Figure 1. The proposed method (framework)

The framework uses the KDD-CUP99 dataset as a source for the training dataset and then chooses 2000 random records for training purposes. These records are divided into two datasets: one contains only normal records, and the other one contains attack records. The attack records are of types: Denial of Service (DoS) attacks, Probe attacks, Remote to Local (R2L) attacks, and User to Root (U2R) attacks. The next step of the proposed method is to apply RNA encoding to the training datasets. The main goal of this step is to convert all records to RNA characters, and each record has 41 features. These features are listed in Figure 2.

Feature name (position)			
Duration (1)	logged_in (12)	Count (23)	dst_host_same_srv_rate (34)
Protocol_type (2)	num_compromised (13)	srv_count (24)	dst_host_diff_srv_rate (35)
Service (3)	root_shell (14)	error_rate (25)	dst_host_same_src_port_rate (36)
Flag (4)	su_attempted (15)	srv_error_rate (26)	dst_host_srv_diff_host_rate (37)
src_bytes (5)	num_root (16)	error_rate (27)	dst_host_error_rate (38)
dst_bytes (6)	num_file_creations (17)	srv_error_rate (28)	dst_host_srv_error_rate (39)
Land (7)	num_shells (18)	same_srv_rate (29)	dst_host_error_rate (40)
wrong_fragment (8)	num_access_files (19)	diff_srv_rate (30)	dst_host_srv_error_rate (41)
urgent (9)	num_outbound_cmds (20)	srv_diff_host_rate (31)	
hot (10)	is_host_login (21)	dst_host_count (32)	
num_failed_logins (11)	is_guest_login (22)	dst_host_srv_count (33)	

Figure 2. Features of KDDCUP 99

Where the example of RNA encoding is as follows:

Example of a record and its 41 features:

From previous examples, it found that key 1 is obtained from features number 1 & 2, and key 2 is obtained from features number 4 and 5. The selected features by using all extracted keys and their locations are listed in Table 2.

Table 2: Extracted keys and their locations

Feature number	Name
1	duration
2	protocol_type
3	service
4	flag
5	src_bytes
6	dst_bytes

After selecting the best six features instead of all 41 features, the Raita algorithm is used to classify the records as either normal records or attack records based on extracted keys. Where this algorithm is applied twice, once for records with 41 features and the second run for records with selected features. The flowchart for the application of the Raita algorithm is shown in Figure 4.

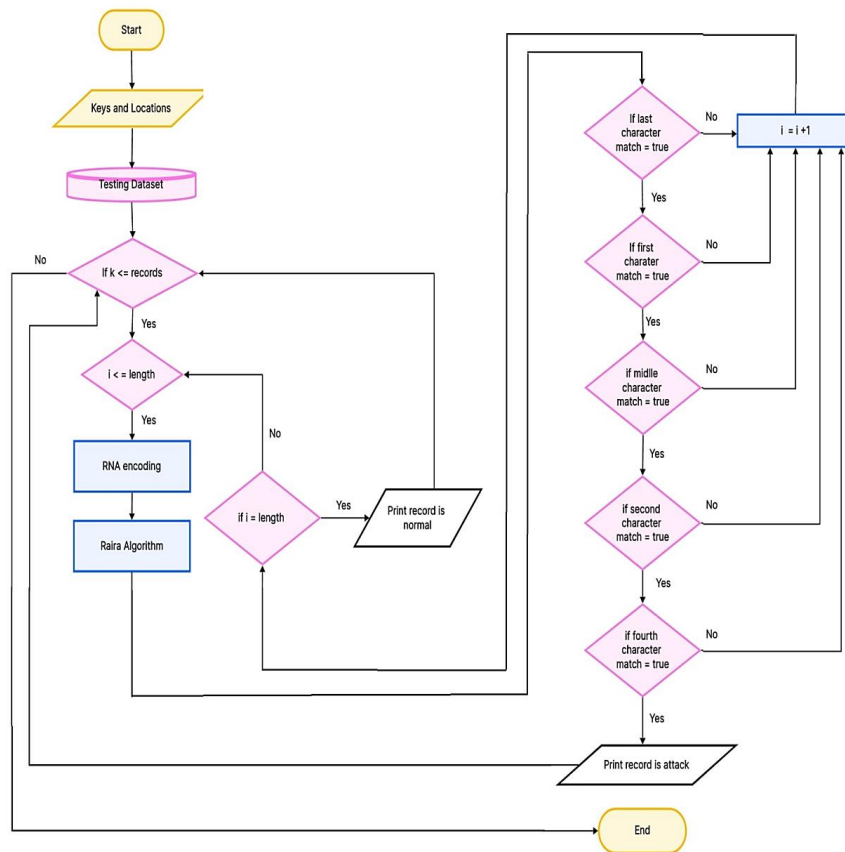


Figure 4. The application of the Raita algorithm

4. Results and discussions

The performance evaluation of the proposed features selection method is computed based on five factors: detection rate (DR), false alarm rate (FAR), accuracy, RNA encoding time, and classification time. The formulae used to calculate these measures are given in the following equations. A new feature selection method for the misuse of IDS is proposed. The proposed method consists of five steps: dataset selection, RNA encoding, keys and their locations extraction, features selection, and matching algorithm. Figure 1 depicts the steps of the proposed method.

$$DR = TP / (TP + FN) \tag{1}$$

$$FAR = FP / (TN + FP) \tag{2}$$

$$Accuracy = (TP + TN) / (TP + TN + FP + FN) \tag{3}$$

The performance is calculated by using a laptop with the following specifications: the OS is Microsoft Windows 10 Professional, the CPU is Intel 2.50GHz, and the memory is 4.00 GB. KDD-Cup99 dataset is used for testing, where 8000 random records are used for this purpose. Firstly, the DR, FAR, and accuracy for all testing dataset records with all features (41 features) are shown in Table 3 and Figure 5.

Table 3: The achieved results before applying the features selection method

Factor	Results
DR	93.31%
FAR	0.002%
Accuracy	94.93%

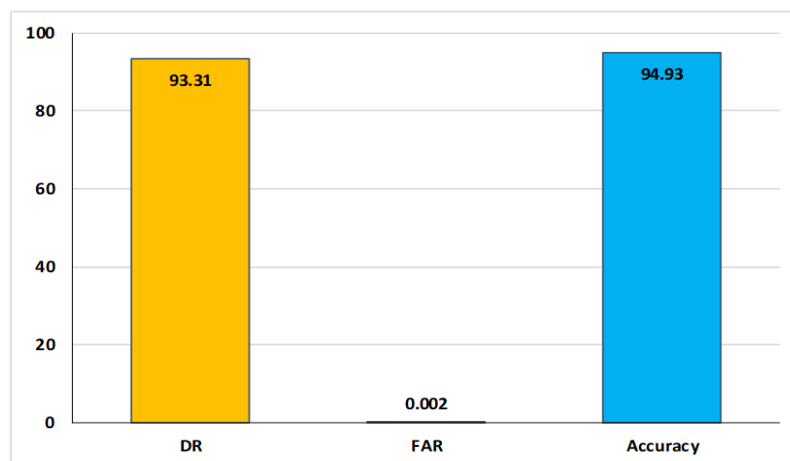


Figure 5. The achieved results before applying the features selection method

The RNA encoding time and classification time for all testing dataset records with all features (41 features) (Table 4 and Figure 6).

Table 4: The obtained time before applied features selection method

Process	Time (seconds)
Encoding Time for 8000 records	1114
Classification Time for 8000 records	27
Encoding Time for 1 record	0.139
Classification Time for 1 record	0.0033

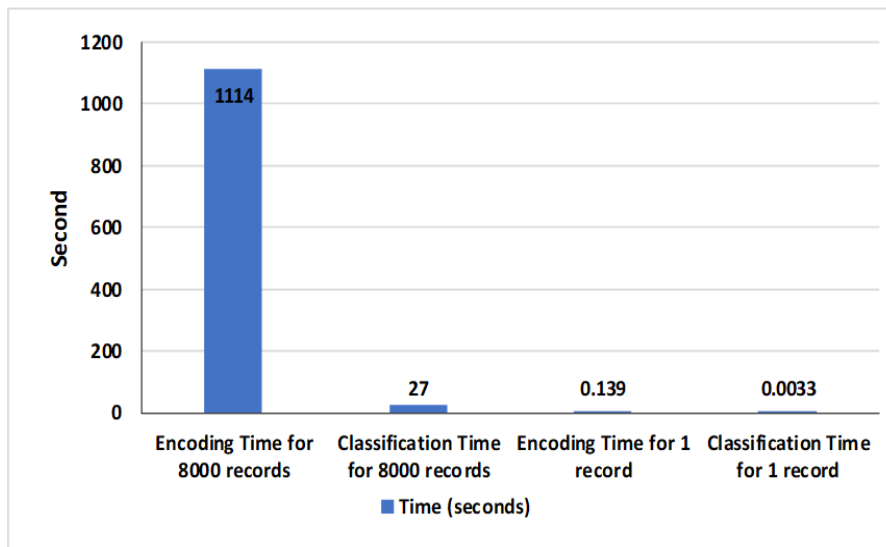


Figure 6. The obtained time before the applied features selection method

As mentioned in Tables 3 and 4, the achieved DR, FAR, accuracy, RNA encoding time, and classification time for all dataset records with all features are equal to 93.31%, 0.002%, 94.93%, 1114 seconds, and 27 seconds respectively.

Secondly, the achieved DR, FAR, and accuracy for all testing dataset records after applying the features selection method are shown in Table 5 and Figure 7.

Table 5: The achieved results after the application of the features selection method

Factor	Results
DR	93.03%
FAR	0.003 %
Accuracy	94.73%

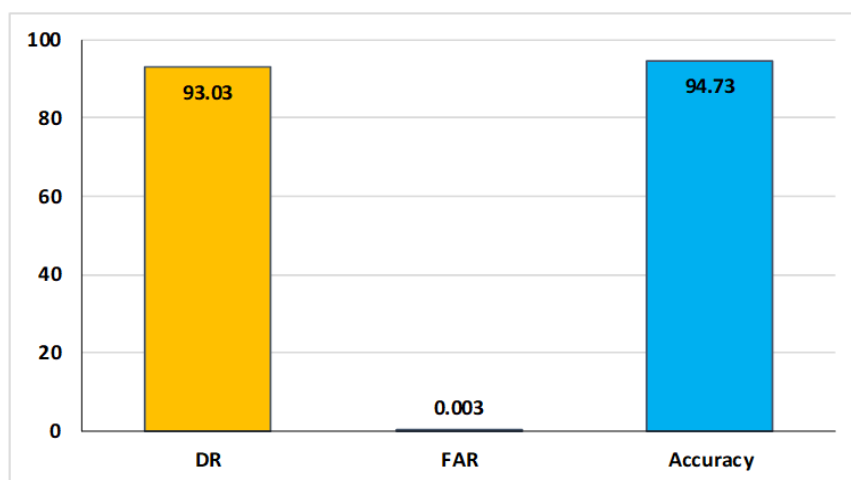
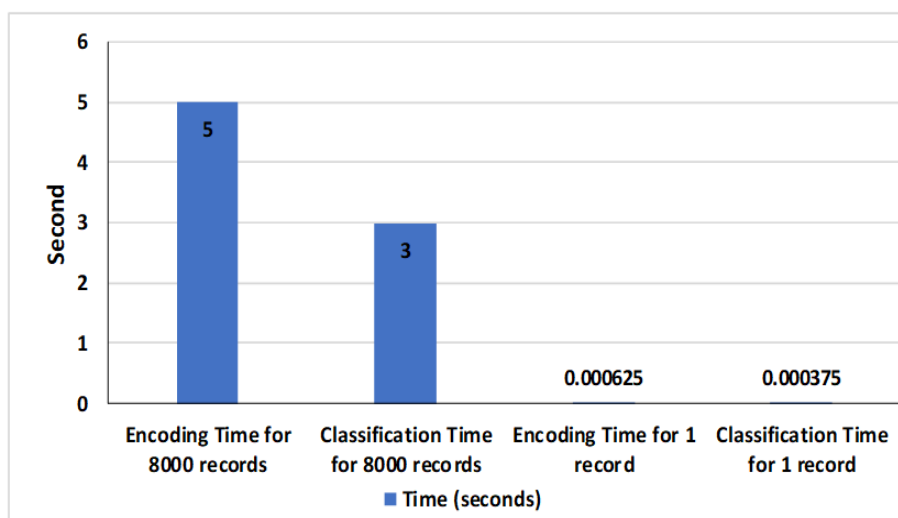


Figure 7. The obtained time after the applied features selection method

On the other hand, the time needed for RNA encoding and classification for all testing dataset records with specific features (6 features) is shown in Table 6 and Figure 8.

Table 6: The obtained time after the application of the features selection method

Process	Time (seconds)
Encoding Time for 8000 records	5
Classification Time for 8000 records	3
Encoding Time per record	0.000625
Classification Time per record	0.000375

**Figure 8.** The obtained time after the applied features selection method

Tables 5 and 6 showed the achieved DR, FAR, accuracy, RNA encoding time, and classification time for all dataset records with all features equal to 93.03%, 0.003%, 94.73%, 5 seconds, and 3 seconds respectively. The achieved results show that the DR, FAR, and accuracy results are not affected within the features selection, where the results are approximately the same. Still, the time for a system based on feature selection is faster than time with full features number, where the RNA encoding time is enhanced from 1114 seconds to 5 seconds only. The classification time becomes 3 seconds instead of 27 seconds.

Table 7 represents the obtained results compared to other IDS methods to evaluate the results obtained by the proposed system.

Table 7: Comparison between the achieved results of proposed methods with published results

Method	Method	Dataset	DR	FAR	Accuracy
Proposed Method	RNA keys and locations	KDDCup99	93.03%	0.003 %	94.73%
Pervez and Farid [23]	Support Vector Machine	KDDCup99	83.4%	9.0%	89.0%
Hooda et al. [24]	Random Forest	KDDCup99	78.5%	9.8%	88.6%
Sheikhi & Kostakos [25]	Genetic Algorithm	KDDCup99	86%	-	94.2%
Sharmila & Nagapadma [26]	Naive Bayes	KDDCup99	81.2%	9.3%	88.4%
Shokeen et al. [27]	Machine Learning	KDDCup99	82.1%	5.3%	87.2%

As mentioned in Table 7, the proposed method achieved the highest DR and Accuracy results over other methods, and these results are equal to 93.03 % and 94.73% respectively. In addition, the obtained FAR by the proposed method is the lowest and equal to 0.003%.

6. Conclusion

This study proposed a new feature selection method using key locations, where this method starts by selecting random records that are used for training and testing, then uses RNA encoding to convert dataset values into RNA characters, and then divides it into blocks. After that, extract the most repeated blocks as classification keys and keep their location. The fourth step is to propose the features selection method, which is done by searching for key locations within which feature numbers are available and then applying the Raita algorithm to perform matching based on extracted keys. For future work, the same procedure can be applied using different datasets, such as the NSL-KDD dataset or the CICIDS2017 dataset. In addition, the proposed method can be enhanced by applying various techniques instead of using a matching algorithm.

Funding: "This research received no external funding"

Conflicts of Interest: "The authors declare no conflict of interest."

References

- [1] O. Ahmed, "Enhancing intrusion detection in wireless sensor networks through machine learning techniques and context awareness integration," *Int. J. Math., Stat., Comput. Sci.*, vol. 2, pp. 244–258, 2024.
- [2] Sharma, N. V., & Yadav, N. S. An optimal intrusion detection system using recursive feature elimination and ensemble of classifiers. *Microprocessors and Microsystems*, 85, 2021.
- [3] Selvakumar, K., Karuppiah, M., SaiRamesh, L., Hafizul Islam, S. K., Hassan, M. M., Fortino, G., & Choo, K. R. Intelligent temporal classification and fuzzy rough set-based feature selection algorithm for intrusion detection system in WSNs. *Information Sciences*, 497, 77-90, 2019.
- [4] Zhou, Y., Cheng, G., Jiang, S., & Dai, M. Building an efficient intrusion detection system based on feature selection and ensemble classifier. *Computer Network*, 174, 2020.
- [5] A. K. Sharma, R. Kumar, and P. Singh, "A Comprehensive Survey on Machine Learning Techniques for Intrusion Detection Systems," *Journal of Information Technology & Software Engineering*, vol. 12, no. 4, pp. 1-10, 2023.
- [6] Hajimirzaei, B., & Navimipour, N. J. Intrusion detection for cloud computing using neural networks and artificial bee colony optimization algorithm. *ICT Express*, 5(1), 56-59, 2019.
- [7] Mohy-eddine, M., Guezzaz, A., Benkirane, S., & Azrou, M. An Intrusion Detection Model using election-Based Feature Selection and K-NN. *Microprocessors and Microsystems*, 2023.
- [8] Turukmane, A. V., & Devendiran, R. M-MultiSVM: An efficient feature selection assisted network intrusion detection system using machine learning. *Computers & Security*, 137, 2024.
- [9] Selvakumar, K., Karuppiah, M., SaiRamesh, L., Islam, S. K. H., Hassan, M. M., Fortino, G., & Choo, K. R. Intelligent temporal classification and fuzzy rough set-based feature selection algorithm for intrusion detection system in WSNs. *Information Sciences*, 497, 77-90, 2019.
- [10] Rashid, O. F., Othman, Z. A., & Zainudin, S. Features Selection for Intrusion Detection System Based on DNA Encoding. In: Piuri, V., Balas, V., Borah, S., Syed Ahmad, S. (eds) *Intelligent and Interactive Computing. Lecture Notes in Networks and Systems*, 67, Springer, Singapore, 2019.
- [11] Thakkar, A., & Lohiya, R. Fusion of statistical importance for feature selection in Deep Neural Network-based Intrusion Detection System. *Information Fusion*, 90, 353-363, 2023.
- [12] Zorarpaci, E. A fast intrusion detection system based on swift wrapper feature selection and speedy ensemble classifier. *Engineering Applications of Artificial Intelligence*, 133, 2024.
- [13] Zhang, L., Liu, K., Xie, X., Bai, W., Wu, B., & Dong, P. A data-driven network intrusion detection system using feature selection and deep learning. *Journal of Information Security and Applications*, 78, 2023.
- [14] Fang, Y., Yao, Y., Lin, X., Wang, J., & Zhai, H. A feature selection based on genetic algorithm for intrusion detection of industrial control systems. *Computers & Security*, 139, 2024.

- [15] Al-Bakaa, A., & Al-Musawi, B. A new intrusion detection system based on using non-linear statistical analysis and features selection techniques. *Computers & Security*, 122, 2022.
- [16] Kunhare, N., Tiwari, R., & Dhar, J. Intrusion detection system using hybrid classifiers with meta-heuristic algorithms for the optimization and feature selection by genetic algorithm. *Computers and Electrical Engineering*, 103, 2022.
- [17] Azimjonov, J., & Kim, T. Designing accurate lightweight intrusion detection systems for IoT networks using fine-tuned linear SVM and feature selectors. *Computers & Security*, 137, 2024.
- [18] Rashid, O. F. DNA Encoding for Misuse Intrusion Detection System based on UNSW-NB15 Data Set. *Iraqi Journal of Science*, 61(12), 3408–3416, 2020.
- [19] Waad, F., & Mohammed, I. J. Hybrid CNN-SMOTE-BGMM Deep Learning Framework for Network Intrusion Detection using Unbalanced Dataset. *Iraqi Journal of Science*, 64(9), 4846-4864, 2023.
- [20] Rashid, O. F., Othman, Z. A., & Zainudin, S. Matching algorithms for intrusion detection system based on DNA encoding. *Journal of Theoretical and Applied Information Technology*, 96 (24), 8410 - 8420, 2018.
- [21] Wei, W., Chen, S., Lin, Q., Ji, J., & Chen, J. A multi-objective immune algorithm for intrusion feature selection. *Applied Soft Computing*, 95, 2020.
- [22] Vijayanand, R., Devaraj, D., & Kannapiran, B. Intrusion detection system for wireless mesh network using multiple support vector machine classifiers with genetic-algorithm-based feature selection. *Computers & Security*, 77, 304-314, 2018.
- [23] Pervez, M. S., & Farid, D. M. Feature selection and intrusion classification in NSL-KDD cup 99 dataset employing SVMs. The 8th *International Conference on Software, Knowledge, Information Management and Applications (SKIMA 2014)*, Dhaka, Bangladesh, 1-6, 2014.
- [24] Hooda, M., Babu, J., Vamsi, P. S., & Gopakumar, G. An Improved Intrusion Detection System Based on KDD Dataset Using Feature Ranking and Data Sampling. *2020 International Conference on Communication and Signal Processing (ICCSP)*, Chennai, India, 1128-1132, 2020.
- [25] Sheikhi, S., & Kostakos, P. A Novel Anomaly-Based Intrusion Detection Model Using PSO-GWO-Optimized BP Neural Network and GA-Based Feature Selection. *Sensors*, 2022.
- [26] Sharmila, B. S., & Nagapadma, R. Intrusion Detection System using Naive Bayes algorithm. *2019 IEEE International WIE Conference on Electrical and Computer Engineering, India*, 1-4, 2019.
- [27] Shokeen, A., Yadav, N., & Sisaudia, V. Performance analysis of different machine learning algorithms for intrusion detection on KDD-CUP-99 dataset. *AIP Conference Proceeding*, 2024.