# Machine Learning for Intelligent Data Analysis and Automation in Cybersecurity: Current and Future Prospects

Iqbal H. Sarker[1,2*]

**Abstract** Due to the digitization and Internet of Things revolutions, the present electronic world has a wealth of cybersecurity data. Efficiently resolving cyber anomalies and attacks is becoming a growing concern in today's cyber security industry all over the world. Traditional security solutions are insufficient to address contemporary security issues due to the rapid proliferation of many sorts of cyber-attacks and threats. Utilizing artificial intelligence knowledge, especially *machine learning* technology, is essential to providing a dynamically enhanced, automated, and up-to-date security system through analyzing security data. In this paper, we provide an extensive view of *machine learning* algorithms, emphasizing how they can be employed for *intelligent data analysis* and *automation* in cybersecurity through their potential to extract valuable insights from cyber data. We also explore a number of potential *real-world use cases* where data-driven intelligence, automation, and decision-making enable next-generation cyber protection that is more proactive than traditional approaches. The *future prospects* of machine learning in cybersecurity are eventually emphasized based on our study, along with relevant research directions. Overall, our goal is to explore not only the current state of machine learning and relevant methodologies but also their applicability for future cybersecurity breakthroughs.

**Keywords** cybersecurity; machine learning; deep learning; artificial intelligence; data-driven decision making; automation; cyber analytics; intelligent systems

## 1 Introduction

We live in the digital age, which, like anything else, has its upsides and downsides. The main drawback is the security risk [96] [105]. As more of our sensitive infor-

[1] Department of Computer Science and Engineering, Chittagong University of Engineering & Technology, Chittagong 4349, Bangladesh.

[2] Swinburne University of Technology, Melbourne, VIC-3122, Australia.

*Correspondece: msarker@swin.edu.au (Iqbal H. Sarker)
ORCID: https://orcid.org/0000-0003-1740-5517

mation transfers to the digital arena, security breaches are becoming more common and catastrophic. Cyber-criminals are growing more adept in their attempts to avoid detection, and many newer malware kits are already incorporating new ways to get out of antivirus and other threat detection systems. Cybersecurity, on the other hand, is at a crossroads, and future research efforts should be focused on cyber-attack prediction systems that can foresee important scenarios and consequences, rather than depending on defensive solutions and focusing on mitigation. Systems that are based on a complete, predictive study of cyber risks are required all around the world. The key functionalities in cybersecurity such as *prediction*, *prevention*, *identification or detection* as well as corresponding *response* should be done *intelligently and automatically*. Artificial intelligence (AI), which is based primarily on *Machine Learning (ML)* [94] [91], is capable of recognizing patterns and predicting future moves based on prior experiences, thereby preventing or detecting potentially malicious activity, which is the primary focus of this study.

ML is one of the most popular current technologies in the fourth industrial revolution ($4IR$ or Industry 4.0) [120] [111] because it allows systems to learn and improve from experience without having to be explicitly programmed [104] [112]. In the cyber security area, machine learning can play a vital role in capturing insights from data. Cybersecurity data can be organized or unstructured, and it can originate from a variety of sources, as explained in Section 3. Intrusion detection, cyber-attack or anomaly detection, phishing or malware detection, zero-day attack prediction, and other intelligent applications can be built by extracting insights from these data. The demand for cybersecurity and protection against cyber anomalies and various sorts of attacks, such as unauthorized access, denial-of-service (DoS), phishing, malware, botnet, spyware, worms, etc. has risen dramatically in recent days. Thus, real-world cyber applications require *intelligent data analysis* tools and approach capable of extracting insights or meaningful knowledge from data in a timely and intelligent manner. Security researchers believe they can utilize attack pattern recognition or detection methods to provide protection against future attacks.

Machine learning technologies are thus used to intelligently analyze cybersecurity data and provide a dynamically upgraded and up-to-date security solution. Learning algorithms can be divided into four categories: supervised, unsupervised, semi-supervised, and reinforcement learning [94]. The nature and quality of the data as well as the effectiveness of the learning algorithms, in general, impact the productivity and efficiency of a machine learning solution. In this paper, we explore various types of machine learning techniques such as classification and regression analysis, security data clustering, rule-based modeling, as well as deep learning approaches, all of which fall within the broad category of machine learning and are capable of building cybersecurity models for different purposes. In addition, we also explore *adversarial* machine learning, which is the study of how machine learning algorithms are attacked and how they are defended. It is challenging to find a suitable learning algorithm for the intended application in a particular domain. This is because different learning algorithms have distinct functions, and even the results of related learning algorithms might vary based on the properties of the input. Therefore, it's crucial to understand the fundamentals of various machine learning algorithms and how they apply to a range of real-world application domains, such as detecting malicious activity, predicting data breaches, intrusion detection, and prevention, as outlined in Section 5.

Based on the aforementioned importance of machine learning, we provide a comprehensive view of machine learning algorithms that can be utilized for intelligent data analysis and automation in cybersecurity due to their ability to capture insights from data in the cyber security domain in this study. Thus data-driven intelligent decision-making and automation allow the next-generation cyber-defense that is more proactive than current approaches. Thus the study's primary strength is exploring the applicability of different machine learning algorithms in the numerous cyber application domains, summarized in Section 5. Therefore, the purpose of this paper is to provide a point of reference for academicians and practitioners from the industry who are interested in learning about, investigating, and creating data-driven automated and intelligent systems in the area of cybersecurity utilizing machine learning techniques.

The key contributions of this paper are listed as follows:

– To define the scope of our study by offering a dynamically improved and automated up-to-date security system using machine learning technologies.
– To provide a comprehensive understanding of machine learning algorithms that can be applied in cybersecurity for intelligent data analysis and automation.
– To explore the applicability of various machine learning approaches in a variety of real-world scenarios in the context of cybersecurity, where data-driven intelligent decision-making and automation allow the next-generation cyber-defense that is more proactive than traditional approaches.
– To emphasize the future prospects of machine learning in cybersecurity, along with relevant research directions.

The rest of the paper is laid out in the following manner. Section 2 motivates and defines the scope of our research by describing why machine learning is relevant in today's cybersecurity research and applications. In Section 3, we look at cybersecurity data in-depth, and in Section Section 4, we go through different machine learning algorithms in detail. In Section 5, many machine learning algorithms-based application fields are explored and summarized. We highlight the future prospects of machine learning in cybersecurity, as well as important research directions in Section 6, and finally, Section 7 concludes this paper.

## 2 Why Machine Learning in Today's Cybersecurity Research and Applications?

Automation is becoming a key tool for overwhelmed security personnel as today's diverse cyber threats become more widespread, sophisticated, and targeted. Malware, phishing, ransomware, denial-of-service (DoS), zero-day attacks, etc. are common as shown in Figure 1. This is because most defense measures are not flawless, and many of today's detection approaches rely on an analyst's manual investigation and decision-making to uncover advanced threats, malicious user behavior, and other major associated risks. When it comes to recognizing and predicting specific patterns, machine learning outperforms humans. Security decisions and policy adaptations have failed to meet security requirements in highly dynamic and sophisticated network systems. Intelligent decision-making utilizing machine learning technology to achieve automation has become possible.
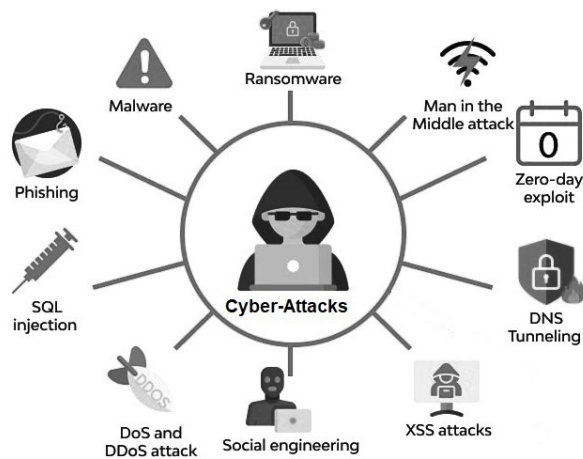
Iqbal H. Sarker[1,2]*



Fig. 1: Several common attacks or threats in the context of cybersecurity.

In Figure 2, we have plotted the global statistical impact of machine learning and cybercrime over the previous five years, where the x-axis indicates timestamp data and the y-axis represents the equivalent value. We can see from the graph that cybercrime is on the rise all over the world. Thus protecting an information system, especially one that is connected to the Internet, from various cyber-threats, attacks, damage, or unauthorized access is a crucial issue that must be addressed immediately. Machine learning techniques, with their outstanding learning capabilities from cyber data, can play a vital part in addressing these issues in accordance with today's needs, which is also a popular technology in recent days, as shown in Figure 2.
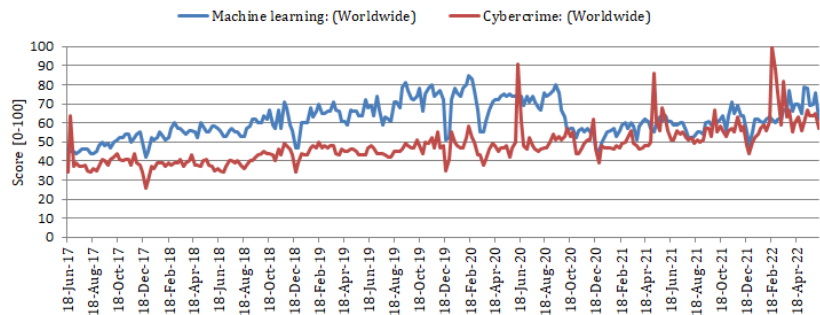


Fig. 2: The global statistical impact of machine learning and cybercrime over time, with the x-axis representing the timestamp information and the y-axis representing the equivalent value, on a scale of 0 (min) to 100 (max).

ML has the potential to revolutionize the planet as well as humans' daily lives through its automated capabilities and ability to learn from experience. All around the world, systems that are based on a comprehensive, predictive analysis of cyber
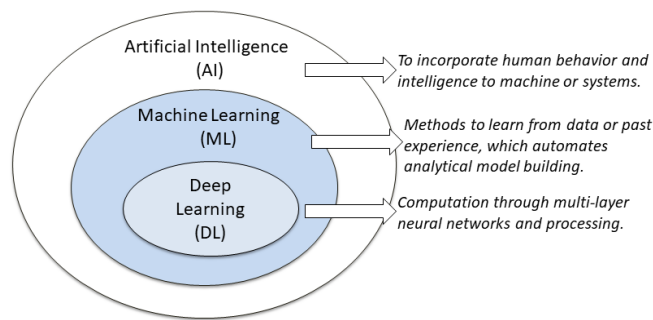
Fig. 3: An illustration of Machine Learning (ML) including Deep Learning (DL) relative to Artificial Intelligence (AI).

risks are expected. Prediction, prevention, identification, and response are all crucial cybersecurity functions that should be handled intelligently and automatically. Thus the knowledge of artificial intelligence (AI) [102], which is mostly based on machine learning (ML), is capable of recognizing patterns and predicting future moves based on recent experiences, thereby preventing or detecting potentially malicious behavior. We also explore machine learning compared with deep learning and artificial intelligence in Figure 3. ML is a subset of AI and DL is a subset of ML, according to Figure 3. In general, AI [95] combines human behavior and intelligence into machines or systems, whereas ML [94] is a method of learning from data or experience that automates the creation of analytical models in a particular application domain, e.g., cybersecurity according to our focused area.

Thus machine learning can be considered a key AI technology, a frontier for artificial intelligence that can be utilized to develop intelligent systems and automate processes, in which we are interested in the context of cybersecurity. Therefore, to have a real influence on increasing an organization's ability to recognize and respond to emerging and ever-evolving cyber threats, it's necessary to deploy machine learning appropriately.

## 3 Understanding Cybersecurity Data

As machine learning algorithms create models from data, understanding cybersecurity data is essential for intelligent analysis and decision-making. Cybersecurity datasets are often collections of information records that contain a variety of attributes or features, as well as related facts, on which machine learning-based modeling is based. A sample of features from the KDD'99 cup dataset [4] is shown in Table 1. Thus understanding the nature of cybersecurity data, which includes various types of cyberattacks as well as key features, is important. Intrusion detection, malware detection, and spam detection are just a few of the datasets available in the realm of cybersecurity [104].

For instance, the KDD'99 Cup dataset [4], the most widely used data set including 41 features attributes and a class identification, with attacks divided into four categories: denial of service (DoS), remote-to-local (R2L) intrusions, and user-to-remote (U2R) intrusions, and PROB as well as conventional data. NSL-KDD [119], an updated version of the KDD'99 cup dataset that removes redundant

Table 1: An example of features of KDD'99 cup dataset [4].

| No. | Features | Types | No. | Features | Types |
|---|---|---|---|---|---|
| 1 | duration | Continuous | 22 | is_guest_login | Symbolic |
| 2 | protocol_type | Symbolic | 23 | count | Continuous |
| 3 | service | Symbolic | 24 | srv_count | Continuous |
| 4 | flag | Symbolic | 25 | serror_rate | Continuous |
| 5 | src_bytes | Continuous | 26 | srv_serror_rate | Continuous |
| 6 | dst_bytes | Continuous | 27 | rerror_rate | Continuous |
| 7 | Land | Symbolic | 28 | srv_rerror_rate | Continuous |
| 8 | wrong_fragment | Continuous | 29 | same_srv_rate | Continuous |
| 9 | urgent | Continuous | 30 | diff_srv_rate | Continuous |
| 10 | hot | Continuous | 31 | drv_diff_host_rate | Continuous |
| 11 | num_failed_logins | Continuous | 32 | dst_host_count | Continuous |
| 12 | logged_in | Symbolic | 33 | dst_host_srv_count | Continuous |
| 13 | num_compromised | Continuous | 34 | dst_host_same_srv_rate | Continuous |
| 14 | root_shell | Continuous | 35 | dst_host_diff_srv_rate | Continuous |
| 15 | su_attempted | Continuous | 36 | dst_host_same_src_port_rate | Continuous |
| 16 | num_root | Continuous | 37 | dst_host_srv_diff_host_rate | Continuous |
| 17 | num_file_creations | Continuous | 38 | dst_host_serror_rate | Continuous |
| 18 | num_shells | Continuous | 39 | dst_host_srv_serror_rate | Continuous |
| 19 | num_access_files | Continuous | 40 | dst_host_rerror_rate | Continuous |
| 20 | num_outbound_cmds | Continuous | 41 | dst_host_srv_rerror_rate | Continuous |
| 21 | is_host_login | Symbolic | | | |

records. Thus a machine learning classification-based security model based on the dataset will not be skewed towards more frequent records. For evaluating computer network intrusion detection systems, the cyber systems and technologies group at MIT Lincoln Laboratory collects and publishes datasets containing traffic and attacks [59]. CAIDA'07 is a dataset that contains anonymized traces of DDoS attack traffic recorded in 2007, with the attack mostly consisting of flooding traffic of SYN, ICMP, and HTTP [1]. ISCX'12 represents network traffic generated in a real-world physical test environment while containing centralized botnets generated by Canadian Institute for Cybersecurity [2]. CTU-13, a botnet traffic dataset collected at CTU University in the Czech Republic containing thirteen separate malware captures, including Botnet, Normal, and Background traffic [3]. UNSW-NB15 was founded in 2015 at the University of New South Wales containing 49 features and roughly 257,700 records, which represent nine various forms of current attacks, including denial-of-service attacks [70]. DDoS intrusion detection system developed by a group of Japanese network research and academic institutions [41]. For the aim of network forensic analytics in the Internet of Things, another dataset Bot-IoT includes legitimate and simulated IoT network traffic, as well as various cyberattacks [49].

A variety of such datasets available on the Internet, along with their various attributes and cyberattacks, could be used to emphasize their usage in various cyber applications through machine learning-based analytical modeling. Analyzing and processing these security elements efficiently, constructing a target machine

learning-based security model based on the needs, and eventually, data-driven decision-making might all help deliver intelligent cybersecurity services. A variety of machine learning approaches, which are briefly mentioned in Section 4, can be employed to achieve our goal.

## 4 Machine Learning Tasks and Algorithms in Cybersecurity

Machine learning is typically known as a methodological approach that automates the formation of analytical models, focusing on the use of data and algorithms to mimic the way humans learn while gradually improving accuracy. A key component of the development of machine learning algorithms and the enhancement of their performance is the loss function [125]. A broad structure for a machine learning-based prediction model is shown in Figure 4, with the model being trained from historical security data containing benign and malware in phase 1, and the output is generated for new test data in phase 2. As shown in Figure 5, machine learning is typically divided into four categories: supervised, unsupervised, semi-supervised and reinforcement learning [94]. Within the broad field of machine learning, we first explore classification and regression analysis, security data clustering, as well as rule-based modeling. We have also explored deep learning methodologies in this section, according to their capabilities to solve real-world issues in cybersecurity.
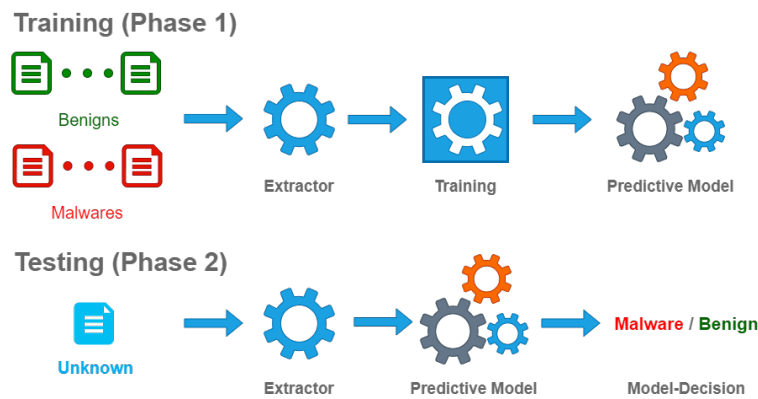


Fig. 4: The training and testing phases of a machine learning-based predictive model (i.e., benign or malware).

4.1 Classification and Regression Analysis in Cybersecurity

Both classification and regression approaches are well-known as supervised learning and are frequently employed in the field of machine learning. Many classification algorithms have been proposed in the machine learning and data science literature that can be used for intelligent data analysis to solve various real-world issues in the context of cybersecurity. The decision tree is the most powerful and
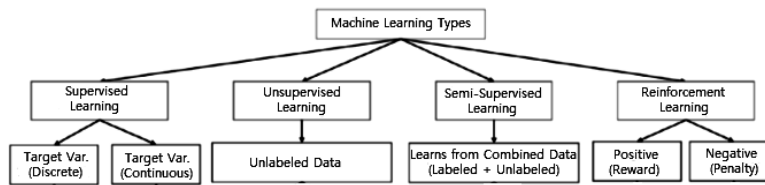
Fig. 5: Traditional machine learning types.

widely used tool for classification and prediction. For instance, an intelligent intrusion detection model for cyber security has been proposed, which is based on the notion of decision trees and takes into account the ranking of security features [8]. In [124] the authors offer a gradient boosting decision tree based on network events records for detecting cyber security concerns. Authors in [81] present an anomaly-based intrusion detection system for the smart grid based on the cart decision tree. Typically ID3 [79], C4.5 [80], and CART [16] are well-known DT algorithms in the area of machine learning. Furthermore, Sarker et al. recently proposed BehavDT [99] that has been designed based on behavior analysis, and IntrudTree [97] by taking into account a generalized decision tree with selected security features, that could be employed for a better outcome in the relevant application domains. K-nearest neighbors [7], support vector machines [46], navies Bayes [43], adaptive boosting [29], logistic regression [53], etc. are also popular techniques in the area. An optimal detection of a phishing attack using SCA-based K-NN has been presented in [69]. To profile abnormal behavior [54] or detect android malware [33] the support vector machine classification technique can be employed. To detect anomalies [116] a naive Bayes based classification model is useful while a logistic regression-based method to detect malicious botnets [14] [77] .

Ensemble learning is another popular approach, typically known as a general meta approach to machine learning that combines the predictions of numerous models to improve predictive performance. For instance, a random forest [15] technique consists of multiple decision trees is used to detect anomalies [18] [76]. Similarly, detecting denial of service attack [24], intrusions [87] [68], smart city anomalies [10] [96] the most popular forest technique can be used. The authors in [139] also offer a Bayesian network-based ensemble learning solution for detecting XSS attacks with domain knowledge and threat intelligence. In [83], authors studied a stacked ensemble learning model for intrusion detection in wireless networks, in which random forest and gradient boost are utilized as base learners for identifying attacks.

A regression model, on the other hand, is beneficial for statistically predicting cyberattacks or predicting the impact of an attack, such as worms, viruses, or other malicious software [39]. Regression techniques could be effective for a quantitative security model, such as phishing in a specific period or network packet parameters [102]. Linear, polynomial, Ridge, Lasso, and other prominent regression techniques [94] can be utilized to develop a quantitative security model based on their machine learning principles. For example, authors in [51] use a linear regression-based model to identify the source of a cyber attack, and [31] use multiple regression analysis to connect human traits with cybersecurity activity intents. Because of the enormous dimensionality of cyber security data, regression regu-

larization methods such as Lasso, Ridge, and ElasticNet can improve security breaches analysis [32]. The authors in [107] look into the profitability of trading strategies supported by ML approaches as well as the predictability of returns on the most well-liked cryptocurrencies. Regression models are employed in their research to forecast the returns of the dependent variable, which in this case is a cryptocurrency, and classification models are utilized to produce binary buy or sell trading recommendations.

Therefore, we can conclude that classification techniques can be used to build the prediction and classification model [91] utilizing relevant data in the domain of cyber security, whereas regression techniques are primarily used to determine the model's impact [39] by determining predictor strength, time-series causes, or the effect of the relations, taking into account the security attributes and the outcome.

4.2 Clustering in Cybersecurity

Clustering, which is classified as unsupervised learning, is another common activity in machine learning for processing cybersecurity data. It can cluster or group a set of data points based on measures of similarity and dissimilarity in security data from a variety of sources. Thus clustering may aid in the uncovering of hidden patterns and structures in data, allowing irregularities or breaches to be detected. Clustering data can be done using partition, hierarchy, fuzzy theory, density, and other perspectives [131].

The popular concepts of clustering algorithms include K-means [63], K-medoids [88], single linkage [113], complete linkage [114], agglomerative clustering, DB-SCAN, OPTICS, Gaussian Mixture Model [131], etc. In [101] Sarker et al. proposed a bottom-up clustering algorithm by taking into account behavior analysis. Various cybersecurity issues can be solved using these clustering strategies. For example, in profiling the anomalous behavior of devices, the k-Means algorithm is utilized [54]. To detect outlier or noisy occurrences in data, authors in [90] use a dynamic threshold-based technique. In intrusion detection, Liu et al. [60] use a fuzzy clustering technique. Overall, clustering methods are beneficial for extracting relevant insights or knowledge from system log data for cybersecurity applications summarized by landauer et al. [52].

Clustering techniques can help solve a variety of security problems, such as outlier detection, anomaly detection, signature extraction, fraud detection, cyber-attack detection, and so on, by revealing hidden patterns and structures in cyber-security data and measuring behavioral similarity or dissimilarity. Thus clustering-based unsupervised learning could be a significant topic to explore more for future research in the context of next-generation cybersecurity.

4.3 Rule-based Modeling in Cybersecurity

A rule-based system that extracts rules from data can be used to simulate human intelligence, which is defined as a system that uses rules to make an intelligent decision [92]. Thus by learning security or policy rules from data, rule-based systems can play a vital role in cybersecurity [102]. In the discipline of machine learning, association rule learning is a popular approach to detecting associations or rules

among a set of available characteristics in a security dataset [5]. Several types of association rules have been proposed in this field, including frequent pattern based [6] [38], [62], logic-based [27], tree-based [35], fuzzy-rules [118], belief rule [140], and so on. AIS [5], Apriori [6], Apriori-TID and Apriori-Hybrid [6] as well as Eclat [137], RARM [21], FP-Tree [35] are some of the rule learning techniques that can be used to solve cybersecurity problems and intelligent decision making due to their rule learning capabilities from data. In [108], for example, an association rule-mining algorithm-based network intrusion detection has been described. Additionally, fuzzy association rules are employed to construct a rule-based intrusion detection system [118]. In [71], an FP-tree association rule-based study was carried out to investigate malware behaviors. A belief rule-based anomaly detection under uncertainty has been presented in [121]. Such belief rules can also be used to build an expert system modeling in a particular application area depending on the problem nature [140].

A rule-based technique is simple to implement, but it has a high temporal complexity since it generates a large number of associations or common patterns based on the support and confidence values, making the model complex [6] [117]. This problem could be mitigated with a good association model. For example, in our previous publication, Sarker et al. [103], we propose a rule learning strategy that effectively identifies non-redundant and dependable association rules, which could be useful in the realm of cybersecurity as well. To solve more complicated challenges in cybersecurity, the rules can be utilized to develop knowledge-based systems or rule-based expert systems [102] [100]. Each of these systems is made up of a set of policy rules that define the scope of what types of activities should be permitted on a network, with each rule clearly allowing or disallowing particular activities. Future zero-day attacks that use rule-driven controls or filters are even blocked by security policy monitoring. Thus various types of security rule-based models can be explored more for future research and deployment according to the needs and the nature of the problem in the context of cybersecurity.

4.4 Deep Learning in Cybersecurity

In many situations, deep learning (DL), a subset of machine learning that emerged from the Artificial Neural Network (ANN), outperforms conventional machine learning algorithms, especially when learning from huge security datasets. The ANN is a type of computational architecture for data-driven learning that incorporates different processing layers such as input, hidden, and output layers into a single network [34]. As deep learning techniques are knowledge-capture techniques in deep architecture, they may learn from cybersecurity data, e.g., intrusion detection, over numerous layers and are known as hierarchical learning methods [26]. According to the taxonomy presented in our earlier paper by Sarker et al. [93], deep learning techniques can be broadly categorized into three types: supervised or discriminative learning, e.g., CNN; unsupervised or generative learning, e.g., Auto-encoder; and hybrid learning combining both with other applicable techniques, can be used to address today's cybersecurity issues. For instance, an intrusion detection model based on the NSL-KDD dataset [28], malware analysis [45], and detecting malicious botnet traffic [40] are all constructed using the MLP network. A CNN-based deep learning model can be used to detect intru-

Table 2: A summary of machine learning tasks in the domain of cybersecurity.

| Used Technique | Purpose | References |
|---|---|---|
| SVM | Classifying cyber-attacks known as DoS, U2R, R2L, and Probing | Kotpalliwar et al. [50] |
| SVM | Selecting security features, detecting and classifying intrusions | Pervez et al. [74], Yan et al. [134], Li et al. [57], Raman et al. [84] |
| SVM-PSO | Developing intrusion detection model | Saxena et al. [106] |
| FCM clustering, ANN and SVM | Building network intrusion detection system and modeling | Chandrasekhar et al. [17] |
| KNN | To build intrusion detection system | Shapoorifard et al. [109], Vishwakarma et al. [123] |
| KNN | Reducing the false alarm rate | Meng et al. [66] |
| SVM and KNN | Building intrusion detection system | Dada et al. [19] |
| K-means and KNN | Building intrusion detection system | Sharifi et al. [110] |
| KNN and Clustering | Building intrusion detection system | Lin et al. [58] |
| Decision Tree | Selecting security features and building an effective network intrusion detection system | Radoglou et al. [81], Malik et al. [64], Relan et al. [86], Rai et al. [82], Sarker et al. [97], Puthran et al. [78] |
| Decision Tree and KNN | To detect anomaly intrusions | Balogun et al. [13] |
| Genetic Algorithm and Decision Tree | Solving the issue of small disjunct while building a tree-based IDS | Azad et al. [12] |
| Decision Tree and ANN | Intrusion detection system | Jo et al. [42] |
| Ensemble learning | Detecting XSS attacks | Zhou et al. [139] |
| RF | Detect cyber anomalies | Chang et al. [18], Alrashdi et al. [10] |
| RF | Detecting DoS attack | Doshi et al. [24] |
| RF | Intrusion detection systems | Resende et al. [87], Mohamed et al. [68] |
| Association Rule | Building and effective network IDS | Tajbakhsh et al. [118] |
| Behavior Rule | Building IDS for safety critical medical cyber physical systems | Mitchell et al. [67] |
| NBC | Detecting anomalies | Swarnkar et al. [116] |
| LR | Detecting malicious botnets | Prokofiev et al. [77], Bapat et al. [14] |
| LR | Predicting the impact of cyber-attacks | Jaganathan et al. [39] |
| Regression Regularization | Handling high dimensions of security data | Hagos et al. [32] |
| PCA | Handling high dimensionality security data | Hoang et al. [37] |
| fuzzy cluster | Building IDS | Liu et al. [60] |
| Semi-supervised | Distributed threat detection system | Rathore et al. [85] |
| FP-tree | Analyzing and detecting malwares | Ozawa et al. [71] |
| Deep Learning Recurrent, RNN, LSTM | Detecting and classifying anomaly intrusions and attacks | Alrawashdeh et al. [11], Yin et al. [135], Kim et al. [47], Almiani et al. [9] |
| Deep Learning Convolutional | Classifying malware traffics | Kolosnjaji et al. [48], Wang et al. [127] |
| multi-CNN | Building IDS | Li et al. [56] |
| LSTM+CNN | Detecting and mitigating phishing and Botnet attacks | Parra et al. [72] |
| Transformer | Autonomous cyberbullying detection | Pericherla et al. [73] |
| Q-Learning | DDoS detection | Liu et al. [61] |

sions such as denial-of-service (DoS) attacks [115], malware detection [133], and android malware detection [65]. Recurrent connections can aid neural networks in detecting security risks when the threat's behavior patterns are time-dependent. In the sphere of security, an LSTM model-based recurrent network can be utilized for a variety of tasks, including intrusion detection [47], detecting and classifying malicious apps [122], backdoor attack classification [20] and so on.

In contrast, generative learning techniques are frequently employed for feature learning, data generating, and representation [22] [23]. Deep neural network algorithms for unsupervised or generative learning such as Autoencoders, Generative Adversarial Networks, and Deep Belief Networks as well as their variants, can be employed to address cybersecurity problems as well. Several examples are - auto-encoder based malware [126] as well as intrusion detection [132], deep belief network-based intrusion detection model [128] and so on. Moreover, in [55], a novel GAN-based adversarial-example attack method was constructed that outperformed the leading technique by a significant amount. A method to improve botnet detection models using generative adversarial networks Bot-GAN was provided in [136], which increases detection effectiveness and lowers the probability of false positives.

Hybrid network models, such as the ensemble of learning models, e.g., CNN and RNN, or others with their optimization can also be used to detect cyberattacks, such as malware detection [133], phishing, and Botnet attack detection and mitigation [72]. In addition, authors in [73] describe a transformer network-based word embeddings approach for autonomous cyberbullying detection. A robust transformer-based intrusion detection system has been presented in [129]. The authors in [138] provide a generative adversarial network for anomaly detection using multiple transformer encoders. Overall, due to the capabilities to effectively learn from a large amount of security data at several levels, deep learning models and their variants as well as their hybridization could also play a key role in the field of cybersecurity.

4.5 Semi-Supervised and Other learning techniques in Cybersecurity

Semi-supervised learning is a significant part of machine learning processes because it increases and enhances the capabilities of machine learning systems by operating on both labeled and unlabeled data. This is a substantial advantage over a fully supervised model, which requires all data to be labeled. As a result, cost and time reductions are associated with semi-supervised learning. When compared to an unsupervised model, a supervised model can save computational resources and improve the model's accuracy when utilized with even a minimal amount of labeled data.

Merging clustering and classification algorithms could be an example of semi-supervised learning, where clustering algorithms are unsupervised machine learning methodologies for grouping data based on similarity. In [85] authors combine a semi-supervised Fuzzy C-Means with the extreme learning classifier to create a semi-supervised learning-based distributed threat detection system for IoT. An intrusion detection system based on semi-supervised learning with an adversarial auto-encoder has been presented in [36]. Authors in [30] present a semi-supervised transfer learning malware categorization for the cloud. In many cases, security

feature engineering and optimization are regarded as crucial issues in the cyber threat landscape for a successful cyber security system based on a machine learning methodology. The reason for this is that security characteristics and associated data have a direct impact on machine learning-based security models, therefore a data dimensionality reduction strategy is essential to comprehending [75]. Thus, while constructing a cybersecurity model with high dimensional data sets, an optimal number of security features selected based on their impact or importance [97] could reduce such issues. Similarly, principal component analysis (PCA) [98], Pearson correlation, regularization, etc. as discussed briefly in our earlier paper Sarker et al. [94] can handle such issues and could give better results for the resultant security model.

Reinforcement learning is another machine learning technique that typically enables an agent to learn in an interactive setting through trial and error while receiving feedback from its own actions and experiences. A Markov decision process is a common way to represent the environment. The most popular reinforcement learning algorithms in the field are Monte Carlo, Q-learning, Deep Q Networks, etc. [44]. For instance, authors in [61] provide CPSS LR-DDoS detection and defense in edge computing using DCNN Q-learning. For the purpose of anomaly detection in intelligent environments, a double deep Q-learning approach with prioritized experience replay has been proposed in [25].

Overall, we have detailed in Table 2 how various machine learning technologies, including deep learning, are utilized to address the main cybersecurity challenges. Accordingly, we can draw the conclusion that the aforementioned machine learning or deep learning techniques, as well as their variants or modified lightweight approaches or newly proposed algorithms, could play a significant role to achieve our goal in the context of security analytics.

4.6 Adversarial Machine Learning in Cybersecurity

In the domain of cybersecurity, ML approaches discussed above are typically employed to detect cyber security issues, where adversaries actively transform their objects to avoid detection. The study of adversarial machine learning focuses on how machine learning algorithms are attacked and how to defend against such attacks. Thus this is considered an emerging threat in learning systems that aims to deceive machine learning models by giving them false information. Machine learning systems can be attacked using a wide range of diverse adversarial strategies. Many of them employ classic machine learning models like linear regression and support vector machines (SVMs) [94], as well as deep learning [93] systems. In a white box attack, the attacker has total control over the target model, including its architecture and parameters. On the other hand, a black box attack is a situation in which the attacker is unable to access the model and is only able to observe the model's outputs. Adversarial attacks can be classified broadly into the following categories:

- *Poisoning Attacks:* This more sophisticated attack aims to affect the learning process by adding false or misleading data that discredits the algorithm's outputs. For instance, intrusion detection systems (IDSs) are often re-trained using collected data. This data may be contaminated by an attacker by inject-

ing malicious samples during operation, which then prevents retraining from taking place.

- *Evasion Attacks:* The most common and most investigated types of attacks are evasion attacks. During deployment, the attacker tampers with the data to mislead classifiers that have already been trained. They are the most common sorts of attacks employed in intrusion and malware scenarios since they are carried out during the deployment phase.
- *Model Extraction:* When a machine learning system is black-boxed, an attacker may analyze it to either reconstruct the model or retrieve the data it was trained on. This process is known as model hijacking or model extraction. This is especially crucial if the training data or the model itself contain sensitive or confidential information.

Defending robustly against adversarial attacks is still an open question. For each attack, a similar form of defense should be available. For instance, if malware targeting a machine learning model is similar to adversarial attacks, then security strategies might be thought of as anti-malware tools. Adversarial defense methods can be categorized as detection and robustness methods [89] defined as below:

- *Detection methods* - that are used to detect the adversarial examples.
- *Robustness methods* - that are used to enhance a classifier's rigidity to adversarial attacks without explicitly attempting to detect them.

Overall, in the field of cybersecurity, adversarial machine learning strives to confuse and trick models by producing special fraudulent inputs that mislead the model and cause it to malfunction. Organizations that implement machine learning technology need to be aware of the risks of adversarial samples, compromised models, and data manipulation. The majority of current adversarial machine learning research focuses on supervised learning [130]. On the other hand, labeling a huge number of data points or samples from the most recent attacks may demand expensive human expertise and turn into a significant bottleneck. Thus it is important to pay more attention to how to recognize adversarial samples in unsupervised and weakly supervised situations. Quantifying the robustness and accuracy trade-off for machine learning algorithms subject to adversarial attacks is crucial. Although certain robustness or uncertainty metrics have been proposed in the area, additional research on the trade-off is required to develop resilient learning algorithms.

## 5 Potential Use Cases of Machine Learning in Cybersecurity

Machine learning techniques have been effectively used to a variety of problems in a variety of application domains in the context of cybersecurity over the last several years. Intrusion detection, malware analysis, and detection, spam filtering, anomaly, and fraud detection, detecting zero-day attacks, cyberbullying detection, IoT attacks, and threat analysis, as well as a wide variety of other applications as shown in Figure 6, have all become commonplace. Defenders can identify and prioritize possible threats more precisely with the aid of machine learning, as discussed in the earlier Section 4. A wide range of specialized tasks, such as various types of vulnerability identification, deception, and attack disruption, could be

entirely or partially automated with the use of machine learning algorithms [94]. Several potential uses of ML in cybersecurity are discussed below.
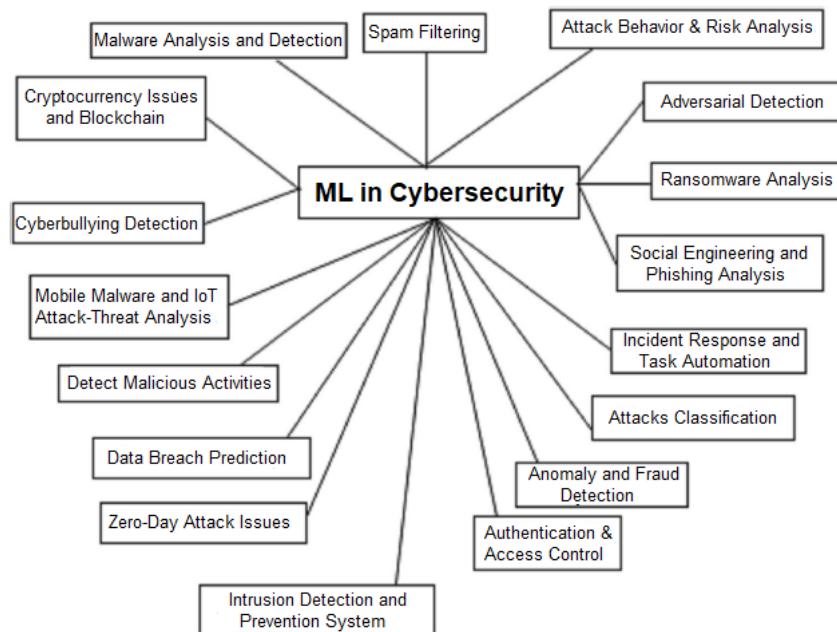


Fig. 6: Potential use cases of machine learning in cybersecurity.

— *Network risk scoring and prioritizing:* To determine which areas of the network have been targeted the most frequently, machine learning is being utilized to evaluate historical cyber threat data sets. Data from previous cyberattacks can be analyzed using machine learning algorithms [94], which can be used to identify the network segments that were most frequently targeted by a certain attack. Additionally, it is being used to identify the network components that, in the event of a breach, would cause the most significant harm to the business. With regard to a specific network area, this score can help estimate the likelihood and impact of an attack assisting organizations to lower their chance of becoming victims of such attacks. Cyber analysts are prioritizing their resources to concentrate on the biggest threats after giving each component of the company network a score.

— *Rapidly detect intrusions and response:* Machine learning is also being used by businesses to automatically and precisely identify malicious activities [91]. Organizations can respond to intrusions as soon as they happen because of the capability of machine learning models to detect, evaluate, and defend against diverse cyber threats in real time.

— *Identifying suspicious behaviors:* Machine learning techniques are also employed to identify suspicious user activity. Organizations use machine learning to distinguish between typical behavior and suspicious behavior that may be

signs of a cyber-attack in order to address vulnerabilities before a data breach occurs. This is done by monitoring users' suspicious activities, such as when they log in at odd hours of the day or download an unusually high volume of data or others.

— *Detecting fraud:* Many businesses use machine learning [94] and deep learning [93] algorithms to anticipate anomalous customer behavior in order to protect themselves against financial fraud. These technologies are assisting companies in identifying potential fraud threats before they materialize, hence minimizing their financial losses.

— *Discovering malware:* Businesses may now forecast future malware attacks with the use of the machine and deep learning as discussed in Section 4. Cyber analysts can predict malware attacks and reduce the risk at a speed that is not possible with manual operations by exploiting patterns observed in past attacks.

— *Detecting and classifying cyber-anomalies and multi-attacks:* Machine learning can quickly and easily analyze huge amounts of data, making it much faster than human threat detection. To find anomalies that might be signs of an attack, machine learning uses behavioral analysis and constantly changing parameters [91]. Intelligent security services can be created by building security models based on machine learning that assesses numerous cyberattacks or anomalies and finally detect or predict the threats.

— *Future predicting and responding to data breaches in real-time:* To predict cyber threats before they materialize, machine learning enables the processing of vast amounts of data from many sources. When a cyber threat is identified, machine learning has the ability to provide alerts and respond without human interference by rapidly building defensive patches in response to the attack.

— *Access control and advanced authentication:* Technology used in authentication validate that a user's credentials match those stored in the system of authorized users or in a data authentication server to enable access control for systems. In order to decide whether to ask users for multi-factor authentication, adaptive authentication leverages machine learning [94]. The procedure is more explicit because it makes use of a wide range of inputs to compute risk scores and choose the best security measure for a particular circumstance. Advanced authentication can be performed by applying machine learning to monitor in real time and find inconsistencies in the user's authentication behavior or even risks in the authentication process.

— *Cryptocurrencies and Blockchain intelligence:* Blockchain instantly produces enormous amounts of data. By analyzing blockchain data, we can discover potential issues, anticipate breakdowns, and pinpoint performance bottlenecks to optimize or improve the performance of blockchain systems. Massive blockchain data sets can be used for ML research to find abnormalities, assess market manipulations, and identify fraudulent users. Automatically identifying and locating exploitable flaws in smart contracts is possible with the support of machine learning techniques. For instance, ML regression models are employed to forecast the returns of the cryptocurrency-based dependent variable, and classification models are utilized to produce binary buy or sell trading recommendations in [107].

— *Automating tasks:* One of the main benefits of machine learning in cyber security is the automation of repetitive and time-consuming tasks including vul-

nerability assessments, malware analysis, network log analysis, and intelligence evaluation. By including machine learning in the security workflow, businesses may finish tasks more quickly and respond to threats and incidents at a rate that would be impossible with solely manual human expertise. By automating repetitive tasks, businesses may simply scale up or down without changing the amount of manpower required, hence reducing expenses.

Overall, we think that machine learning can be applied to improve security procedures and make it automate and intelligent for security analysts to recognize, prioritize, respond to, and address emerging attacks and threats in a variety of cyber security application areas. We have also listed various machine learning tasks and approaches in Table 2 that are used to solve various cybersecurity challenges. As illustrated in Table 2, machine learning modeling has a wide range of applications in real-world application domains, and there are various opportunities to work and conduct research in the context of cybersecurity. In the following section, we will look at the future aspect of machine learning, as well as research concerns in automation and intelligent decision-making in the cybersecurity area.

## 6 Future Aspects and Research Directions

In the cyber security world, machine learning has become a popular buzzword. As cyber-attacks become more widespread, sophisticated, and targeted, automation is becoming a crucial tool for overloaded security professionals. More automated methods for detecting risks and malicious user behavior are desperately needed by security teams, and machine learning provides a promising future.

Cybersecurity is considered a 'zero-tolerance field', meaning that one successful attack results in the security system failing. In their efforts to escape detection, cyber adversaries are growing more sophisticated, and many modern malware tools are already adding new ways to get around antivirus and other threat detection measures. Cybersecurity, on the other hand, is in a crisis, and future research efforts should be focused on cyber-threat intelligent systems that can predict crucial scenarios and consequences, rather than depending on defensive measures and mitigation. Systems that are based on a complete, predictive study of cyber risks are required all around the world. Machine learning [94] enables round-the-clock monitoring and can manage much more data than a human can. Thus the necessary functions may be beneficial to a successful cybersecurity system that achieves the desired results. These are:

- *Prediction:* To predict most likely attacks, targets, and methods.
- *Prevention:* To prevent or deter attacks so no loss is experienced.
- *Detection:* In order to respond quickly and thoroughly, it is necessary to identify attacks that could not be prevented.
- *Response:* To promptly address issues in order to reduce losses and get back to normal.

In summary, machine learning has the potential to improve cybersecurity by making it intelligent, more proactive, economical, and efficient. There are a number of machine learning methods that are frequently categorized as supervised or unsupervised learning. Since supervised learning requires annotated training datasets

[92], it is less suited for cyber security. Unsupervised learning, on the other hand, is more appropriate for discovering unusual activities, such as attacks that have never been seen before because it does not require labeled training data. So it can be difficult to choose a learning algorithm that is suitable for the intended application. This is because, depending on the quality of the data, different learning algorithms may produce different outcomes [91] [97]. The techniques presented in Section 4 can be utilized directly to tackle various real-world issues in the context of cybersecurity, as outlined in Section 5. However, a future study in the field could include a hybrid learning model, such as an ensemble of methods, updating with an improvement, or designing novel algorithms or models.

The nature and quality of the data, as well as the general success of the learning algorithms, have an impact on how effective and efficient a machine learning-based solution is. One of the most challenging concerns is gathering data from endpoints, networks, and clouds, standardizing it, and then using it effectively for machine learning [104]. Furthermore, historical data may include a sizable number of ambiguous values, missing values, outliers, and other data that is otherwise worthless [90] [92]. As a result, it can be challenging to clean and pre-process various data from various sources. Therefore, both data and learning algorithms are necessary for a machine learning-based solution to be effective over the long term and for its applications. If the data is unsuitable for learning, such as having non-representative, low-quality, irrelevant features, or not enough for training, machine learning models may become useless or deliver less accurate results.

Overall, machine learning has emerged as a crucial tool for cybersecurity. Nowadays, deploying good cybersecurity solutions without relying substantially on machine learning is nearly difficult. However, machine learning is challenging to deploy successfully without a thorough, in-depth, and complete approach to the underlying data. Cybersecurity systems built on machine learning can identify patterns and learn from them to help deter reoccurring cyberattacks and adapt to changeable behavior. It has the ability to make cybersecurity teams more proactive in terms of preventing threats and responding to active attacks in real time. Thus this can help organizations better allocate their resources by minimizing the amount of time they spend on routine tasks.

## 7 Conclusion

We have provided a comprehensive view of machine learning techniques for intelligent data analysis and automation in cybersecurity in this paper. For this, we have explored briefly the potentiality of various machine learning techniques to solve practical issues across a range of cyber application fields covered in the paper. The success of a machine learning model depends on how well the data and learning algorithms perform. Prior to the system being able to enable intelligent decision-making, the sophisticated learning algorithms should be trained utilizing real-world cyber data and information particular to the target application. Finally, we discussed the challenges as well as potential future research direction in the field. Overall, we believe that our study on machine learning-based modeling and security solutions is useful and points in the right direction for further research and application by academics and professionals in the domain of cybersecurity.

## Compliance with ethical standards

**Conflict of interest** The author declares no conflict of interest.

## References

1. Caida ddos attack 2007 dataset. http://www.caida.org/data/ passive/ddos-20070804-dataset.xml/ (accessed on 20 october 2019).
2. Canadian institute of cybersecurity, university of new brunswick, iscx dataset, url http://www.unb.ca/cic/datasets/index.html/ (accessed on 20 october 2019).
3. The ctu-13 dataset. available online: https://stratosphereips.org/category/datasets-ctu13 (accessed on 20 october 2019).
4. Kdd cup 99. available online: http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html (accessed on 20 october 2019).
5. Rakesh Agrawal, Tomasz Imieliński, and Arun Swami. Mining association rules between sets of items in large databases. In *ACM SIGMOD Record*, volume 22, pages 207–216. ACM, 1993.
6. Rakesh Agrawal, Ramakrishnan Srikant, et al. Fast algorithms for mining association rules. In *Proc. 20th int. conf. very large data bases, VLDB*, volume 1215, pages 487–499, 1994.
7. David W Aha, Dennis Kibler, and Marc K Albert. Instance-based learning algorithms. *Machine learning*, 6(1):37–66, 1991.
8. Mohammad Al-Omari, Majdi Rawashdeh, Fadi Qutaishat, Mohammad Alshira'H, and Nedal Ababneh. An intelligent tree-based intrusion detection model for cyber security. *Journal of Network and Systems Management*, 29(2):1–18, 2021.
9. Muder Almiani, Alia AbuGhazleh, Amer Al-Rahayfeh, Saleh Atiewi, and Abdul Razaque. Deep recurrent neural network for iot intrusion detection system. *Simulation Modelling Practice and Theory*, page 102031, 2019.
10. Ibrahim Alrashdi, Ali Alqazzaz, Esam Aloufi, Raed Alharthi, Mohamed Zohdy, and Hua Ming. Ad-iot: Anomaly detection of iot cyberattacks in smart city using machine learning. In *2019 IEEE 9th Annual Computing and Communication Workshop and Conference (CCWC)*, pages 0305–0310. IEEE, 2019.
11. Khaled Alrawashdeh and Carla Purdy. Toward an online anomaly intrusion detection system based on deep learning. In *2016 15th IEEE International Conference on Machine Learning and Applications (ICMLA)*, pages 195–200. IEEE, 2016.
12. Chandrashekhar Azad and Vijay Kumar Jha. Genetic algorithm to solve the problem of small disjunct in the decision tree based intrusion detection system. *International Journal of Computer Network and Information Security (IJCNIS)*, 7(8):56, 2015.
13. Abdullateef Oluwagbemiga Balogun and Rasheed Gbenga Jimoh. Anomaly intrusion detection using an hybrid of decision tree and k-nearest neighbor. 2015.
14. Rohan Bapat, Abhijith Mandya, Xinyang Liu, Brendan Abraham, Donald E Brown, Hyojung Kang, and Malathi Veeraraghavan. Identifying malicious botnet traffic using logistic regression. In *2018 Systems and Information Engineering Design Symposium (SIEDS)*, pages 266–271. IEEE, 2018.
15. Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
16. Leo Breiman, Jerome H Friedman, Richard A Olshen, and Charles J Stone. *Classification and regression trees*. Routledge, 2017.
17. AM Chandrasekhar and K Raghuveer. Confederation of fcm clustering, ann and svm techniques to implement hybrid nids using corrected kdd cup 99 dataset. In *2014 International Conference on Communication and Signal Processing*, pages 672–676. IEEE, 2014.
18. Yaping Chang, Wei Li, and Zhongming Yang. Network intrusion detection based on random forest and support vector machine. In *2017 IEEE international conference on computational science and engineering (CSE) and IEEE international conference on embedded and ubiquitous computing (EUC)*, volume 1, pages 635–638. IEEE, 2017.
19. EG Dada. A hybridized svm-knn-pdapso approach to intrusion detection system. In *Proc. Fac. Seminar Ser.*, pages 14–21, 2017.
20. Jiazhu Dai, Chuanshuai Chen, and Yufeng Li. A backdoor attack against lstm-based text classification systems. *IEEE Access*, 7:138872–138878, 2019.

21. Amitabha Das, Wee-Keong Ng, and Yew-Kwong Woon. Rapid association rule mining. In *Proceedings of the tenth international conference on Information and knowledge management*, pages 474–481. ACM, 2001.

22. Aminu Da'u and Naomie Salim. Recommendation system based on deep learning methods: a systematic review and new directions. *Artificial Intelligence Review*, 53(4):2709–2748, 2020.

23. Li Deng. A tutorial survey of architectures, algorithms, and applications for deep learning. *APSIPA Transactions on Signal and Information Processing*, 3, 2014.

24. Rohan Doshi, Noah Apthorpe, and Nick Feamster. Machine learning ddos detection for consumer internet of things devices. In *2018 IEEE Security and Privacy Workshops (SPW)*, pages 29–35. IEEE, 2018.

25. Daniel Fährmann, Nils Jorek, Naser Damer, Florian Kirchbuchner, and Arjan Kuijper. Double deep q-learning with prioritized experience replay for anomaly detection in smart environments. *IEEE Access*, 2022.

26. Mohamed Amine Ferrag, Leandros Maglaras, Sotiris Moschoyiannis, and Helge Janicke. Deep learning for cyber security intrusion detection: Approaches, datasets, and comparative study. *Journal of Information Security and Applications*, 50:102419, 2020.

27. Peter A Flach and Nicolas Lachiche. Confirmation-guided discovery of first-order rules with tertius. *Machine Learning*, 42(1-2):61–95, 2001.

28. Felipe De Almeida Florencio, Edward David Moreno Ordonez, Hendrik Teixeira Macedo, Ricardo José Paiva De Britto Salgueiro, Filipe Barreto Do Nascimento, and Flavio Arthur Oliveira Santos. Intrusion detection via mlp neural network using an arduino embedded system. In *2018 VIII Brazilian Symposium on Computing Systems Engineering (SBESC)*, pages 190–195. IEEE, 2018.

29. Yoav Freund, Robert E Schapire, et al. Experiments with a new boosting algorithm. In *Icml*, volume 96, pages 148–156. Citeseer, 1996.

30. Xianwei Gao, Changzhen Hu, Chun Shan, Baoxu Liu, Zequn Niu, and Hui Xie. Malware classification for the cloud via semi-supervised transfer learning. *Journal of Information Security and Applications*, 55:102661, 2020.

31. Margaret Gratian, Sruthi Bandi, Michel Cukier, Josiah Dykstra, and Amy Ginther. Correlating human traits and cyber security behavior intentions. *computers & security*, 73:345–358, 2018.

32. Desta Haileselassie Hagos, Anis Yazidi, Øivind Kure, and Paal E Engelstad. Enhancing security attacks analysis using regularized machine learning techniques. In *2017 IEEE 31st International Conference on Advanced Information Networking and Applications (AINA)*, pages 909–918. IEEE, 2017.

33. Hyo-Sik Ham, Hwan-Hee Kim, Myung-Sup Kim, and Mi-Jung Choi. Linear svm-based android malware detection for reliable iot services. *Journal of Applied Mathematics*, 2014, 2014.

34. Jiawei Han, Jian Pei, and Micheline Kamber. *Data mining: concepts and techniques*. Elsevier, 2011.

35. Jiawei Han, Jian Pei, and Yiwen Yin. Mining frequent patterns without candidate generation. In *ACM Sigmod Record*, volume 29, pages 1–12. ACM, 2000.

36. Kazuki Hara and Kohei Shiomoto. Intrusion detection system using semi-supervised learning with adversarial auto-encoder. In *NOMS 2020-2020 IEEE/IFIP Network Operations and Management Symposium*, pages 1–8. IEEE, 2020.

37. Dang Hai Hoang and Ha Duong Nguyen. A pca-based method for iot network traffic anomaly detection. In *2018 20th International Conference on Advanced Communication Technology (ICACT)*, pages 381–386. IEEE, 2018.

38. Maurice Houtsma and Arun Swami. Set-oriented mining for association rules in relational databases. In *Data Engineering, 1995. Proceedings of the Eleventh International Conference on*, pages 25–33. IEEE, 1995.

39. Venkatesh Jaganathan, Priyesh Cherurveettil, and Premapriya Muthu Sivashanmugam. Using a prediction model to manage cyber security threats. *The Scientific World Journal*, 2015, 2015.

40. Yousra Javed and Navid Rajabi. Multi-layer perceptron artificial neural network based iot botnet traffic classification. In *Proceedings of the Future Technologies Conference*, pages 973–984. Springer, 2019.

41. Xuyang Jing, Zheng Yan, Xueqin Jiang, and Witold Pedrycz. Network traffic fusion and analysis against ddos flooding attacks with a novel reversible sketch. *Information Fusion*, 51:100–113, 2019.

42. Seongrae Jo, Haengnam Sung, and Byunghyuk Ahn. A comparative study on the performance of intrusion detection using decision tree and artificial neural network models. *Journal of the Korea Society of Digital Industry and Information Management*, 11(4):33–45, 2015.

43. George H John and Pat Langley. Estimating continuous distributions in bayesian classifiers. In *Proceedings of the Eleventh conference on Uncertainty in artificial intelligence*, pages 338–345. Morgan Kaufmann Publishers Inc., 1995.

44. Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. *Journal of artificial intelligence research*, 4:237–285, 1996.

45. ElMouatez Billah Karbab, Mourad Debbabi, Abdelouahid Derhab, and Djedjiga Mouheb. Maldozer: Automatic framework for android malware detection using deep learning. *Digital Investigation*, 24:S48–S59, 2018.

46. S. Sathiya Keerthi, Shirish Krishnaj Shevade, Chiranjib Bhattacharyya, and Karuturi Radha Krishna Murthy. Improvements to platt's smo algorithm for svm classifier design. *Neural computation*, 13(3):637–649, 2001.

47. Jihyun Kim, Jaehyun Kim, Huong Le Thi Thu, and Howon Kim. Long short term memory recurrent neural network classifier for intrusion detection. In *2016 International Conference on Platform Technology and Service (PlatCon)*, pages 1–5. IEEE, 2016.

48. Bojan Kolosnjaji, Apostolis Zarras, George Webster, and Claudia Eckert. Deep learning for classification of malware system call sequences. In *Australasian Joint Conference on Artificial Intelligence*, pages 137–149. Springer, 2016.

49. Nickolaos Koroniotis, Nour Moustafa, Elena Sitnikova, and Benjamin Turnbull. Towards the development of realistic botnet dataset in the internet of things for network forensic analytics: Bot-iot dataset. *Future Generation Computer Systems*, 100:779–796, 2019.

50. Manjiri V Kotpalliwar and Rakhi Wajgi. Classification of attacks using support vector machine (svm) on kddcup'99 ids database. In *2015 Fifth International Conference on Communication Systems and Network Technologies*, pages 987–990. IEEE, 2015.

51. Mohammed Lalou, Hamamache Kheddouci, and Salim Hariri. Identifying the cyber attack origin with partial observation: a linear regression based approach. In *2017 IEEE 2nd International Workshops on Foundations and Applications of Self* Systems (FAS* W)*, pages 329–333. IEEE, 2017.

52. Max Landauer, Florian Skopik, Markus Wurzenberger, and Andreas Rauber. System log clustering approaches for cyber security applications: A survey. *Computers & Security*, 92:101739, 2020.

53. Saskia Le Cessie and Johannes C Van Houwelingen. Ridge estimators in logistic regression. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 41(1):191–201, 1992.

54. Soo-Yeon Lee, Sa-rang Wi, Eunil Seo, Jun-Kwon Jung, and Tai-Myoung Chung. Profiot: Abnormal behavior profiling (abp) of iot devices based on a machine learning approach. In *2017 27th International Telecommunication Networks and Applications Conference (ITNAC)*, pages 1–6. IEEE, 2017.

55. Heng Li, ShiYao Zhou, Wei Yuan, Jiahuan Li, and Henry Leung. Adversarial-example attacks toward android malware detection system. *IEEE Systems Journal*, 14(1):653–656, 2019.

56. Yanmiao Li, Yingying Xu, Zhi Liu, Haixia Hou, Yushuo Zheng, Yang Xin, Yuefeng Zhao, and Lizhen Cui. Robust detection for network intrusion of industrial iot based on multi-cnn fusion. *Measurement*, 154:107450, 2020.

57. Yinhui Li, Jingbo Xia, Silan Zhang, Jiakai Yan, Xiaochuan Ai, and Kuobin Dai. An efficient intrusion detection system based on support vector machines and gradually feature removal method. *Expert Systems with Applications*, 39(1):424–430, 2012.

58. Wei-Chao Lin, Shih-Wen Ke, and Chih-Fong Tsai. Cann: An intrusion detection system based on combining cluster centers and nearest neighbors. *Knowledge-based systems*, 78:13–21, 2015.

59. Richard P Lippmann, David J Fried, Isaac Graf, Joshua W Haines, Kristopher R Kendall, David McClung, Dan Weber, Seth E Webster, Dan Wyschogrod, Robert K Cunningham, et al. Evaluating intrusion detection systems: The 1998 darpa off-line intrusion detection evaluation. In *Proceedings DARPA Information Survivability Conference and Exposition. DISCEX'00*, volume 2, pages 12–26. IEEE, 2000.

60. Liqun Liu, Bing Xu, Xiaoping Zhang, and Xianjun Wu. An intrusion detection method for internet of things based on suppressed fuzzy clustering. *EURASIP Journal on Wireless Communications and Networking*, 2018(1):113, 2018.

61. Zengguang Liu, Xiaochun Yin, and Yuemei Hu. Cpss lr-ddos detection and defense in edge computing utilizing dcnn q-learning. *IEEE Access*, 8:42120–42130, 2020.

62. Bing Liu Wynne Hsu Yiming Ma. Integrating classification and association rule mining. In *Proceedings of the fourth international conference on knowledge discovery and data mining*, 1998.

63. James MacQueen. Some methods for classification and analysis of multivariate observations. In *Fifth Berkeley symposium on mathematical statistics and probability*, volume 1, 1967.

64. Arif Jamal Malik and Farrukh Aslam Khan. A hybrid technique using binary particle swarm optimization and decision tree pruning for network intrusion detection. *Cluster Computing*, 21(1):667–680, 2018.

65. Niall McLaughlin, Jesus Martinez del Rincon, BooJoong Kang, Suleiman Yerima, Paul Miller, Sakir Sezer, Yeganeh Safaei, Erik Trickel, Ziming Zhao, Adam Doupé, et al. Deep android malware detection. In *Proceedings of the Seventh ACM on Conference on Data and Application Security and Privacy*, pages 301–308, 2017.

66. Weizhi Meng, Wenjuan Li, and Lam-For Kwok. Design of intelligent knn-based alarm filter using knowledge-based alert verification in intrusion detection. *Security and Communication Networks*, 8(18):3883–3895, 2015.

67. Robert Mitchell and Ray Chen. Behavior rule specification-based intrusion detection for safety critical medical cyber physical systems. *IEEE Transactions on Dependable and Secure Computing*, 12(1):16–30, 2014.

68. TagyAldeen Mohamed, Takanobu Otsuka, and Takayuki Ito. Towards machine learning based iot intrusion detection service. In *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems*, pages 580–585. Springer, 2018.

69. Rajalakshmi Shenbaga Moorthy and P Pabitha. Optimal detection of phising attack using sca based k-nn. *Procedia Computer Science*, 171:1716–1725, 2020.

70. Nour Moustafa and Jill Slay. Unsw-nb15: a comprehensive data set for network intrusion detection systems (unsw-nb15 network data set). In *2015 military communications and information systems conference (MilCIS)*, pages 1–6. IEEE, 2015.

71. Seiichi Ozawa, Tao Ban, Naoki Hashimoto, Junji Nakazato, and Jumpei Shimamura. A study of iot malware activities using association rule learning for darknet sensor data. *International Journal of Information Security*, 19(1):83–92, 2020.

72. Gonzalo De La Torre Parra, Paul Rad, Kim-Kwang Raymond Choo, and Nicole Beebe. Detecting internet of things attacks using distributed deep learning. *Journal of Network and Computer Applications*, page 102662, 2020.

73. Subbaraju Pericherla and E Ilavarasan. Transformer network-based word embeddings approach for autonomous cyberbullying detection. *International Journal of Intelligent Unmanned Systems*, 2021.

74. Muhammad Shakil Pervez and Dewan Md Farid. Feature selection and intrusion classification in nsl-kdd cup 99 dataset employing svms. In *The 8th International Conference on Software, Knowledge, Information Management and Applications (SKIMA 2014)*, pages 1–6. IEEE, 2014.

75. Morteza Safaei Pour, Elias Bou-Harb, Kavita Varma, Nataliia Neshenko, Dimitris A Pados, and Kim-Kwang Raymond Choo. Comprehending the iot cyber threat landscape: A data dimensionality reduction technique to infer and characterize internet-scale iot probing campaigns. *Digital Investigation*, 28:S40–S49, 2019.

76. Rifkie Primartha and Bayu Adhi Tama. Anomaly detection using random forest: A performance revisited. In *2017 International conference on data and software engineering (ICoDSE)*, pages 1–6. IEEE, 2017.

77. Anton O Prokofiev, Yulia S Smirnova, and Vasiliy A Surov. A method to detect internet of things botnets. In *2018 IEEE Conference of Russian Young Researchers in Electrical and Electronic Engineering (EIConRus)*, pages 105–108. IEEE, 2018.

78. Shubha Puthran and Ketan Shah. Intrusion detection using improved decision tree algorithm with binary and quad split. In *International Symposium on Security in Computing and Communication*, pages 427–438. Springer, 2016.

79. J. Ross Quinlan. Induction of decision trees. *Machine learning*, 1(1):81–106, 1986.

80. J. Ross Quinlan. C4.5: Programs for machine learning. *Machine Learning*, 1993.

81. Panagiotis I Radoglou-Grammatikis and Panagiotis G Sarigiannidis. An anomaly-based intrusion detection system for the smart grid based on cart decision tree. In *2018 Global Information Infrastructure and Networking Symposium (GIIS)*, pages 1–5. IEEE, 2018.

82. Kajal Rai, M Syamala Devi, and Ajay Guleria. Decision tree based algorithm for intrusion detection. *International Journal of Advanced Networking and Applications*, 7(4):2828, 2016.

83. Hariharan Rajadurai and Usha Devi Gandhi. A stacked ensemble learning model for intrusion detection in wireless network. *Neural computing and applications*, pages 1–9, 2020.

84. MR Gauthama Raman, Nivethitha Somu, Sahruday Jagarapu, Tina Manghnani, Thirumaran Selvam, Kannan Krithivasan, and VS Shankar Sriram. An efficient intrusion detection technique based on support vector machine and improved binary gravitational search algorithm. *Artificial Intelligence Review*, pages 1–32, 2019.

85. Shailendra Rathore and Jong Hyuk Park. Semi-supervised learning based distributed attack detection framework for iot. *Applied Soft Computing*, 72:79–89, 2018.

86. Neha G Relan and Dharmaraj R Patil. Implementation of network intrusion detection system using variant of decision tree algorithm. In *2015 International Conference on Nascent Technologies in the Engineering Field (ICNTE)*, pages 1–5. IEEE, 2015.

87. Paulo Angelo Alves Resende and André Costa Drummond. A survey of random forest based methods for intrusion detection systems. *ACM Computing Surveys (CSUR)*, 51(3):1–36, 2018.

88. Lior Rokach. A survey of clustering algorithms. In *Data Mining and Knowledge Discovery Handbook*, pages 269–298. Springer, 2010.

89. Ishai Rosenberg, Asaf Shabtai, Yuval Elovici, and Lior Rokach. Adversarial machine learning attacks and defense methods in the cyber security domain. *ACM Computing Surveys (CSUR)*, 54(5):1–36, 2021.

90. Iqbal H Sarker. A machine learning based robust prediction model for real-life mobile phone data. *Internet of Things*, 5:180–193, 2019.

91. Iqbal H Sarker. Cyberlearning: effectiveness analysis of machine learning security modeling to detect cyber-anomalies and multi-attacks. *Internet of Things*, 14:100393, 2021.

92. Iqbal H Sarker. Data science and analytics: An overview from data-driven smart computing, decision-making and applications perspective. *SN Computer Science*, 2021.

93. Iqbal H Sarker. Deep learning: A comprehensive overview on techniques, taxonomy, applications and research directions. *SN Computer Science*, 2021.

94. Iqbal H Sarker. Machine learning: Algorithms, real-world applications and research directions. *SN Computer Science*, 2(3):1–21, 2021.

95. Iqbal H Sarker. Ai-based modeling: Techniques, applications and research issues towards automation, intelligent and smart systems. *SN Computer Science*, 3(2):1–20, 2022.

96. Iqbal H Sarker. Smart city data science: Towards data-driven smart cities with open research issues. *Internet of Things*, 19:100528, 2022.

97. Iqbal H Sarker, Yoosef B Abushark, Fawaz Alsolami, and Asif Irshad Khan. Intrudtree: A machine learning based cyber security intrusion detection model. *Symmetry*, 12(5):754, 2020.

98. Iqbal H Sarker, Yoosef B Abushark, and Asif Irshad Khan. Contextpca: Predicting context-aware smartphone apps usage based on machine learning techniques. *Symmetry*, 12(4):499, 2020.

99. Iqbal H Sarker, Alan Colman, Jun Han, Asif Irshad Khan, Yoosef B Abushark, and Khaled Salah. Behavdt: A behavioral decision tree learning to build user-centric context-aware predictive model. *Mobile Networks and Applications*, pages 1–11, 2019.

100. Iqbal H Sarker, Alan Colman, Jun Han, and Paul A Watters. *Context-aware machine learning and mobile data analytics: automated rule-based services with intelligent decision-making*. Springer Nature, 2021.

101. Iqbal H Sarker, Alan Colman, Muhammad Ashad Kabir, and Jun Han. Individualized time-series segmentation for mining mobile phone user behavior. *The Computer Journal*, 61(3):349–368, 2018.

102. Iqbal H Sarker, Md Hasan Furhad, and Raza Nowrozy. Ai-driven cybersecurity: an overview, security intelligence modeling and research directions. *SN Computer Science*, 2(3):1–18, 2021.

103. Iqbal H Sarker and ASM Kayes. Abc-ruleminer: User behavioral rule-based machine learning method for context-aware intelligent services. *Journal of Network and Computer Applications*, 168:102762, 2020.

104. Iqbal H Sarker, ASM Kayes, Shahriar Badsha, Hamed Alqahtani, Paul Watters, and Alex Ng. Cybersecurity data science: an overview from machine learning perspective. *Journal of Big Data*, 7(1):1–29, 2020.

105. Iqbal H Sarker, Asif Irshad Khan, Yoosef B Abushark, and Fawaz Alsolami. Internet of things (iot) security intelligence: a comprehensive overview, machine learning solutions and research directions. *Mobile Networks and Applications*, pages 1–17, 2022.

106. Harshit Saxena and Vineet Richariya. Intrusion detection in kdd99 dataset using svm-pso and feature reduction with information gain. *International Journal of Computer Applications*, 98(6), 2014.

107. Helder Sebastião and Pedro Godinho. Forecasting and trading cryptocurrencies with machine learning under changing market conditions. *Financial Innovation*, 7(1):1–30, 2021.

108. Devaraju Sellappan and Ramakrishnan Srinivasan. Association rule-mining-based intrusion detection system with entropy-based feature selection: Intrusion detection system. In *Handbook of Research on Intelligent Data Processing and Information Security Systems*, pages 1–24. IGI Global, 2020.

109. Hossein Shapoorifard and Pirooz Shamsinejad. Intrusion detection using a novel hybrid method incorporating an improved knn. *Int. J. Comput. Appl.*, 173(1):5–9, 2017.

110. Aboosaleh M Sharifi, Saeed K Amirgholipour, and Alireza Pourebrahimi. Intrusion detection based on joint of k-means and knn. *Journal of Convergence Information Technology*, 10(5):42, 2015.

111. Yong Shi. *Advances in big data analytics: theory, algorithms and practices*. Springer Nature, 2022.

112. Beata Ślusarczyk. Industry 4.0: Are we ready? *Polish Journal of Management Studies*, 17, 2018.

113. Peter HA Sneath. The application of computers to taxonomy. *Journal of General Microbiology*, 17(1), 1957.

114. Thorvald Sorensen. method of establishing groups of equal amplitude in plant sociology based on similarity of species. *Biol. Skr.*, 5, 1948.

115. Bambang Susilo and Riri Fitri Sari. Intrusion detection in iot networks using deep learning algorithm. *Information*, 11(5):279, 2020.

116. Mayank Swarnkar and Neminath Hubballi. Ocpad: One class naive bayes classifier for payload based anomaly detection. *Expert Systems with Applications*, 64:330–339, 2016.

117. Syeda Manjia Tahsien, Hadis Karimipour, and Petros Spachos. Machine learning based solutions for security of internet of things (iot): A survey. *Journal of Network and Computer Applications*, 161:102630, 2020.

118. Arman Tajbakhsh, Mohammad Rahmati, and Abdolreza Mirzaei. Intrusion detection using fuzzy association rules. *Applied Soft Computing*, 9(2):462–469, 2009.

119. Mahbod Tavallaee, Ebrahim Bagheri, Wei Lu, and Ali A Ghorbani. A detailed analysis of the kdd cup 99 data set. In *2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications*, pages 1–6. IEEE, 2009.

120. James M Tien. Internet of things, real-time decision making, and artificial intelligence. *Annals of Data Science*, 4(2):149–178, 2017.

121. Raihan Ul Islam, Mohammad Shahadat Hossain, and Karl Andersson. A novel anomaly detection algorithm for sensor data under uncertainty. *Soft Computing*, 22(5):1623–1639, 2018.

122. R Vinayakumar, KP Soman, and Prabaharan Poornachandran. Deep android malware detection and classification. In *2017 International conference on advances in computing, communications and informatics (ICACCI)*, pages 1677–1683. IEEE, 2017.

123. Satyendra Vishwakarma, Vivek Sharma, and Ankita Tiwari. An intrusion detection system using knn-aco algorithm. *Int. J. Comput. Appl.*, 171(10):18–23, 2017.

124. Quang Hieu Vu, Dymitr Ruta, and Ling Cen. Gradient boosting decision trees for cyber security threats detection based on network events logs. In *2019 IEEE International Conference on Big Data (Big Data)*, pages 5921–5928. IEEE, 2019.

125. Qi Wang, Yue Ma, Kun Zhao, and Yingjie Tian. A comprehensive survey of loss functions in machine learning. *Annals of Data Science*, 9(2):187–212, 2022.

126. Wei Wang, Mengxue Zhao, and Jigang Wang. Effective android malware detection with a hybrid model based on deep autoencoder and convolutional neural network. *Journal of Ambient Intelligence and Humanized Computing*, 10(8):3035–3043, 2019.

127. Wei Wang, Ming Zhu, Xuewen Zeng, Xiaozhou Ye, and Yiqiang Sheng. Malware traffic classification using convolutional neural network for representation learning. In *2017 International Conference on Information Networking (ICOIN)*, pages 712–717. IEEE, 2017.

128. Peng Wei, Yufeng Li, Zhen Zhang, Tao Hu, Ziyong Li, and Diyang Liu. An optimization method for intrusion detection classification model based on deep belief network. *IEEE Access*, 7:87593–87605, 2019.

129. Zihan Wu, Hong Zhang, Penghai Wang, and Zhibo Sun. Rtids: a robust transformer-based approach for intrusion detection system. *IEEE Access*, 2022.

130. Bowei Xi. Adversarial machine learning for cybersecurity and computer vision: Current developments and challenges. *Wiley Interdisciplinary Reviews: Computational Statistics*, 12(5):e1511, 2020.

131. Dongkuan Xu and Yingjie Tian. A comprehensive survey of clustering algorithms. *Annals of Data Science*, 2(2):165–193, 2015.

132. Binghao Yan and Guodong Han. Effective feature extraction via stacked sparse autoencoder to improve intrusion detection system. *IEEE Access*, 6:41238–41248, 2018.

133. Jinpei Yan, Yong Qi, and Qifan Rao. Detecting malware with an ensemble method based on deep neural network. *Security and Communication Networks*, 2018, 2018.

134. Manfu Yan and Zhifang Liu. A new method of transductive svm-based network intrusion detection. In *International Conference on Computer and Computing Technologies in Agriculture*, pages 87–95. Springer, 2010.

135. Chuanlong Yin, Yuefei Zhu, Jinlong Fei, and Xinzheng He. A deep learning approach for intrusion detection using recurrent neural networks. *Ieee Access*, 5:21954–21961, 2017.

136. Chuanlong Yin, Yuefei Zhu, Shengli Liu, Jinlong Fei, and Hetong Zhang. An enhancing framework for botnet detection using generative adversarial networks. In *2018 International Conference on Artificial Intelligence and Big Data (ICAIBD)*, pages 228–234. IEEE, 2018.

137. Mohammed Javeed Zaki. Scalable algorithms for association mining. *IEEE transactions on knowledge and data engineering*, 12(3):372–390, 2000.

138. Zhenfei Zhao, Weina Niu, Xiaosong Zhang, Runzi Zhang, Zhenqi Yu, and Cheng Huang. Trine: Syslog anomaly detection with three transformer encoders in one generative adversarial network. *Applied Intelligence*, 52(8):8810–8819, 2022.

139. Yun Zhou and Peichao Wang. An ensemble learning approach for xss attack detection with domain knowledge and threat intelligence. *Computers & Security*, 82:261–269, 2019.

140. Zhi-Jie Zhou, Guan-Yu Hu, Chang-Hua Hu, Cheng-Lin Wen, and Lei-Lei Chang. A survey of belief rule-base expert system. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 51(8):4944–4958, 2019.