

## *Retraction*

# **Retracted: Research on Boruta-ET-Based Anomalous Traffic Detection Model**

### **Security and Communication Networks**

Received 5 December 2023; Accepted 5 December 2023; Published 6 December 2023

Copyright © 2023 Security and Communication Networks. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

This article has been retracted by Hindawi, as publisher, following an investigation undertaken by the publisher [1]. This investigation has uncovered evidence of systematic manipulation of the publication and peer-review process. We cannot, therefore, vouch for the reliability or integrity of this article.

Please note that this notice is intended solely to alert readers that the peer-review process of this article has been compromised.

Wiley and Hindawi regret that the usual quality checks did not identify these issues before publication and have since put additional measures in place to safeguard research integrity.

We wish to credit our Research Integrity and Research Publishing teams and anonymous and named external researchers and research integrity experts for contributing to this investigation.

The corresponding author, as the representative of all authors, has been given the opportunity to register their agreement or disagreement to this retraction. We have kept a record of any response received.

### **References**

- [1] H. He, G. Huang, B. Zhang, and L. Qin, "Research on Boruta-ET-Based Anomalous Traffic Detection Model," *Security and Communication Networks*, vol. 2022, Article ID 9169266, 8 pages, 2022.

## Research Article

# Research on Boruta-ET-Based Anomalous Traffic Detection Model

Hongyan He <sup>1,2</sup>, Guoyan Huang <sup>1,2</sup>, Bing Zhang <sup>1,2</sup> and Lewei Qin <sup>1,2</sup>

<sup>1</sup>School of Information Science and Engineering, Yanshan University, Qinhuangdao, China

<sup>2</sup>The Key Laboratory for Software Engineering of Hebei Province, Qinhuangdao, Hebei, China

Correspondence should be addressed to Guoyan Huang; hgy@ysu.edu.cn

Received 13 April 2022; Revised 19 May 2022; Accepted 27 May 2022; Published 22 November 2022

Academic Editor: Muhammad Arif

Copyright © 2022 Hongyan He et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

With the rapid development of networks, intrusion detection has received increasing attention. In order to solve the problems of large dimensionality of intrusion detection data, unbalanced data samples, and large dispersion of datasets, which seriously affect the classification performance, this study proposes an anomaly detection based on Boruta and extreme tree (Boruta-ET) model. First, the network traffic data are preprocessed, which includes data cleaning, numerical and normalization processes, as well as equalization of the attack categories for a small number of samples by random oversampling at the data level; second, the traffic features are dimensionality reduced using the Boruta-based algorithm. The goal of Boruta dimensionality reduction is to extract all the features related to the dependent variable with a global dimension and find the optimal subset of features containing the most information; finally, the optimal feature subset is used as the input parameters of the extreme tree (ET) algorithm model for training and testing. Experiments were conducted on the real network traffic dataset CICIDS2017, and by evaluating the classification performance of several different machine learning algorithms, the experimental results show that the Boruta-ET model has the best performance with an accuracy rate of 99.80%, which can effectively improve the detection rate and achieve an effective recall rate for attack types with a small number of samples.

## 1. Introduction

In recent years, as the Internet has continued to grow, it has been integrated into all areas of people's daily lives, such as electronic communication, teaching, business, and entertainment. However, the massive expansion of the network has obviously led to an increase in network traffic data. As a result, this expansion has led to a number of security issues, such as a variety of known and unknown Internet attacks on network security. The need to develop network security has attracted a great deal of attention from industry and academia worldwide [1], and for this reason, the use of intrusion detection systems has become a necessary option for ensuring network security. Intrusion detection is an indispensable and very important line of defense in terms of security systems, which collects information from a number of critical nodes in a computer network security system, looks at the network for signs of violations of security policies and attacks, identifies threats in

the network and generates alerts, thus providing protection against internal attacks, external attacks, and misuse of implementation. Network intrusion detection systems (IDSs) are tools commonly used to detect network intrusions by collecting data on the current operational state of the network and analyzing network traffic using system preprogrammed algorithms and historical experience [2].

The study of intrusion detection has been the focus of national and international research scholars. Network traffic anomaly detection refers to the application of various anomaly detection techniques to analyze network traffic and detect network attacks in a timely manner. In order to achieve network anomaly detection and improve the accuracy of detection, various traditional and emerging techniques have been applied to network anomaly detection. Harish and Kumar [3] designed a fuzzy clustering-based network anomaly detection method. The method first eliminates duplicate samples from the sample set, based on which

principal component analysis is applied to select the most discriminative features, and finally, a fuzzy C-means algorithm is used to cluster the network samples. Mazini et al. [4] designed a network anomaly detection system combining reliable artificial bee colony and AdaBoost algorithms, with the artificial ant colony algorithm for feature selection and the AdaBoost algorithm for feature evaluation and classification, and validated it on the NSL-KDD and ISCXID2012 datasets. The accuracy and detection rate of the method were improved compared to traditional algorithms. Basati and Faghieh [5] proposed a novel lightweight architecture-parallel deep autoencoder (PDAE) that aims to construct nearest neighbor values and nearest neighbor information for each feature vector. The effectiveness of the proposed architecture was evaluated using the KDDCup99, UNSW-NB15, and CICIDS2017 datasets, and the evaluation results showed that the proposed model was effective in improving accuracy and performance. Zavrak and Iskefiyeli [6] proposed an anomaly detection model based on a variational autoencoder. The reconstruction error of the autoencoder is used as the anomaly score criterion to detect anomalies in network traffic. This model can only distinguish whether data traffic is intrusive or not and cannot detect specific types of intrusion attacks. Alkadi et al. [7] proposed a collaborative intrusion detection system based on a deep blockchain network, which is practical for identifying network traffic attacks on IoT networks. The study also focuses on privacy-preserving aspects by combining a trusted execution environment with blockchain technology for the purpose of providing confidentiality to smart contracts. The model was evaluated on the UNSW-NB15 dataset and the results showed that the system has high accuracy and detection rates when performing classification, especially for attacks that exploit cloud networks. Popoola et al. [8] proposed to reduce feature dimension through the encoding stage of long short-term memory autoencoder (LAE). By analyzing the association changes of the low-dimensional feature sets generated by LAE, in order to confirm the effectiveness of the method, a deep bidirectional long and short-term memory method (BLSTM) was used to achieve an improved classification accuracy of network traffic samples.

From the above think-aloud work, we found that the combination of feature selection and intrusion detection is a successful approach, as feature selection can assist in selecting the optimal subset of features with the most information and the least number of features from the entire feature set. When the distribution of class samples is unbalanced, it can affect the performance of the classification algorithm and thus reduce the detection rate, especially for a small number of classes. In network traffic, intrusions are much less common than normal behavior. Aiming at the problem of class imbalance in network intrusion traffic data, this study uses random oversampling to balance the data. Inspired by existing research, the use of feature selection and integrated classifiers has been highly successful in network traffic analysis and intrusion attack detection. We have designed the Boruta-ET model to address the problem of low accuracy and high false alarm rates, thus improving the efficiency of anomalous traffic detection.

The rest of the study is organized as follows: the second section describes the overall framework of the study and the sources of the experimental data. The third section specifies the key techniques studied in this study. The fourth section conducts various experimental validation studies and evaluates the model approach. The fifth section concludes the whole study as well as future perspectives.

## 2. Overall Architecture and Data Sources

*2.1. Overall Architecture.* In this section, the model proposed in this study, Boruta-ET, will be described in detail. The flowchart of this model is shown in Figure 1. First, the raw network traffic data are preprocessed, which includes data cleaning, character numerical normalization of the network traffic, and slicing of the network traffic dataset. Second, the Boruta [9] feature selection is performed on the training set of the network traffic data, and then the selected feature subsets are counted and the training set is randomly oversampled to expand the attack types of a small number of samples for the purpose of balancing the dataset. Finally, the optimal feature subset is used as the input data for the ET algorithm model for training, and the performance of the model is evaluated using the testing dataset data to obtain the final classification results of the model.

*2.2. Data Sources.* The CICIDS2017 [10] dataset used in this study was published by the Canadian Cyber Security Institute, which spans eight different files, and a short description of all of them is listed in Table 1. The CICIDS2017 dataset is the largest intrusion detection dataset currently available on the Internet, and the dataset contains 11 of the most important features, namely, attack diversity, available protocols, complete captures, metadata, complete interactions, heterogeneity, complete network configurations, feature sets, complete traffic, anonymity, and tagging [11]. In addition, it contains necessary and newer examples of attacks such as botnets, distributed DoS (DDoS), port scanning, and SQL injection [12]. In the previous publicly available dataset, there were fewer types of traffic, less capacity, various anonymous traffic packets, and payloads of information, and also there were many limitations on the various types of traffic attacks. However, the CICIDS2017 dataset has overcome the problems mentioned above, and the dataset contains various protocols such as FTP, HTTP, SSH, HTTPS, and e-mail that are not available in the previous dataset. The dataset has a total of 2830743 tagged network flows, each with 79 characteristics, which are distributed in 8 files, including SYN flag count, stream duration, destination port, etc.

## 3. Methodology

*3.1. Boruta Feature Selection.* Boruta aims to select the set of all features that are relevant to the dependent variable and is a wrapper algorithm that uses a random forest as a classifier to filter out the features that are relevant to the dependent variable across all features to construct a new subset of features, primarily by reducing the average precision value.

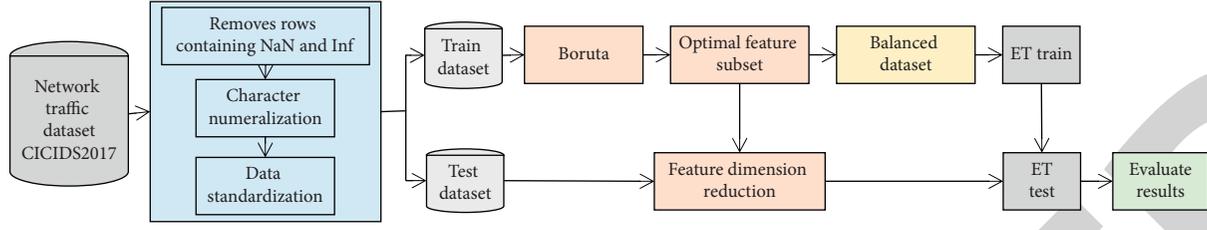


FIGURE 1: Flowchart of the Boruta-ET model framework.

TABLE 1: CICIDS2017 dataset file description.

Name of files	Attacks	Day activity
Monday-WorkingHours.pcap_ISCX.csv	Benign (normal human activities)	Monday
Tuesday-WorkingHours.pcap_ISCX.csv	Benign, FTP-patator, SSH-patator	Tuesday
Wednesday-WorkingHours.pcap_ISCX.csv	Benign, DoS GoldenEye, DoS hulk, DoS slowhttptest, DoS slowloris, Heartbleed	Wednesday
Thursday-WorkingHours-Morning-WebAttacks.pcap_ISCX.csv	Benign, web attack—brute force, web attack—Sql injection, web attack—XSS	Thursday
Thursday-WorkingHours-afternoon-infiltration.pcap_ISCX.csv	Benign, infiltration	Thursday
Friday-WorkingHours-morning.pcap_ISCX.csv	Benign, bot	Friday
Friday-WorkingHours-Afternoon-PortScan.pcap_ISCX.csv	Benign, port scan	Friday
Friday-WorkingHours-afternoon-DDoS.pcap_ISCX.csv	Benign, DDoS	Friday

The Boruta algorithm obtains the importance of all features in the dataset with respect to the target variable, selects the important features, removes the redundant ones, and features a black box predictive model with good predictive accuracy to obtain the importance indicators associated with the target variable. The flowchart of the Boruta algorithm is shown in Figure 2.

Boruta's algorithm consists of the following steps:

- (1) The individual features of the feature matrix  $X$  are shuffled, and the original features are spliced with the shuffled features to construct a new feature matrix, that is, a matrix with two times the number of features.
- (2) Randomly disrupt the added attributes to remove their correlation with the response.
- (3) Run a random forest classifier on the expanded feature matrix, using the newly constructed feature matrix as the input of the classifier, and the feature\_importance of each feature can be output through the training of the model.
- (4) Calculate the  $Z\_Score$  for original features and shadow features.

The importance score in Boruta's algorithm is defined based on the out-of-bag error of the RF model and is given by the following equation:

$$MSE_{OOB} = \frac{(y_i - \hat{y}_{i_{OOB}})^2}{N}. \quad (1)$$

Here,  $MSE_{OOB}$  is the out-of-bag error of the random forest,  $y_i$  is the sample value, and  $\hat{y}_{i_{OOB}}$  is the predicted value of the out-of-bag sample of the sample  $y_i$ .

$$Z\_Score = \frac{\overline{MSE_{OOB}}}{SDMSE_{OOB}}. \quad (2)$$

Here,  $Z\_Score$  is the z-score,  $\overline{MSE_{OOB}}$  is the mean of the out-of-bag error, and  $SDMSE_{OOB}$  is the standard deviation of the out-of-bag error.

- (5) Find the maximum  $Z\_Score$  in the shadow features matrix, which is  $S\_max$ , and use  $S\_max$  as the screening index.
- (6) Original features with  $Z\_Score$  higher than  $S\_max$  are regarded as "important" and reserved. Original features with  $Z\_score$  lower than  $S\_max$  are considered "unimportant" and permanently removed from the feature set.
- (7) Repeat this process until all features are assigned importance.

**3.2. Extreme Trees.** Extreme trees are an integrated learning prediction method based on decision trees. The extreme tree algorithm is based on the traditional top-down approach of building a series of unpruned decision trees. It has two main features: first, each decision tree is built using the full training sample; second, each decision tree completes the node splitting by choosing the splitting threshold completely randomly. Algorithm 1 is the limit random tree algorithm pseudocode.

**3.3. Evaluation Metrics.** In order to verify the performance of each algorithm, the experiments in this study mainly use precision, recall, F1, and accuracy (Acc) as the evaluation metrics for anomaly detection effectiveness

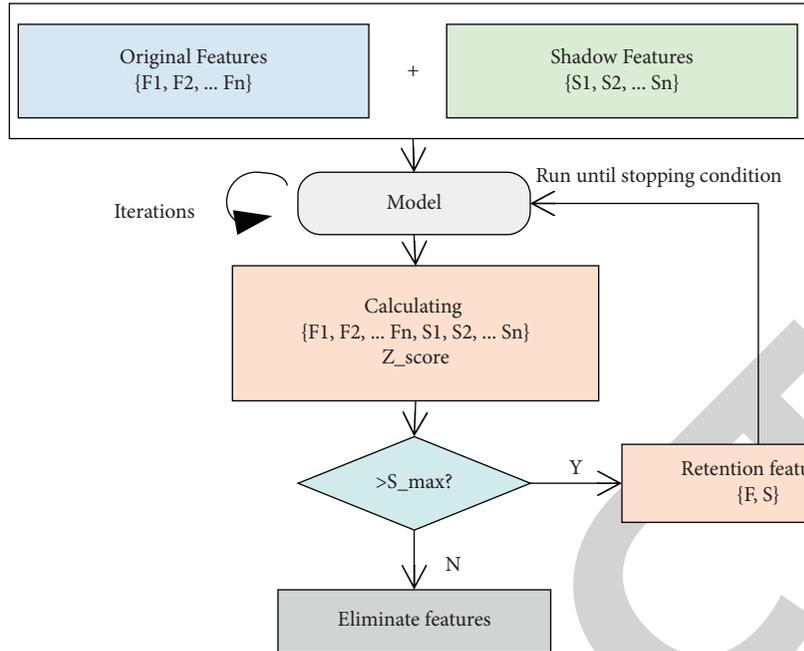


FIGURE 2: Boruta algorithm flowchart.

Build\_an\_extra\_tree\_ensemble(D)

Input: Train set  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$

Output: Extreme Trees  $T = \{t_1, t_2, \dots, t_M\}$

(1) for  $i = 1$  to  $M$  do

(2) Generating decision trees,  $t_i = \text{Build\_an\_extra\_tree}(D)$

(3) Return Extreme Trees  $T$

(4) end for

Build\_an\_extra\_tree(D)

Input: Train data  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$

Output: Decision Tree  $t$

(1) if  $|D| < n_{\min}$  or all candidate attributes in  $D$  are constant or output variables in  $D$  are constant then

(2) Return a leaf node

(3) else

(4) Randomly select  $K$  attributes from all candidate attributes  $\{a_1, a_2, \dots, a_k\}$

(5) Generate  $K$  split thresholds  $\{d_1, d_2, \dots, d_k\}$ , Among them  $d_i = \text{Pick\_a\_random\_split}(D, a_i)$

(6) According to  $\text{Score}(d_*, D) = \max_{i=1,2,\dots,K} \text{Score}(d_i, D)$ , Selecting the best test split threshold  $d_*$

(7) According to test split thresholds  $d_*$ , Divide the sample set  $D$  into two sub-sample sets  $D_l$  and  $D_r$ .

(8) Construct a left subtree  $t_l = \text{Build\_an\_extra\_tree}(D_l)$  and a right subtree  $t_r = \text{Build\_an\_extra\_tree}(D_r)$  using subsets  $D_l$  and  $D_r$  respectively

(9) Create a tree node based on  $d_*$ , with  $t_l$  and  $t_r$  as its left and right subtrees respectively, and return a decision tree  $t$

(10) end if

Pick\_a\_random\_split(D, a)

Input: Train data  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ , Attributes  $a$

Output: Divided attributes

(1) Calculate the minimum and maximum values of attribute  $a$  in the training set  $D$ , denoted respectively as  $a_{\min}^D$  and  $a_{\max}^D$

(2) Select a random splitting attribute  $a_c$  from  $[a_{\min}^D, a_{\max}^D]$

(3) Return to Split attributes  $[a < a_c]$

ALGORITHM 1: Extreme trees.

[13]. When conducting a multicategory classification anomaly detection study, we mainly use recall as the evaluation metric. It is not a good description of the performance of the classifier because the accuracy is high

for categories with many data samples and low for categories with few data samples but still gives a high overall accuracy. The confusion matrix of classification results is listed in Table 2.

TABLE 2: Confusion matrix of classification results.

		Predicted value		Indicator
		Benign	Attack	
True value	Benign	TP	FN	Recall=TP/(TP+FN)
	Attack	FP	TN	Precision=TP/(TP+FP)
		F1 = 2/[(1/precision) + (1/recall)]		Acc=(TP+TN)/(TP+TN+FP+FN)

## 4. Experimental Results and Analysis

**4.1. Experimental Environment.** The algorithm in the study is implemented in Python language. The operating system used for the experiments is Windows 10, 64 bit. The hardware environment is an Inter(R) Core (TM) i5-7200U CPU@ 2.50 GHz with 8G RAM.

**4.2. Dataset Processing.** In this study, the 14 attack types are divided into 6 domains, namely, DoS, PortScan, Bot, Brute Force, Web Attack, and Infiltration, and the detailed division is listed in Table 3. By counting the number of each attack domain, this study uses a pie chart to visualize the overall distribution of the data, as shown in Figure 3.

**4.2.1. Dataset Cleaning.** The rows in the CICIDS dataset where the NaN and Inf values were located were removed. The number of samples after deletion is listed in Table 4.

**4.2.2. Numerical Characters.** The dataset was marked with “benign” as “0” and the six attack types were marked as “1–6,” as in the new label column in Table 3.

**4.2.3. Data Normalization.** In order to reduce the problem of inconsistent impact weights between different dimensions of the data, this study uses a min-max normalization method to normalize the traffic data. The aim is to perform a linear transformation on the original data so that the results fall into the interval [0, 1]. The conversion function for the min-max normalization method is as follows:

$$X^+ = \frac{X - X_{\min}}{X_{\max} - X_{\min}}. \quad (3)$$

Here,  $X_{\min}$  is the minimum value of all the sample data and  $X_{\max}$  is the maximum value of all the sample data.  $X$  is the original sample data before conversion.  $X^+$  is the data after the conversion [14].

**4.3. Feature Selection Results.** To facilitate experimental validation, the CICIDS2017 dataset is divided into a training dataset and a testing dataset in the ratio of 7 : 3 in this study. The number of training and test sets after the division is listed in Table 5. The statistics on the dataset in Table 5 show that the number of the three attack types “bot,” “network attack,” and “infiltration” is relatively small compared to the other attack types. In order to avoid unbalanced distribution of samples, which would affect the performance of the

TABLE 3: Distribution of classes in the CICIDS2017 dataset.

Traffic type	New label	Label	No. of instances
Benign	0	Benign	2273097
		DDoS	128027
		Slowloris	5796
		Slowhttptest	5499
		Hulk	231073
DoS	1	GoldenEye	10293
		Heartbleed	11
Port scan	2	PortScan	158930
Bot	3	Bot	1966
		FTP	7938
Brute force	4	SSH	5897
		Brute Force	1507
Web attack	5	XSS	652
		SQLInjection	21
Infiltration	6	Infiltration	36
Total	—	—	2830743

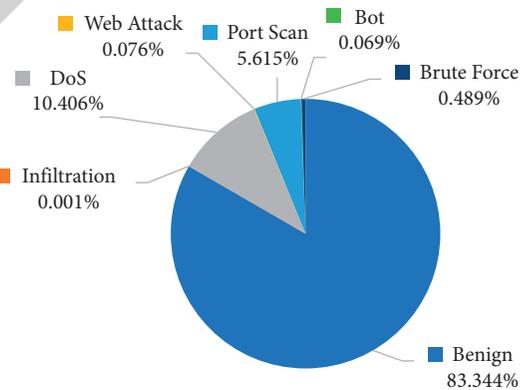


FIGURE 3: The overall distribution of the CICIDS2017 dataset.

TABLE 4: Number of CICIDS2017 data before and after cleaning.

	Before cleaning (original)	After cleaning
0	2273097	2271313
1	380699	379748
2	158930	158804
3	1966	1956
4	13835	13832
5	2180	2180
6	36	36

classification algorithm and thus degrade the detection, we used random oversampling to rebalance the dataset. The three types of attack types with a small number of samples

TABLE 5: Distribution of the training and testing sets after random oversampling.

Label	Traffic type	Train data	Expand train data	Test data
0	Benign	1589821	1589821	681492
1	DoS	265887	265887	113861
2	Port scan	111254	111254	47550
3	Bot	1956	6956	603
4	Brute force	9633	9633	4199
5	Web attack	2180	7180	646
6	Infiltration	36	5036	10

'Flow Duration', 'Total Fwd Packets', 'Total Backward Packets', 'Total Length of Fwd Packets', 'Total Length of Bwd Packets', 'Fwd Packet Length Max', 'Fwd Packet Length Min', 'Fwd Packet Length Mean', 'Fwd Packet Length Std', 'Bwd Packet Length Max', 'Bwd Packet Length Min', 'Bwd Packet Length Mean', 'Bwd Packet Length Std', 'Flow Bytes/s', 'Flow Packets/s', 'Flow IAT Mean', 'Flow IAT Std', 'Flow IAT Max', 'Flow IAT Min', 'Fwd IAT Total', 'Fwd IAT Mean', 'Fwd IAT Std', 'Fwd IAT Max', 'Fwd IAT Min', 'Bwd IAT Total', 'Bwd IAT Mean', 'Bwd IAT Std', 'Bwd IAT Max', 'Bwd IAT Min', 'Fwd PSH Flags', 'Fwd Packets/s', 'Bwd Packets/s', 'Min Packet Length', 'Max Packet Length', 'Packet Length Mean', 'Packet Length Std', 'Packet Length Variance', 'SYN Flag Count', 'PSH Flag Count', 'ACK Flag Count', 'URG Flag Count', 'Down/Up Ratio', 'Average Packet Size', 'Avg Fwd Segment Size', 'Avg Bwd Segment Size', 'Subflow Fwd Packets', 'Subflow Fwd Bytes', 'Subflow Bwd Packets', 'Subflow Bwd Bytes', 'Init\_Win\_bytes\_forward', 'Init\_Win\_bytes\_backward', 'act\_data\_pkt\_fwd', 'Active Mean', 'Active Std', 'Active Max', 'Active Min', 'Idle Mean', 'Idle Max', 'Idle Min'

FIGURE 4: Display the names of the features selected by Boruta algorithm.

TABLE 6: The performance of the methods in the study.

Method		Benign	DoS	Port scan	Bot	Brute force	Web attack	Infiltration
NB	Precision	0.87	0.81	0.76	0	0	0	0.01
	Recall	0.98	0.58	0	0	0	0	0.70
	F1	0.92	0.68	0	0	0	0	0.02
LGBM	Precision	0.88	0.61	0.83	0	0.08	0.10	0
	Recall	0.90	0.64	0.14	0	0.03	0.71	0
	F1	0.89	0.14	0.23	0	0.05	0.18	0
DT	Precision	0.99	0.99	0.99	0.41	0.95	0.66	0.17
	Recall	1	0.94	1	0.63	0.88	0.93	1
	F1	0.99	0.97	1	0.050	0.91	0.77	0.29
Xgboost	Precision	0.95	0.99	1	0.41	1	0.99	0.24
	Recall	1	0.93	0.46	0.25	0.46	0.86	1
	F1	0.97	0.96	0.63	0.31	0.63	0.92	0.38
RF	Precision	0.96	1	1	0.55	0.99	0.99	1
	Recall	1	0.94	0.56	0.26	0.77	0.89	1
	F1	0.98	0.97	0.71	0.35	0.87	0.94	1
DNN	Precision	0.92	0.85	0.98	0	0.98	0	0
	Recall	0.97	0.77	0.50	0	0.53	0	0
	F1	0.95	0.81	0.66	0	0.69	0	0
The method proposed Boruta-ET	Precision	1	1	0.99	0.87	1	0.98	1
	Recall	1	1	1	0.26	0.98	0.95	1
	F1	1	1	1	0.40	0.99	0.97	1

were randomly replicated, then the dataset obtained from each random sampling was superimposed by setting the "Sample\_strategy" parameter to the specified number, and we expanded the number by another 5000, thus obtaining a new balanced dataset, and the number of the extended training set is listed in Table 5.

In this study, by using the Boruta algorithm feature selection, the Borutapy software wrapper package in the python language was used to perform 100 iterations by filtering the features related to the dependent variable, and

finally 59 features were selected. The selected feature names are shown in Figure 4.

*4.4. Classification Performance Evaluation.* To validate the model Boruta-ET proposed in this study, we compared Boruta-ET with five other machine learning algorithms in terms of three metrics: precision, recall, and F1 value, and the results are listed in Table 6. We can see from the metrics in the table that our proposed model has a slightly lower

TABLE 7: Comparison of the accuracy of the model proposed in this study with other algorithms.

Author	Method	Acc (%)
Ahmim et al. [15]	Rep + RF	96.67
Wang [16]	PCA + SVM	92.91
Ustebay et al. [17]	AutoEncoder + DNN	96.71
Di and Li [18]	SVM + DBN	92.56
Zhang et al. [19]	Confidence + DNN	93.80
Other comparison algorithms in this article	NB	86.24
	LGBM	82.09
	DT	98.87
	Xgboost	95.58
	RF	96.44
	DNN	92.02
The method proposed in this article	Boruta-ET	99.80

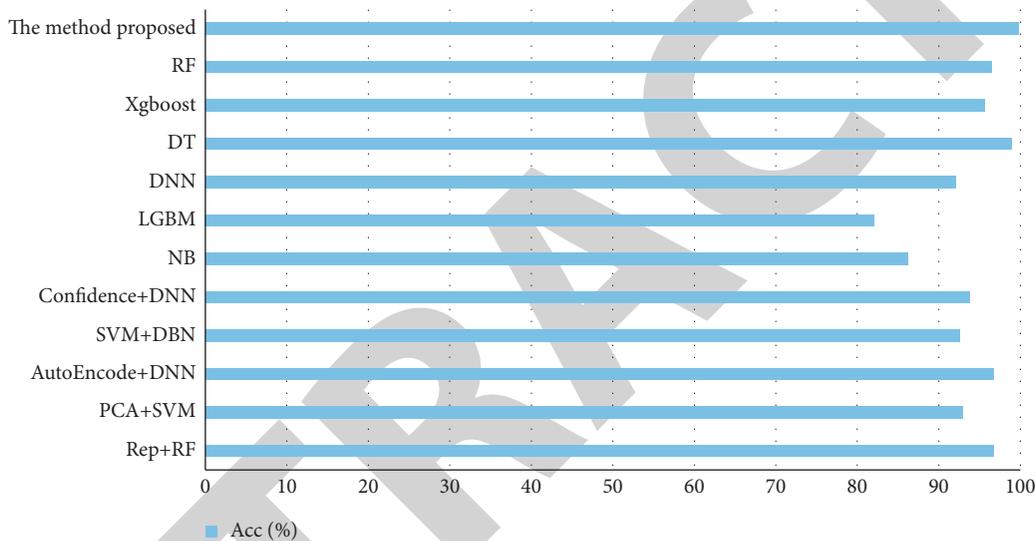


FIGURE 5: Comparison of the accuracy of the method proposed in this study with other methods.

recall when detecting Bot attack types, but the overall performance is excellent. We also conducted experiments on deep neural networks (DNNs) and the results show that the results are not as good as our proposed model. We also compared the overall accuracy of the model with published literature, and as can be seen from the accuracy rates in Table 7, the model in this study achieves an accuracy rate of 99.8%, which is the highest accuracy rate and the highest detection rate compared to other models proposed in the literature. In order to demonstrate the high performance of the method proposed in this study more visually, we use bar charts for this purpose. This is shown in Figure 5. In summary, the feasibility of the model proposed in this study for the detection of abnormal traffic is also very efficient.

## 5. Conclusion and Future Work

Through the analysis of the current state of research on network traffic anomaly detection technology, the problem of high traffic feature dimensionality is very common and a key issue that has attracted attention; however, not all features have a positive correlation on the results of anomaly

detection, and many useless and redundant features not only increase the computational complexity of traffic anomaly detection but also have a significant impact on the accuracy of detection. Boruta algorithm's aim is to select all feature sets associated with the dependent variable, as opposed to the traditional minimization of feature sets using a model-specific cost function. Boruta algorithm enables a global view of the impact of the dependent variable, leading to an increase in the efficiency of feature selection. In this study, we use a randomly oversampled balanced dataset, which can make the information learned by the model too specific and not general enough. We used the CICIDS2017 dataset to evaluate and compare existing models under similar experimental conditions. The model outperformed other existing methods in terms of accuracy, false positives, and recall. The results show that the model can be used effectively for intrusion detection, improving the accuracy of intrusion detection and the ability to identify the type of intrusion.

This study uses a random oversampling method to equalize the number of samples, and other sampling methods such as smote oversampling, undersampling, and hybrid sampling methods will be considered for

experimentation in future research. The Boruta algorithm is very comprehensive in terms of feature selection to find relevant features, but it is also expensive to train as it has to extend the dataset, is computationally expensive, and cannot be reduced by parallelization. In future research, the use of GUP acceleration will be considered to reduce the training time of the model. In future research, we plan to extend this work by deploying the experimental results to corresponding software systems to observe the performance of the software in real network environments.

## Data Availability

The datasets used and/or analyzed during the current study are available from the corresponding author upon reasonable request.

## Conflicts of Interest

The authors declare that they have no conflicts of interest.

## Acknowledgments

This work was supported by Central government guided local science and Technology Development Fund Project (no. 226Z0701G), the Natural Science Foundation of Hebei Province (no. F2022203026), Science and Technology Project of Hebei Education Department (nos. BJK2022029, QN2021145) and Innovation Capability Improvement Plan Project of Hebei Province (no. 22567637H).

## References

- [1] Y. Zhou, G. Cheng, and S. Jiang, "Building an efficient intrusion detection system based on feature selection and ensemble classifier," *Computer Networks*, vol. 174, pp. 1–17, 2020.
- [2] A. Muna, N. Moustafa, and E. Sitnikova, "Identification of malicious activities in industrial internet of things based on deep learning models," *Journal of Information Security and Applications*, vol. 41, pp. 1–11, 2018.
- [3] B. S. Harish and S. V. A. Kumar, "Anomaly based intrusion detection using modified fuzzy clustering," *International Journal of Interactive Multimedia and Artificial Intelligence*, vol. 4, no. 6, pp. 54–59, 2017.
- [4] M. Mazini, B. Shirazi, and I. Mahdavi, "Anomaly network-based intrusion detection system using a reliable hybrid artificial bee colony and AdaBoost algorithms," *Journal of King Saud University - Computer and Information Sciences*, vol. 31, no. 4, pp. 541–553, 2019.
- [5] A. Basati and M. M. Faghieh, "PDAE: efficient network intrusion detection in IoT using parallel deep auto-encoders," *Information Sciences*, vol. 598, pp. 57–74, 2022.
- [6] S. Zavrak and M. Iskefiyeli, "Anomaly-based intrusion detection from network flow features using variational autoencoder," *IEEE Access*, vol. 8, no. 99, pp. 108346–108358, 2020.
- [7] O. Alkadi, N. Moustafa, and B. Turnbull, "A deep blockchain framework-enabled collaborative intrusion detection for protecting IoT and cloud networks[J]," *IEEE Internet of Things Journal*, vol. 8, no. 12, 2020.
- [8] S. I. Popoola, B. Adebisi, and M. Hammoudeh, "Hybrid Deep Learning for Botnet Attack Detection in the Internet of Things Networks[J]," *IEEE Internet of Things Journal*, vol. 8, no. 6, pp. 2327–4662, 2020.
- [9] M. B. Kursu and W. R. Rudnicki, "Feature selection with Boruta package[J]," *Journal of Statistical Software*, vol. 36, no. 11, pp. 1–13, 2010.
- [10] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "Toward generating a new intrusion detection dataset and intrusion traffic characterization [J]," in *Proceedings of the 4th International Conference on Information Systems Security and Privacy (ICISSP)*, vol. 1, pp. 108–116, IEEE, Madeira, Portugal, January 2018.
- [11] P. Torres, C. Catania, S. Garcia, and C. G. Garino, "An Analysis of Recurrent Neural Networks for Botnet Detection behavior[C]," in *Proceedings of the 2016 IEEE Biennial Congress of Argentina (ARGENCON)*, pp. 1–6, IEEE, Buenos Aires, Argentina, June 2016.
- [12] S. Thakur, A. Chakraborty, R. De, N. Kumar, and R. Sarkar, "Intrusion detection in cyber-physical systems using a generic and domain specific deep autoencoder model," *Computers & Electrical Engineering*, vol. 91, no. 4, Article ID 107044, 2021.
- [13] Z. Chiba, N. Abghour, K. Moussaid, A. El omri, and M. Rida, "Intelligent approach to build a Deep Neural Network based IDS for cloud environment using combination of machine learning algorithms," *Computers & Security*, vol. 86, pp. 291–317, 2019.
- [14] B. Wei, L. Hu, Y. Zhang, and Y. Zhang, "Parts classification based on PSO-BP," in *Proceedings of the 2020 IEEE 4th Information Technology, Networking, Electronic and Automation Control Conference (ITNEC)*, pp. 1113–1117, IEEE, Chongqing, China, June 2020.
- [15] A. Ahmim, L. Maglaras, M. A. Ferrag, M. Derdour, and H. Janicke, "A novel hierarchical intrusion detection system based on decision tree and rules-based models," in *Proceedings of the 2019 15th International Conference on Distributed Computing in Sensor Systems (DCOSS)*, pp. 228–233, IEEE, Santorini, Greece, May 2019.
- [16] H. B. Wang, *Application Research of SVM Integration and Incremental Algorithm in Intrusion detection[D]*, Wuhan University of Technology, Wuhan, China, 2018.
- [17] S. Ustebay, Z. Turgut, and M. A. Aydin, "Cyber attack detection by using neural network approaches: shallow neural network, deep neural network and AutoEncoder," in *Proceedings of the International Conference on Computer Networks*, pp. 144–155, Springer, Gliwice, Poland, December 2019.
- [18] C. Di and T. Li, "Deep learning method for unknown network attack detection[J]," *Computer Engineering and Applications*, vol. 56, no. 22, Article ID 109116, 2020.
- [19] H. Zhang, Y. Li, Z. Lv, A. K. Sangaiyah, and T. Huang, "A real-time and ubiquitous network attack detection based on deep belief network and support vector machine," *IEEE/CAA Journal of Automatica Sinica*, vol. 7, no. 3, pp. 790–799, 2020.