



IJIRCCCE

e-ISSN: 2320-9801 | p-ISSN: 2320-9798



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

Volume 10, Issue 5, May 2022

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

Impact Factor: 8.165



9940 572 462



6381 907 438



ijircce@gmail.com



www.ijircce.com

Machine Learning Approach for Attack Detection on Network Traffic Data Using NSL-KDD Dataset

Prof. Chetan Kumar G S¹, Suraj Panduranga Vernekar²

Assistant Professor, Department of MCA, University BDT College of Engineering, Davangere, Karnataka, India¹

MCA Student, Department of MCA, University BDT College of Engineering, Davangere, Karnataka, India²

ABSTRACT: Distributed Denial of Service (DDoS) is a network cyber-attack designed to interrupt a targeted server's regular activity. Although advanced Machine Learning (ML) techniques were developed to classify DDoS, the assault remains a major Internet threat. Most new DDoS recognition methods are in two categories: Managed and Unsupervised. Availability of named network traffic datasets relies on managed DDoS recognition approaches. Network traffic analysis is used to identify attacks, even if the ML technique is unattended. Due to the large amounts of network traffic data, low identification accuracy, and a high number of false positives, both approaches are problematic. This study proposes a semi-controlled DDoS-detection method that tests entropy, co-clustering, knowledge benefit ratio and algorithm of the Random Forest network. The unregulated functionality of the device enables the usual meaningless traffic data for DDoS identification to be eliminated, helping to minimize wrong positives and improve precision. The monitored portion tends to reduce the unregulated part's fake optimistic rates and correctly distinguish DDoS traffic.

KEYWORDS: DDoS attacks, Machinelearning, Semi-supervised

I. INTRODUCTION

The DDoS assault is still a serious Internet hazard, despite the rapid advancement of information security solutions in recent years. A distributed Denial-of-Service (DDoS) attack is a deliberate effort to interrupt normal domain, service or network traffic by an enormous cascade of Internet traffic on the aim or surrounding networks. The key purpose of the attack is to deprive legal Internet users. How successful the assault will be will be determined by how quickly and how much data will be provided to the targeted victim.

In the field of machine learning (ML), data and knowledge are fed into computers in order to make them learn and behave like humans. The computer receives data as input and formulates replies based on an algorithm. There are three types of DDoS-based methodologies for machine learning: regulated, uncontrolled, and semi-monitored. Consistency and false positive rates may be improved by combining both regulated and unmonitored approaches, which work on both classified and unlisted data sets. In the Semi-supervised ML approach, entropy is used to determine the header characteristics of network traffic knowledge. The Unsupervised co-clustering technique breaks incoming network traffic into three groups. The average function header entropy between traffic data and the cluster is then used to determine the information-gain ratio for each cluster. For example, a Random Forest method uses a preprocessing procedure called anomalous to choose the data cluster that delivers the best information-gain ratio for preprocessing and classification. Using the NSL-KDD network traffic dataset, this technique may be evaluated.

II. PROBLEM STATEMENT

The inclusion of significant volumes of meaningless data in the incoming network traffic data for DDoS identification limits the efficiency of the supervised method. The curse of dimensionality issue arises due to high dimensional network traffic info, which prevents the unsupervised method from accurately detecting the attacks.

1.2.1 Existing System

DDoS detection models are built using named network traffic datasets in the present supervised ML technique. Unlike the first group, no classified dataset is needed in the unsupervised approaches to construct the model of detection. Based on the study of their underlying delivery features, the DDoS and the regular traffic are separated.

Disadvantages:

- Supervised ML methods do not anticipate new legal actions and assault behaviors. The existence of noisy data decreases the classifiers' efficiency.
- High false positive rates are the biggest downside of the unsupervised ML strategy.

1.2.2 Proposed System

In the semi-supervised technique, the data may be both labelled and unlabeled, allowing for the application of both supervised and unsupervised methodologies. In the unsupervised part, entropy computation, co-clustering, and the info-gain ratio are used. The Random forest ensemble classifier is the supervised part.

Advantages:

- Unregulated component of our plan to minimize meaningless and noisy daily traffic outcomes, decreasing false positive rates and improving the consistency of the monitored portion.
- The controlled part removes the unsupervised portion's false positive rate and correctly classifies DDoS traffic.

III. LITERATURE SURVEY

Bhuyan MH, Bhattacharyya DK, KalitaJK[1] is an empirical analysis of many of the major metrics of the results, including entropy of Hartley and Shannon, entropy, widespread entropy, Kullback-Leibler divergence and widespread assessment of the gap between the details.

Akilandeswari V. et al. utilize a Probabilistic Neural Flash crowd events may be used to respond to DDoS assaults. With fewer false positives, the system has a high detection rate.

DDoS detection based on ANN was proposed by Alan S. and colleagues (DDMA). Protocol Data Distribution and Multiple Access (DDMA) (DDMA). In order to identify three types of DDOS attacks, the authors employed three different MLP topologies, one for each context protocol, namely TCP, UDP, and ICMP. Detection of DDoS assaults by an unknown and recorded zero-day device is accurate[3].

Entropy-based approaches were used by Lui T, Wang Z, Wang H, and Lu K[4] to examine and identify actual IDS alerts. IP address dispersion, destination address invasion, source assaults, and IDS data alert time are all measured using Shannon entropy, which is used in conjunction with Reyni cross entropy to identify network attacks.

Boro D. Et al.[5] suggested a DyProSDdefense approach incorporating the merits of a predictive feature-based solution for coping with DDoS assaults through floods. The math module labels the traffic of the attacker and sends it to classifiers to label traffic as dangerous or normal.

A managed do-detection approach focused on the neural feed forward network with Mohamed I et al.[6] was proposed. This process comprises three key steps: (1) the compilation of incoming network traffic, (2) the selection of DoS identification features utilizing an unmonitored CFS, (3) the sortation of incoming network traffic in DoS or typical traffic.

A two-stage classification was introduced, based on RepTree algorithms and network intrusion detection subsets[7]. They are theoretically liable for splitting inbound traffic into three types: TCP, UDP or Other, and labeling it into regular or irregular traffic. A second level multi-class algorithm is used to identify the attack class to select the right behaviour. Two public sources are used for analysis, UNSW-NB15 and NSL-KDD.

Ali S.B. A groundbreaking Sugeno-style adaptive neuro-fuzzy classifier community is recommended[8] for utilizing Marlboost effective DDoS recognition boosting technologies. The proposed approach was tested for fair efficiency on the NSL-KDD dataset.

Mohiuddin A. AbdunNaser M. Presence [9] implemented a DDoS co-clustering recognition technique. The co-clustering algorithm was generalized by writers to adopt categorical characteristics. The technique was tested and the KDD cup 99 data set was successful.

The Van C. Van C. [10] introduced a modern one-class learning approach for combining measurements of the anomaly detection density and vehicle encoders. Authors have tested their NSL-KDD dataset framework and provided satisfactory results.

V. Jaiganesh, Dr.P. Sumathi, S. Mangayarkarasi[11], classed attacks as machine-learning and BPN technics in 4 groups: DoS, Demo, U2R, R2L. The detection rate for DoS risks is 78.15%.

The 8-style BPN attack data was qualified by ChangjunHan[12], Yi Lv, and Dan Yang, Yu Hao. 1325 instructions and 1245 theoretical relations. Their findings are: 80.5% identification rate, 7.4% false alarm rate, 11.3% absence. Sufyan T. Faraj et al. are eligible for usual and abnormal BPN detection and distinction in[13] initial instances. Abnormal events are split down into five distinct groups. Identification rate and false positive rate in various

cases are measured. The test collection recognition levels for the detection of usual and pathological cases is approximately 90% and for the classification into DoS, U2R and R2L is approximately 60-85%.

Eth. Eth. Eth. Eth. and Mukhopadhyay. For DoS, U2R, Inquire, U2L, and BPN neural network versions, al[14]. The method has an overall performance of 73.9% for the latest test range and 95.6% for stage 1. In neural SVM and MLP anomaly analysis Hua TANG and Zhuolin CAO are included. They compared precision for DoS, U2R, Study, U2L attack groups and found that the performance of the neural network is greater than SVM. Vladimir Bukhtoyarovf and Eugene Semenkin used the ensemble approach of a neural network. Their research centered on classifying sample attacks using combined use of specialized neural networks. They noticed a 99.87% ID rate for test attacks, but one of the IDS problems needed significant planning time.

The performance of the network for intrusion prevention typically focuses on the propagation properties of the data used to forecast network traffic. Two main classes of unsupervised approaches and controlled approaches represent DDoS literature recognition strategies. Unmonitored approaches frequently suffer from high false positives and managed approaches cannot handle vast amounts of network traffic data based on benchmark datasets used, and their performance is often restricted by noisy, irrelevant network data. It is also necessary to combine regulated and unregulated methods to address DDoS recognition concerns.

IV. SYSTEM ARCHITECTURE

One way to conceptualise and describe the behaviour of a whole system is via the use of a model called System Architecture. As a formally conceived idea, a system's formal notion and representation is represented by an architectural overview.

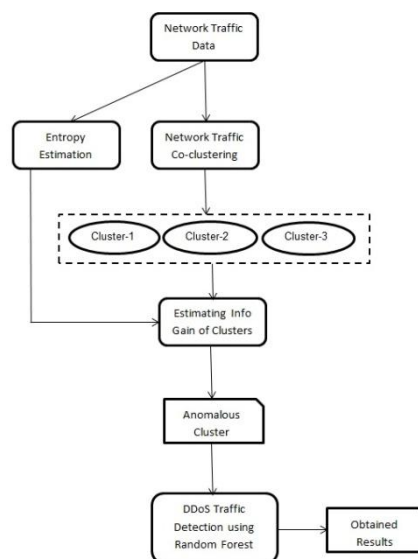


Fig. 4.1.1 System Architecture

Fig 4.1.1 Represents the device design for the solution suggested. It consists of several modules that are interrelated and function together to execute the framework.

Network Traffic Data

The proposed work contains traffic details from NSL-KDD[13]. NSL-KDD is a compilation of data proposed to resolve some of the essential issues of the KDD'99 dataset. Although the current iteration of KDD is not completely representative of real networks, it can also be used as a big dataset to help researchers validate numerous nuanced detection methods due to the scarcity of public data sets for network-based IDSs.

NSL-KDD dataset includes descriptions of attack. It has 42 functions: main features, content and traffic features, grouped into three groups. This data collection contains a total of 125973 training records and 22554 study records.

Compared with the original KDD data collection, the NSL-KDD data set presents the following benefits:

- It does not have duplicate records in the train collection, so more regular records would not be skewed against the classifiers.
- The proposed test sets contain no redundant data; hence, the performance of the learner is not influenced by techniques for higher frequency detection thresholds.

- As the percentage of records in the original data set increases, the number of records in each category picked decreases. A variety of computer-aided teaching techniques may be evaluated more readily because of this.
- The amount of train data and test sets is fair enough that the whole range of experiments can be conducted economically without the need to pick a single item arbitrarily. As a consequence, the findings of analyzing several research articles will be accurate and equivalent.

Entropy Estimation

Entropy is first calculated for FSD traffic data. Entropy estimate for the flow size distribution (FSD), the source/destination packet count, and the source/destination byte count are all performed using various functions and functions. Like the NSL-KDD dataset, it contains source bytes and destination bytes as part of its two FSD characteristics. In the event of a DDoS assault, zombie hosts would flood the target with a massive volume of packets, necessitating the FSD capability.

Network Traffic Co-clustering

The next move is to split network traffic data into three clusters, i.e. the Spectral co-clustering algorithm. Network traffic separation is targeted at reducing the amount of data to be categorized by removing the usual sorting cluster. At times, the latest unseen intermittent traffic accidents lead to raising the false positive rate and reducing classification accuracy. Eliminating irregular network traffic disruptive data for classification is also useful for low false positive rates and consistency of classifications.

Estimating Info Gain of Clusters

Based on the FSD functionality, calculating the data gain ratio allows it possible to differentiate between the two clusters that maintain more DDoS assault details and the cluster of regular traffic. The lower data acquisition ratio is then considered normal, and the other clusters are regarded as strange.

DDoS Traffic Detection using Random Forest

The data in the anomalous cluster is preprocessed for classification by taking care of missing information, encoding categorical data, and function scaling.

The representational problem of the unvarying decision tree is addressed by ensemble-based trees like Random Forest, which better reflect attack data. Classification relies on composite trees.

V. RESULTS AND DISCUSSION

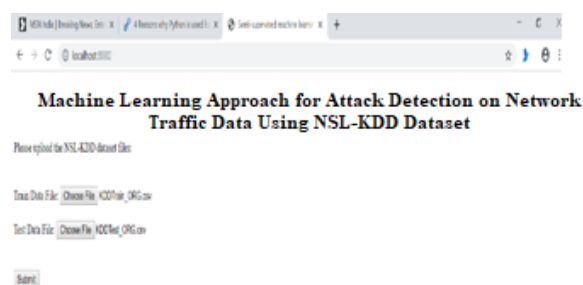


Fig.5.1.1 Screenshot of web page to upload Train and Test dataset

NSL-KDD training and testing datasets may be found here: Fig 5.1.1. There is a "Submit" button, which may be used to submit a dataset and proceed with the rest of the calculations.

The size of the data in each cluster after clustering is shown in Figure 5.1.2. Finally, each cluster's information gains are estimated to assist reduce the unnecessary traffic data.

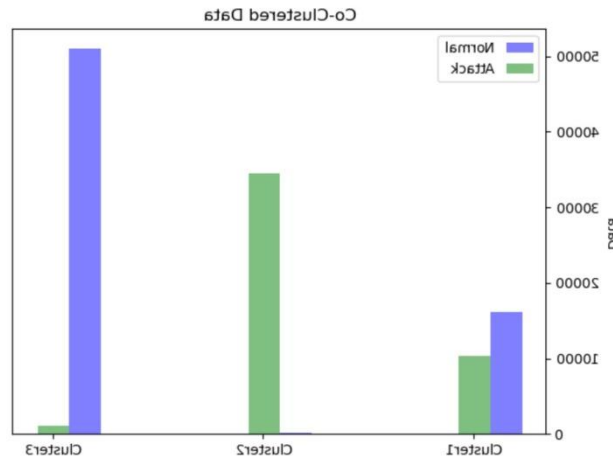


Fig.5.1.2 Screenshot of web page to upload Train and Test dataset

Average Entropy:0.492034015598514

X	shannon_entropy	normalized_entropy
src_bytes	6.1434157820153334	0.5248062578263939
dst_bytes	6.056304622502039	0.45926177337063406

Average Entropy:0.46945053754513505

X	shannon_entropy	normalized_entropy
src_bytes	5.193116096657044	0.493659283632879
dst_bytes	4.162444471237995	0.4452417914573911

Average Entropy:0.006201544992091686

X	shannon_entropy	normalized_entropy
src_bytes	0.04174545855797314	0.00670198646521364
dst_bytes	0.03381174750526495	0.005701103518969732

Average Entropy:0.7765967650261083

X	shannon_entropy	normalized_entropy
src_bytes	8.76464699816996	0.7591748808924171
dst_bytes	10.43786249240959	0.7940186491597994

Fig.5.1.3 Screenshot of text file having entropy values

Fig 5.1.3 illustrates the Shannon entropy, normalised entropy, and average entropy values of the whole dataset and the three clusters. A text file is used to keep track of these values.

Cluster ID	Info Gain
1	0.3823283942617106
2	0.4901388405309652
3	0.13424552629960462

Fig.5.1.4 Screenshot of text file having info-gain values of clusters

Fig 8.1.4 represents the calculated information gain ratio value of each cluster which is stored in a text file.

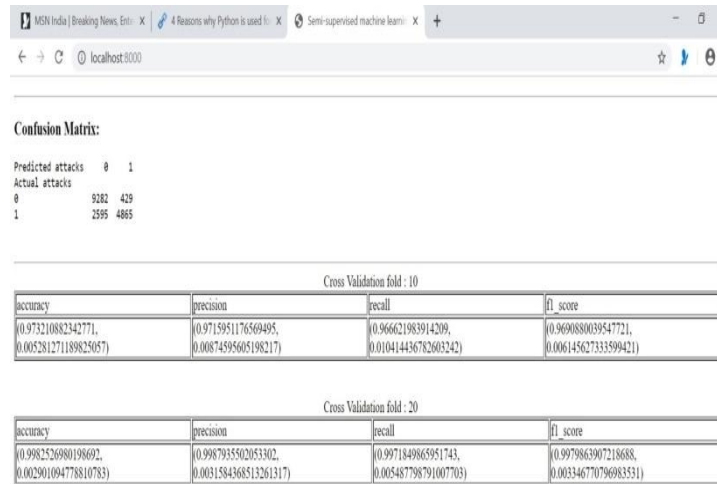


Fig.5.1.5 Screenshot of result page

Fig 5.1.5 illustrates the result page where a confusion matrix and k-fold cross validation results are created. Tables of test data are utilised to identify the actual values and explain the output model using an uncertainty matrix. We see the four basic terms: real positives, real negatives, phoney positive, and false negative, all in capital letters.

For estimating the accuracy of machine learning models, cross-validation is a useful mathematical technique. K-fold cross validation is a method for estimating the model's capacity to handle fresh data. The accuracy, consistency, recall, and f1-score in each fold were all assessed, as shown in Fig. 5.1.55.

- **Accuracy** - The most intuitive indicator of success is precision and it is literally a proportion of correctly expected measurement to overall observations.

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+FN+TN}$$
- **Precision** - Accuracy is the percentage of positive observations correctly estimated to the overall positive observations predicted.

$$\text{Precision} = \frac{TP}{TP+FP}$$
- **Recall (Sensitivity)** - Recall is the percentage of optimistic findings accurately forecast to all observations in the real class - yes.

$$\text{Recall} = \frac{TP}{TP+FN}$$
- **F1 score** - The F1 Score is the Precision and Recall Weighted Average.

$$\text{F1 Score} = 2 * (\text{Recall} * \text{Precision}) / (\text{Recall} + \text{Precision})$$

Classification Report

	precision	recall	f1-score	support
Normal Data	0.78	0.96	0.86	9711
Attack Data	0.92	0.65	0.76	7460
micro avg	0.82	0.82	0.82	17171
macro avg	0.85	0.80	0.81	17171
weighted avg	0.84	0.82	0.82	17171

Fig.5.1.6 Classification report

Fig 5.1.6 represents the Classification report of the proposed approach.

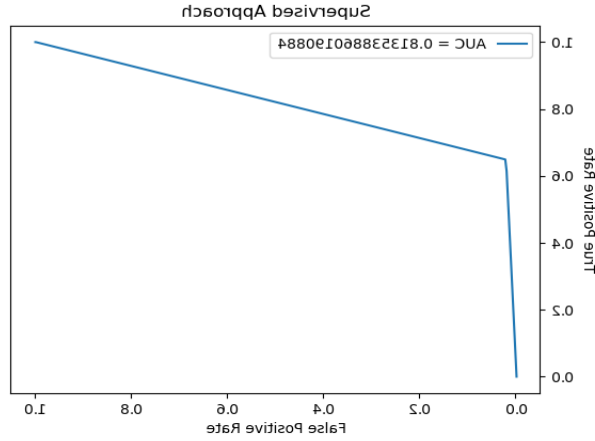


Fig.5.1.7 Graph of Supervised approach

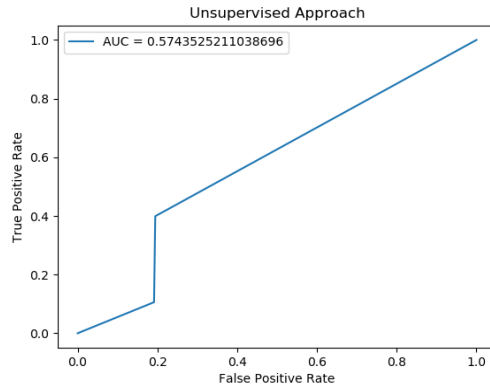


Fig5.1.8 Graph of Unsupervised approach

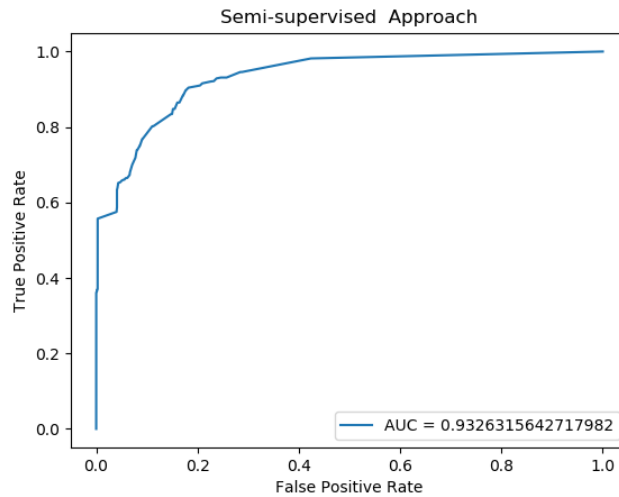


Fig.5.1.9 Graph of Semi-supervised approach

Figure 5.1.9 compares the suggested strategy to the supervised and unsupervised approaches given in Figures 8.1.7 and 8.1.8, respectively, and exhibits an improved accuracy and reduced false positive rate.

VI. CONCLUSION

Semi-Supervised DDoS Detection ML Methodology is the primary purpose of this research. The entropy of network traffic is analysed by an estimate of entropy. Clusters of traffic data are created using the co-clustering technique. The average entropy of the network header functions for the current dataset and each cluster is then used to construct a knowledge-benefit ratio. Data clusters with a high gain ratio are dubbed anomalous and are selected from the Random Forest method utilising ensemble classifiers prior to and during classification. The results are adequate in terms of accuracy and false positive rate when compared to specialist DDoS tactics. Since the good performance of the offered solution using publicly available benchmark data sets must be verified in real-world circumstances, the solution We need to test the proposed technique against a variety of DDoS instruments in the real world.

REFERENCES

- [1] Bhuyan MH, Bhattacharyya DK, Kalita JK (2015) An empirical evaluation of information metrics for low-rate and high-rate ddoattack detection. *Pattern RecognLett* 51:1–7
- [2] Akilandeswari V, Shalinie SM (2012) Probabilistic neuralnetwork based attack traffic classification. In: 2012 fourthinternational conference on advanced computing (ICoAC). IEEE, pp 1–8
- [3] Saied A, Overill RE, Radzik T (2016) Detection of knownand unknown ddo attacks using artificial neural networks. *Neurocomputing* 172:385–393
- [4] Liu T, Wang Z, Wang H, Lu K (2014) An entropy-based method for attack detection in large scale network. *Int J ComputCommunControl* 7(3):509–517
- [5] Boro D, Bhattacharyya DK (2016) Dyprosd: a dynamic protocol specific defense for high-rate ddo flooding attacks. *MicrosystTechnol* 23:1–19
- [6] Idhammad M, Afdel K, Belouch M (2017) Dos detection methodbased on artificial neural networks. *Int J AdvComputSciAppl(ijacsa)* 8(4):465–471
- [7] Mustapha B, Salah EH, Mohamed I (2017) A two-stage classifierapproach using reptree algorithm for network intrusion detection. *Int J AdvComputSciAppl(ijacsa)* 8(6):389–394
- [8] Boroujerdi AS, Ayat S (2013) A robust ensemble of neurofuzzyclassifiers for ddoattack detection. In: 2013 3rdinternational conference on computer science and networktechnology (ICCSNT). IEEE, pp 484–487
- [9] Ahmed M, Mahmood AN (2015) Novel approach for networktraffic pattern analysis using clustering-based collective anomalydetection. *Ann Data Sci* 2(1):111–130
- [10] Nicolau M, McDermott J et al (2016) A hybrid autoencoder anddensity estimation model for anomaly detection. In: *Internationalconference on parallel problem solving from nature*. Springer, pp717–726
- [11] Jaiganesh V., Sumathi P. and Mangayarkarasi S., "An Analysis of Intrusion Detection System using Back Propagation Neural Network", IEEE 2013 publication
- [12] Han C., Yi Lv, Yang D., Hao Y., "An Intrusion Detection System Based on Neural Network", 2011 International Conference on Mechatronic Science, Electric Engineering and Computer, August 19-22, 2011, Jilin, China, IEEE Publication
- [13] Faraj S, Al-Janabi and Saeed H, "A Neural Network Based Anomaly Intrusion Detection System", 2011 Developments in E-systems Engineering, DOI 10.1109/DeSE.2011.19, IEEE publication
- [14] Mukhopadhyay I, Chakraborty M, Chakrabarti S, Chatterjee T, "Back Propagation Neural Network Approach to Intrusion Detection System", 2011 International Conference on Recent Trends in Information Systems, IEEE publication



INNO  SPACE
SJIF Scientific Journal Impact Factor

Impact Factor: 8.165

 **doi**[®]
CROSS **ref**

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

 9940 572 462  6381 907 438  ijircce@gmail.com



www.ijircce.com

Scan to save the contact details