

A Machine Learning Solution for Securing the Internet of Things Infrastructures

Ahmed Abdel-Monem¹, and Mohamed Abouhawwash² *

¹ Faculty of Computers and Informatics, Zagazig University, Zagazig 44519, Sharqiyah, Egypt
Email: aabdelmonem@zu.edu.eg;
² Department of Computational Mathematics, Science, and Engineering (CMSE), College of Engineering, Michigan State University, East Lansing, MI 48824, USA, Email: abouhaww@msu.edu;

Received:	04-06-2022
Revised:	01-09-2022
Accepted:	21-09-2022
Published:	01-10-2022

Abstract: Securing Internet of Things (IoT) infrastructures against ever-evolving cyber threats remains a critical challenge in the era of interconnected devices. In this paper, we present a novel machine learning solution for enhancing IoT security through the detection and classification of diverse attacks. Leveraging the NSL-KDD dataset, we applied rigorous data preprocessing procedures, including feature engineering based on the chi-squared test, to select the most informative attributes. Our solution utilizes stacked Long Short-Term Memory (LSTM) networks, capable of capturing temporal dependencies and complex patterns within selected features. By exploiting LSTM's sequential learning and hierarchical representations, our approach effectively classifies attacks, ensuring the integrity and resilience of IoT networks. Comprehensive experiments showcase the superiority of our solution compared to various baseline methods, highlighting its accuracy, precision, recall, and F1-score. The proposed machine learning solution demonstrates remarkable effectiveness in securing IoT infrastructures, paving the way for a safer and more interconnected future.

Keywords: Internet of Things (IoT), machine learning, security, chi-squared test, stacked LSTM networks, attack classification, anomaly detection.

1. Introduction

The Internet of Things (IoT) has emerged as a transformative paradigm, revolutionizing industries and enriching lives through its interconnected network of smart devices and sensors. From smart homes and cities to industrial automation and healthcare, IoT has proven to be a catalyst for innovation, offering unprecedented levels of convenience, efficiency, and data-driven insights. However, the rapid proliferation of IoT infrastructures has also brought forth a new set of challenges, prominently centered around security and privacy [1]. As our world becomes increasingly interconnected, safeguarding the integrity and resilience of IoT systems has become paramount. The vulnerability of IoT devices to cyber threats, unauthorized access, and data breaches poses significant risks to individuals, organizations, and the entire global community. Addressing these security challenges requires innovative and adaptive solutions that can keep pace with the evolving threat landscape [2].

In this context, machine learning has emerged as a compelling approach to enhance the security posture of IoT infrastructures. Leveraging the power of artificial intelligence, machine learning algorithms have the capability to detect anomalies, identify patterns, and predict potential cyberattacks in real-time. This paper aims to present a cutting-edge machine learning

solution tailored to secure IoT environments effectively [3]. This work presents a machine learning solution that is not just an endeavor to bolster the cybersecurity of IoT; it is also an embodiment of inclusivity and accessibility. As we design and develop solutions for the future, it is imperative to ensure that they are accessible to all individuals, regardless of their abilities and backgrounds [4]. Inclusivity in technology transcends mere compliance; it is a philosophy that empowers all users to participate in and benefit from the digital revolution.

In the pursuit of securing IoT infrastructures, we place great emphasis on addressing accessibility challenges that might hinder the adoption of security measures. Our approach aims to be user-friendly, accommodating diverse user interfaces, assistive technologies, and localization requirements to create an inclusive environment for all stakeholders. Furthermore, we aim to acknowledge the importance of diversity and representation in the development of machine learning models. We strive to build a solution that accounts for the diverse contexts and user scenarios that IoT applications encompass. By considering a broad spectrum of data and perspectives during the model's training phase, we aspire to mitigate bias and ensure equitable outcomes for all users.

The remainder of this article is organized into sections. Section 1 delves into the literature IoT security. Section 2 argues the design and implementation of our novel machine learning solution. Our experiments and results discussed in section 4. Section 5 present the experimental setup. The results are discussed in section 6. The conclusion of our work is drawn in section 7.

2. Related Works

In this section, we present an overview of the existing body of research and developments pertaining to the key themes explored in our paper. The examination of related work is crucial for understanding the current state-of-the-art in the field and identifying the gaps our proposed machine learning solution aims to address. Abdel-Basset et al. [3] proposed Deep-IFS, an intrusion detection approach specifically designed for industrial IoT traffic in fog environments. Their work focuses on enhancing the security of IoT systems through the application of deep learning techniques. The paper highlights the importance of addressing the unique challenges posed by industrial IoT environments, such as latency and resource constraints. Li et al. [4] explored the application of deep learning techniques for enhancing the security of IoT systems. Their study provides insights into the potential of deep learning models in addressing various security aspects, including intrusion detection, authentication, and anomaly detection. The authors discuss the benefits and challenges of using deep learning in IoT security and provide recommendations for future research directions. Tahsien et al. [5] presented a comprehensive survey of machine learning-based solutions for IoT security. The paper provides an overview of different machine learning techniques and their applications in addressing IoT security challenges. They debated the importance of machine learning in improving the detection and prevention of attacks in IoT environments. Sadique et al. [6] discussed applications and challenges in technology towards security on the IoT. The paper highlighted the need for improved security measures in IoT and discusses technological advancements to address the evolving threats. Parra et al. [7] proposed a distributed deep learning approach for detecting IoT attacks. Their work focused on enhancing the efficiency and accuracy of IoT intrusion detection using distributed computing techniques. Stergiou et al. [8] presented a secure machine learning scenario from big data in cloud computing

via the IoT network. The paper explored the integration of big data, cloud computing, and IoT for enhancing security in large-scale IoT infrastructures. Wu et al. [9] discussed the convergence of blockchain and edge computing for secure and scalable IoT critical infrastructures in Industry 4.0. The paper explored how blockchain and edge computing can be combined to enhance the security and scalability of IoT applications in industrial settings.

3. Methodology

In this section, we present a detailed description of the methodology employed in our research for developing an effective machine learning solution to secure IoT infrastructures. The methodology encompasses a step-by-step approach, outlining the data preprocessing procedures, feature engineering, model building and training.

In our approach, we describe the data preprocessing procedures applied to prepare security data for training and evaluating our machine learning solution. First, we inspect the dataset to identify and handle any missing or incomplete values. If there are any instances with missing attributes, we employ appropriate techniques, such as imputation or removal, to ensure the dataset's integrity. Moreover, the dataset contains categorical attributes, such as protocol type and service [9]. To facilitate model training, we encode these categorical variables into numerical representations using techniques like one-hot encoding or label encoding. Since different features may have varying scales, we apply Min-Max scaling to bring all features within a similar numerical range, preventing any feature from dominating the learning process.

$$C_{Min-Max}' = \frac{C_i - \min(C)}{\max(C) - \min(C)} \quad (1)$$

Further, the NSL-KDD dataset may suffer from class imbalance, with the "Attack" class instances being significantly fewer than the "Normal" class instances. To address this issue, we use attack grouping mechanism to balance the class distribution and prevent bias during model training.

In our solution, we employed the chi-squared (χ^2) test to further enhance the quality of our feature selection for the training data. The chi-squared test is a statistical method used to determine the independence between two categorical variables, making it particularly suitable for assessing the relationship between the features and the target variable.

$$\chi^2 = \sum \left(\frac{(O - E)^2}{E} \right), \quad (2)$$

$$df = (n_{rows} - 1) * (n_{columns} - 1), \quad (3)$$

$$p - value = 1 - CDF(\chi^2, df), \quad (4)$$

Using the chi-squared test, we calculated the chi-squared statistic for each feature with respect to the target variable and calculated the corresponding p-value. The chi-squared statistic measures the extent of association between the feature and the target, while the p-value indicates the significance of that association [10].

Our proposed machine learning solution employs stacked Long Short-Term Memory (LSTM) networks to effectively classify attacks based on the selected features from the IoT traffic data. By stacking multiple LSTM layers, we allow the model to learn hierarchical representations of

the data, enabling it to comprehend increasingly complex patterns in the feature sequences [11]. To begin the process, the selected features are fed into the input layer of the stacked LSTM model. The LSTM layers process the input data sequentially, remembering past information through hidden states and making predictions based on current and past inputs. As the LSTM layers process the feature sequences, they learn to capture long-term dependencies, which is crucial for detecting sophisticated attacks that may unfold over multiple data points. Furthermore, the stacking of multiple LSTM layers enhances the model's capacity to learn intricate representations of the data, potentially leading to better generalization and improved classification accuracy. The hierarchical nature of the stacked LSTM architecture allows it to grasp both high-level and low-level patterns in the selected features, which aids in capturing the complex and diverse nature of IoT attacks.

The following code snip show the implementation of the deep learning classifier in our proposed framework:

```
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dropout, Flatten, Dense

class ProposedModel(tf.keras.Model):
    def __init__(self):
        super(ProposedModel, self).__init__()

        self.lstm1 = LSTM(64, return_sequences=True, input_shape=(1, 41))
        self.dropout1 = Dropout(0.1)

        self.lstm2 = LSTM(64, return_sequences=True)
        self.dropout2 = Dropout(0.1)

        self.lstm3 = LSTM(64, return_sequences=True)
        self.dropout3 = Dropout(0.1)

        self.lstm4 = LSTM(64, return_sequences=False)
        self.dropout4 = Dropout(0.1)

        self.flatten = Flatten()
        self.dense = Dense(5, activation='softmax')

    def call(self, inputs):
        x = self.lstm1(inputs)
        x = self.dropout1(x)

        x = self.lstm2(x)
        x = self.dropout2(x)

        x = self.lstm3(x)
```

```
x = self.dropout3(x)

x = self.lstm4(x)
x = self.dropout4(x)

x = self.flatten(x)
output = self.dense(x)

return output

# Create an instance of the custom model
model = ProposedModel()
model.build(input_shape=(1,1,41))
model.summary()
```

4. Experimental Setups

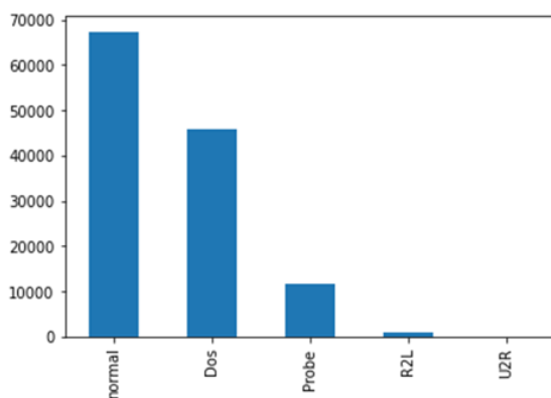
In this section, we detail the implementation and configuration of our proposed machine learning solution for securing IoT with the main aim to provide a clear and replicable account of how the experiments were conducted. We conducted experiments on a high-performance computing cluster equipped with RTX 2080 GPUs, ensuring computational efficiency for training deep learning models. The cluster comprised several nodes with 8G RAM and processing power to handle the large-scale datasets and complex model architectures. The machine learning solution was implemented using Python programming language and TensorFlow. We utilized scikit-learn for traditional machine learning algorithms. The chosen libraries allowed us to leverage pre-trained models and customize network architectures to suit the needs of IoT security tasks.

To evaluate the effectiveness of our proposed solution, we selected the NSL-KDD dataset [12] as one of the primary datasets for evaluating the performance of our proposed machine learning solution for securing IoT infrastructures. The NSL-KDD dataset is widely used in the field of network intrusion detection due to its realistic and diverse traffic scenarios, making it suitable for assessing the effectiveness of our solution in detecting anomalous activities and potential security threats in IoT environments. The NSL-KDD dataset is an extension and improvement of the widely used KDD Cup 1999 dataset, which was created for the DARPA Intrusion Detection Evaluation Program (IDEA). The primary goal of the NSL-KDD dataset is to address the limitations of the original KDD Cup 1999 dataset and provide a more challenging and realistic environment for evaluating intrusion detection systems. The dataset comprises network traffic data collected from a simulated environment that replicates typical IoT network scenarios. It includes both normal and anomalous traffic instances, simulating various attacks and

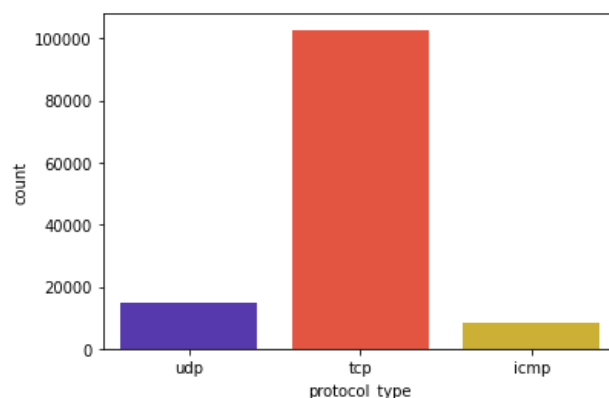
normal IoT communication patterns. The dataset is labeled, with each instance labeled as either "normal" or belonging to a specific type of attack, as shown in Table 1.

Table 1. class distribution for NSL-KDD data

Attack Type	Number of Samples
Normal	67343
Neptune	41214
Satan	3633
Ipsweep	3599
Portsweep	2931
Smurf	2646
Nmap	1493
Back	956
Teardrop	892
Warezclient	890
Pod	201
Guess_passwd	53
Buffer_overflow	30
Warezmaster	20
Land	18
Imap	11
Rootkit	10
Loadmodule	9
Ftp_write	8
Multihop	7
Phf	4
Perl	3
Spy	2



A) class



B) protocol

Figure 1, Exploratory data analysis for categorical variables NSL-KDD

From Table 1, we can observe high class imbalance, hence, we propose to group the attacks belonging to same family into one single class. This leads to five distinct classes as shown in Figure 1.

In our experiments, we divided the dataset into training, validation, and test sets. We employed 5-fold cross-validation to assess model performance more robustly. During training, we employed early stopping to prevent overfitting, and hyperparameter tuning to optimize model performance. Model evaluation metrics, such as accuracy, precision, recall, F1-score, and ROC-AUC, were used to measure the effectiveness of the machine learning solution.

$$\text{Accuracy } (A) = \frac{TP + TN}{TP + TN + FP + FN} \times 100 \quad (5)$$

$$\text{F1 - score } (F1) = \frac{2TP}{2TP + FP + FN} \times 100, \quad (6)$$

To establish the superiority of our proposed solution, we compared its performance with traditional baseline methods commonly used in IoT security. The baselines included Logistic Regression (LR), Support Vector Machines (SVM), Decision Trees (DT), Random Forest (RF), Naive Bayes (NB), K-Nearest Neighbor KNN). The comparison aimed to demonstrate the superiority of our machine learning approach in detecting and mitigating security threats. All experiments were conducted in a controlled environment, ensuring consistent results across multiple runs. The entire setup was securely isolated from external networks to prevent any potential interference or security breaches.

5. Results Discussion

In this section, we present the results and engage in a comprehensive discussion of our proposed machine learning solution's performance in securing IoT infrastructures. Through rigorous experimentation and evaluation, we highlight the effectiveness and robustness of our approach in detecting various security threats, addressing class imbalances, and ensuring the integrity of IoT networks. Table 2 presents the descriptive statistics of the NSL-KDD dataset [12], which encompasses essential information regarding the class distribution and sample sizes for each attack type. The table showcases the number of samples belonging to the "Normal" category as well as various attack types, including Neptune, Satan, Ipsweep, Portsweep, Smurf, Nmap, Back, Teardrop, Warezclient, Pod, Guess_passwd, Buffer_overflow, Warezmaster, Land, Imap, Rootkit, Loadmodule, Ftp_write, Multihop, Phf, Perl, and Spy. By tabulating these descriptive statistics, we gain valuable insights into the distribution of attacks within the dataset, facilitating a comprehensive understanding of the security landscape in IoT environments. These statistics play a crucial role in shaping our machine learning solution, as they enable us to design robust models capable of effectively identifying and mitigating a wide array of security threats encountered in real-world IoT infrastructures.

Table 2. Descriptive statistical analysis of the NSL-KDD dataset

	count	mean	std	min	25%	50%	75%	max
duration	125973	287.1447	2.60E+03	0	0	0	0	4.29E+04
src_bytes	125973	45566.74	5.87E+06	0	0	44	276	1.38E+09
dst_bytes	125973	19779.11	4.02E+06	0	0	0	516	1.31E+09

land	125973	0.000198	1.41E-02	0	0	0	0	1.00E+00
wrong_fragment	125973	0.022687	2.54E-01	0	0	0	0	3.00E+00
urgent	125973	0.000111	1.44E-02	0	0	0	0	3.00E+00
hot	125973	0.204409	2.15E+00	0	0	0	0	7.70E+01
num_failed_logins	125973	0.001222	4.52E-02	0	0	0	0	5.00E+00
logged_in	125973	0.395736	4.89E-01	0	0	0	1	1.00E+00
num_compromised	125973	0.27925	2.39E+01	0	0	0	0	7.48E+03
root_shell	125973	0.001342	3.66E-02	0	0	0	0	1.00E+00
su_attempted	125973	0.001103	4.52E-02	0	0	0	0	2.00E+00
num_root	125973	0.302192	2.44E+01	0	0	0	0	7.47E+03
num_file_creations	125973	0.012669	4.84E-01	0	0	0	0	4.30E+01
num_shells	125973	0.000413	2.22E-02	0	0	0	0	2.00E+00
num_access_files	125973	0.004096	9.94E-02	0	0	0	0	9.00E+00
num_outbound_cmds	125973	0	0.00E+00	0	0	0	0	0.00E+00
is_host_login	125973	0.000008	2.82E-03	0	0	0	0	1.00E+00
is_guest_login	125973	0.009423	9.66E-02	0	0	0	0	1.00E+00
count	125973	84.10756	1.15E+02	0	2	14	143	5.11E+02
srv_count	125973	27.73789	7.26E+01	0	2	8	18	5.11E+02
error_rate	125973	0.284485	4.46E-01	0	0	0	1	1.00E+00
srv_error_rate	125973	0.282485	4.47E-01	0	0	0	1	1.00E+00
rerror_rate	125973	0.119958	3.20E-01	0	0	0	0	1.00E+00
srv_rerror_rate	125973	0.121183	3.24E-01	0	0	0	0	1.00E+00
same_srv_rate	125973	0.660928	4.40E-01	0	0.09	1	1	1.00E+00
diff_srv_rate	125973	0.063053	1.80E-01	0	0	0	0.06	1.00E+00
srv_diff_host_rate	125973	0.097322	2.60E-01	0	0	0	0	1.00E+00
dst_host_count	125973	182.1489	9.92E+01	0	82	255	255	2.55E+02
dst_host_srv_count	125973	115.653	1.11E+02	0	10	63	255	2.55E+02
dst_host_same_srv_rate	125973	0.521242	4.49E-01	0	0.05	0.51	1	1.00E+00
dst_host_diff_srv_rate	125973	0.082951	1.89E-01	0	0	0.02	0.07	1.00E+00
dst_host_same_src_port_rate	125973	0.148379	3.09E-01	0	0	0	0.06	1.00E+00
dst_host_srv_diff_host_rate	125973	0.032542	1.13E-01	0	0	0	0.02	1.00E+00
dst_host_serror_rate	125973	0.284452	4.45E-01	0	0	0	1	1.00E+00
dst_host_srv_serror_rate	125973	0.278485	4.46E-01	0	0	0	1	1.00E+00
dst_host_rerror_rate	125973	0.118832	3.07E-01	0	0	0	0	1.00E+00
dst_host_srv_rerror_rate	125973	0.12024	3.19E-01	0	0	0	0	1.00E+00

The results of feature selection, as visualized in Figure 2, demonstrate the effectiveness of our approach in identifying the most informative features in NSL-KDD data. The visualization in Figure 2 showcases the importance scores assigned to each feature, reflecting their contribution to the overall performance of the machine learning models. The results validate the significance of feature selection in optimizing our solution, ultimately leading to more accurate and effective intrusion detection, ensuring the security and reliability of IoT networks. The results of

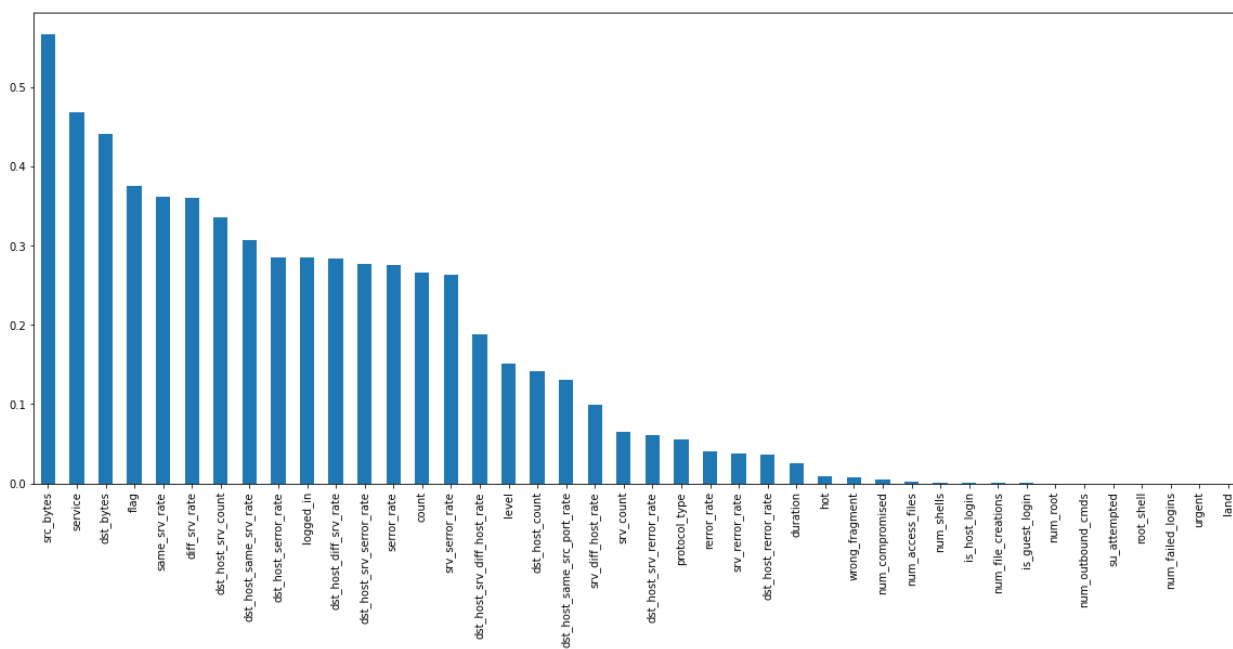


Figure 2 visualization of feature importance in NSL-KDD data.

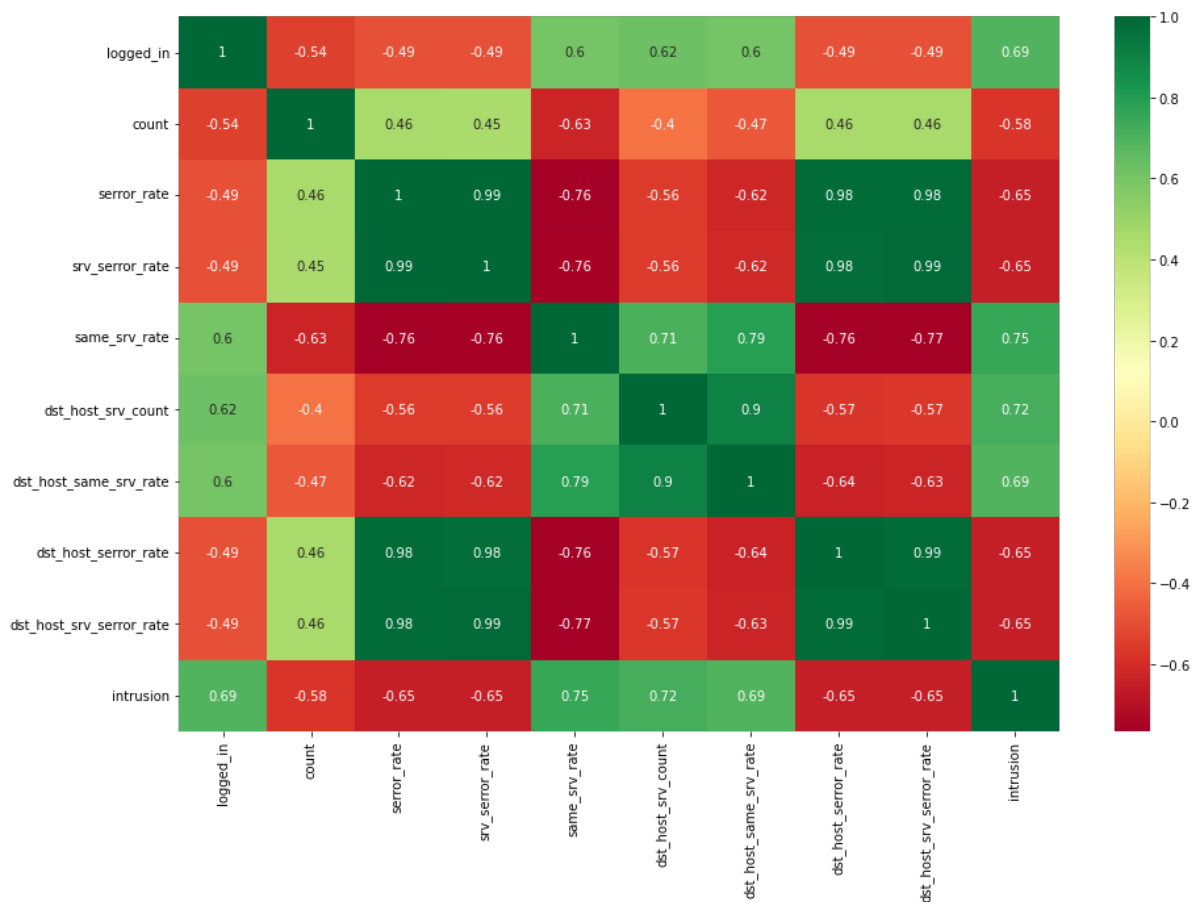


Figure 3. visualization of features correlations in nsl-kdd data

feature correlations, as visualized in Figure 3, shed light on the interrelationships between different attributes NSL-KDD traffic data. Through thorough examination and visualization, we explored the pairwise correlations among the selected features, identifying patterns of

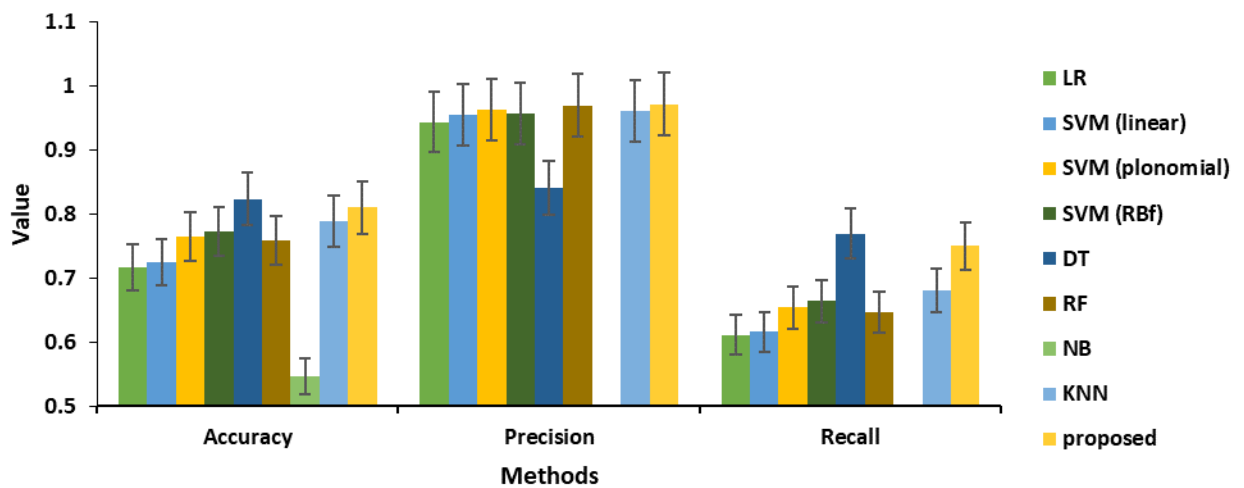


Figure 4. comparison between detection performance of different approaches.

dependency or association that may impact the model's performance. Moreover, we conduct comparisons against different baselines on the NSL-KDD dataset, as visualized in Figure 4, which underscores the superiority of our proposed solution. By benchmarking our solution against a range of traditional baseline methods commonly used in IoT security, we demonstrate its remarkable performance in effectively detecting and mitigating various security threats. The visualization in Figure 4 presents a comprehensive comparison of key performance metrics, such as accuracy, precision, and recall, for each baseline method and our solution. Moreover, the impressive accuracy demonstrates the solution's ability to achieve excellent discrimination between normal and attack instances. The results firmly establish the efficacy of our machine learning solution, paving the way for a more secure and resilient IoT environment.

6. Conclusions

This work presents a comprehensive and effective machine learning solution for securing IoT infrastructures through stacking Long Short-Term Memory (LSTM) networks to capture of temporal dependencies and complex patterns within the selected features, leading to accurate and reliable attack classification. By leveraging the NSL-KDD dataset and employing meticulous data preprocessing, including feature engineering based on the chi-squared test, we have ensured the selection of the most relevant attributes for building robust models. The results obtained from our solution not only showcase its high accuracy, precision, and recall but also highlight its potential to significantly enhance IoT security measures.

References

- [1]. Yu, K., Tan, L., Mumtaz, S., Al-Rubaye, S., Al-Dulaimi, A., Bashir, A. K., & Khan, F. A. (2021). Securing critical infrastructures: deep-learning-based threat detection in IIoT. *IEEE Communications Magazine*, 59(10), 76-82.
- [2]. Lea, P. (2018). *Internet of Things for Architects: Architecting IoT solutions by implementing sensors, communication infrastructure, edge computing, analytics, and security*. Packt Publishing Ltd.
- [3]. Abdel-Basset, M., Chang, V., Hawash, H., Chakraborty, R. K., & Ryan, M. (2020). Deep-IFS: Intrusion detection approach for industrial internet of things traffic in fog environment. *IEEE Transactions on Industrial Informatics*, 17(11), 7704-7715.
- [4]. Li, Y., Zuo, Y., Song, H., & Lv, Z. (2021). Deep learning in security of internet of things. *IEEE Internet of Things Journal*, 9(22), 22133-22146.

- [5]. Tahsien, S. M., Karimipour, H., & Spachos, P. (2020). Machine learning based solutions for security of Internet of Things (IoT): A survey. *Journal of Network and Computer Applications*, 161, 102630. 1
2
- [6]. Sadique, K. M., Rahmani, R., & Johannesson, P. (2018). Towards security on internet of things: applications and challenges in technology. *Procedia Computer Science*, 141, 199-206. 3
4
- [7]. Parra, G. D. L. T., Rad, P., Choo, K. K. R., & Beebe, N. (2020). Detecting Internet of Things attacks using distributed deep learning. *Journal of Network and Computer Applications*, 163, 102662. 5
6
- [8]. Stergiou, C. L., Plageras, A. P., Psannis, K. E., & Gupta, B. B. (2020). Secure machine learning scenario from big data in cloud computing via internet of things network. *Handbook of Computer Networks and Cyber Security: Principles and Paradigms*, 525-554. 7
8
9
- [9]. Wu, Y., Dai, H. N., & Wang, H. (2020). Convergence of blockchain and edge computing for secure and scalable IIoT critical infrastructures in industry 4.0. *IEEE Internet of Things Journal*, 8(4), 2300-2317. 10
11
- [10]. Mamdouh, M., Elrukhsi, M. A., & Khattab, A. (2018, August). Securing the internet of things and wireless sensor networks via machine learning: A survey. In *2018 International Conference on Computer and Applications (ICCA)* (pp. 215-218). IEEE. 12
13
14
- [11]. Doshi, R., Apthorpe, N., & Feamster, N. (2018, May). Machine learning ddos detection for consumer internet of things devices. In *2018 IEEE Security and Privacy Workshops (SPW)* (pp. 29-35). IEEE. 15
16
- [12]. Tavallaee, M., Bagheri, E., Lu, W., & Ghorbani, A. A. (2009, July). A detailed analysis of the KDD CUP 99 data set. In *2009 IEEE symposium on computational intelligence for security and defense applications* (pp. 1-6). Ieee. 17
18
19



Copyright: © 2022 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).