



INTERNATIONAL JOURNAL OF ADVANCE RESEARCH, IDEAS AND INNOVATIONS IN TECHNOLOGY

ISSN: 2454-132X

Impact Factor: 6.078

(Volume 9, Issue 6 - V9I6-1235)

Available online at: <https://www.ijariit.com>

Phishing mail detection using bidirectional LSTM

Vusa Vamsi Krishna

vamsikrishnavusa654@gmail.com

Sathyabama Institute of Science and Technology, Chennai, Tamil Nadu

ABSTRACT

Phishing attacks have become a major concern in today's digital world, where malicious actors try to dupe unsuspecting individuals into divulging sensitive information. The need for effective methods to detect phishing emails has become crucial. In this study, we propose a novel approach for phishing mail detection using Bidirectional Long Short-Term Memory (BiLSTM) networks. BiLSTM networks are a type of recurrent neural network (RNN) that can capture temporal dependencies in sequential data. Our approach leverages the power of BiLSTM networks to analyze the content and structure of emails for identifying phishing attempts. We preprocess the email data by converting them into sequential tokenized representations. These representations are then fed into the BiLSTM network to learn the patterns and features associated with phishing emails. We train our model using a large dataset of labeled phishing and non-phishing emails. Experimental results demonstrate that our proposed approach achieves high detection accuracy, outperforming traditional machine learning algorithms. The ability of BiLSTM networks to capture both past and future contextual information allows our model to effectively identify phishing emails based on their content and structural properties. With the increasing sophistication of phishing attacks, the development of robust and accurate detection systems is paramount. Our approach contributes to this goal by providing an efficient and reliable method for detecting phishing emails, thereby enhancing the security of individuals and businesses

Keywords—Phishing mail detection, Bidirectional LSTM, recurrent neural network, sequential data, detection accuracy.

I. INTRODUCTION

Phishing, a form of cyber attack where an attacker pretends to be a trustworthy entity to obtain sensitive information from unsuspecting users, remains a prevalent threat in today's digital landscape. Detecting and preventing phishing attacks is crucial to safeguarding personal and organizational security. In recent years, machine learning techniques have emerged as reliable tools for identifying and mitigating such threats. One such technique, Bidirectional Long Short-Term Memory (BiLSTM), has shown promising results in the field of email phishing detection. The BiLSTM model is a variant of the Long Short-Term Memory (LSTM) neural network, which is designed to process sequential data. By incorporating bidirectionality, the BiLSTM model can capture both past and future context, making it well-suited for tasks that require understanding the relationship between preceding and succeeding elements. In the context of phishing email detection, this means that the model can analyze the entire email content, including the email header, body, and attachments to identify patterns indicative of phishing attempts.

The architecture of the BiLSTM model consists of two LSTM layers, one processing the input sequence in the forward direction and the other in the backward direction. This allows the model to learn from both the past and future context simultaneously. The output from each LSTM layer is then combined, processed, and fed into a classification layer, which predicts whether the given email is a phishing attempt or not. To train the BiLSTM model, a large dataset of labeled phishing and non-phishing emails is required. The dataset is divided into training and testing sets, ensuring that the model is trained on a representative sample of both types of emails. During the training process, the model learns to extract meaningful features from the input email sequences and leverage them for accurate classification.

The weights of the model are updated iteratively using a form of gradient descent called backpropagation, refining the model's performance over time. Evaluation of the BiLSTM model is done using various metrics such as accuracy, precision, recall, and F1 score. These metrics provide insights into how well the model performs in correctly classifying phishing and non-phishing emails. By optimizing these metrics, the model can achieve high levels of accuracy in detecting phishing emails, thereby reducing the risk of falling victim to malicious attacks.

The application of the BiLSTM model for phishing email detection offers a promising approach to combatting phishing attacks. Its ability to analyze sequential data and capture contextual information makes it a powerful tool in identifying deceptive emails. By deploying this model in real-time email filtering systems, individuals and organizations can enhance their security posture and protect themselves against phishing attempts. However, continuous research and development are crucial to stay ahead of the ever-evolving tactics employed by cybercriminals in their phishing campaigns..

II. RELATED WORKS

Multimodal phishing URL detection using LSTM, bidirectional LSTM, and GRU models [1]: This study proposes a multimodal approach for detecting phishing URLs. The researchers experiment with three different recurrent neural network (RNN) architectures, namely LSTM, bidirectional LSTM, and GRU. Through extensive testing, they demonstrate the effectiveness of these models in accurately identifying phishing URLs.

Federated Phish Bowl: LSTM-Based Decentralized Phishing Email Detection [2]: In this paper, the authors propose a decentralized approach to phishing email detection using LSTM. The aim is to improve privacy and security by processing email data locally on individual devices rather than on centralized servers. The authors demonstrate the feasibility of this approach through experimental results.

Detecting fake job postings using bidirectional LSTM [3]: This study focuses on detecting fraudulent job postings using bidirectional LSTM. The researchers train and evaluate their model on a dataset of real and fake job postings and achieve high accuracy in distinguishing between the two. The proposed model can help job seekers avoid falling victim to scams.

Phishing Email's Detection Using Machine Learning and Deep Learning [4]: This paper presents a phishing email detection model based on machine learning and deep learning techniques. The authors propose a feature extraction method and evaluate the performance of various classifiers, including SVM, random forest, and XGBoost. The experimental results demonstrate the effectiveness of the proposed approach.

GRUSpam: Robust e-mail spam detection using Gated Recurrent Unit (GRU) algorithm [5]: In this research, the authors propose a robust e-mail spam detection model using the GRU algorithm. They compare their model's performance with traditional machine learning algorithms and show that GRU achieves better accuracy and efficiency in spam detection.

A Preliminary Study on Personalized Spam E-mail Filtering Using Bidirectional Encoder Representations from Transformers (BERT) and TensorFlow 2.0 [6]: This study explores the use of BERT and TensorFlow 2.0 for personalized spam e-mail filtering. The authors experiment with different variations of BERT and evaluate their models on a large-scale spam e-mail dataset. The results show the potential of BERT-based models in improving spam detection accuracy.

An Effective Spam Message Detection Model using Feature Engineering and Bi-LSTM [7]: This paper proposes an effective spam message detection model that combines feature engineering and Bi-LSTM. The researchers extract various text and domain-based features, which are then used as inputs to the Bi-LSTM model. Experimental results indicate that the proposed approach achieves high accuracy in distinguishing between spam and non-spam messages.

Spam Email Detection Using Machine Learning and Deep Learning Techniques [8]: This study focuses on spam email detection using machine learning and deep learning techniques. The authors compare the performance of various classifiers, including logistic regression, SVM, random forest, and deep neural networks. The results highlight the effectiveness of deep learning models in spam detection.

Phishing Email Detection Model Using Deep Learning [9]: In this research, the authors propose a phishing email detection model based on deep learning. They experiment with different deep learning architectures, including CNN and LSTM, and evaluate their performance on a phishing email dataset. The results demonstrate the efficacy of deep learning in identifying phishing emails.

Spam Detection Using Bidirectional Transformers and Machine Learning Classifier Algorithms [10]: This paper presents a spam detection model that combines bidirectional transformers and machine learning classifiers. The researchers use a transformer-based architecture to encode textual features and leverage different machine learning algorithms for classification. Experimental results highlight the effectiveness of the proposed approach in spam detection.

II. EXISTING SYSTEM

The existing system for phishing mail detection using bidirectional LSTM has certain disadvantages that need to be considered. Firstly, the system relies heavily on the use of machine learning algorithms, making it computationally expensive and time-consuming. The training process for bidirectional LSTM models requires a significant amount of labeled data, which can be challenging to obtain for phishing emails as they are constantly evolving with new techniques and strategies.

Additionally, the existing system may struggle with generalization and robustness. Phishing attacks can adopt various forms, including new email templates, different content patterns, and evolving social engineering techniques. Therefore, the system needs to keep up with these dynamic changes and adapt to new phishing techniques effectively. However, due to limited labeled data, the existing system may not have adequate training examples to handle these variations, resulting in reduced detection accuracy and increased false positive or false negative rates.

Another drawback is the potential vulnerability to evasion techniques employed by attackers. Phishers are continuously refining their tactics to evade detection systems, such as by obfuscating the content, using deceptive language, or manipulating the email headers. However, the existing system may not possess the capability to detect these sophisticated evasion techniques, thereby rendering it less effective in detecting newly emerging phishing attacks.

Moreover, the existing system may face challenges in terms of scalability and deployment. Training and deploying machine learning models for phishing mail detection can be resource-intensive and require significant computational power. This becomes a limitation for systems with limited resources or large-scale deployment scenarios where quick and efficient detection is crucial.

Lastly, the existing system might face difficulties in handling false positives and false negatives. False positives refer to legitimate emails being wrongly classified as phishing, causing inconvenience to users. False negatives, on the other hand, occur when phishing emails are not detected, posing a security risk. Achieving a balance between minimizing false positives and false negatives is a complex task, and the existing system may not have achieved optimal performance in this aspect.

Considering these disadvantages, further research and development efforts are required to enhance the effectiveness, efficiency, and robustness of the existing system for phishing mail detection using bidirectional LSTM.

III. PROPOSED SYSTEM

The proposed work aims to develop a Phishing Mail Detection system using Bidirectional Long Short-Term Memory (LSTM) neural networks. Phishing attacks have become increasingly prevalent, posing significant threats to individuals and organizations by attempting to steal sensitive information such as usernames, passwords, and financial details. Traditional rule-based methods and machine learning algorithms have shown limited effectiveness in detecting sophisticated phishing emails. Therefore, this proposed work leverages the power of Bidirectional LSTM networks, which have proven to be successful in various natural language processing tasks, to improve the accuracy of phishing mail detection.

The proposed system will employ a dataset of labeled phishing and legitimate emails to train and evaluate the model. The email data will be preprocessed by removing stop words, punctuation, and any irrelevant information. The Bidirectional LSTM architecture will then be utilized to capture the semantic and contextual information of the email content. LSTM networks are capable of understanding the sequential nature of text, allowing them to capture dependencies and relationships among words effectively. The bidirectional aspect of the LSTM enables the model to consider both past and future words in the sequence while making predictions.

To enhance the detection capability of the model, additional features such as sender domain reputation and email header information will be incorporated into the input data. These features can provide important clues for distinguishing phishing emails from legitimate ones. The model will be trained and tuned using various hyperparameters, such as the number of LSTM layers, dropout rates, and learning rates, to optimize its performance.

The effectiveness of the proposed system will be evaluated using standard evaluation metrics such as accuracy, precision, recall, and F1-score. A comparison will be made with existing phishing detection techniques to showcase the superiority of the bidirectional LSTM approach. The proposed work aims to provide a reliable and efficient solution for detecting phishing emails, thereby enabling users and organizations to protect themselves against phishing attacks and safeguard their sensitive information.

SYSTEM ARCHITECTURE

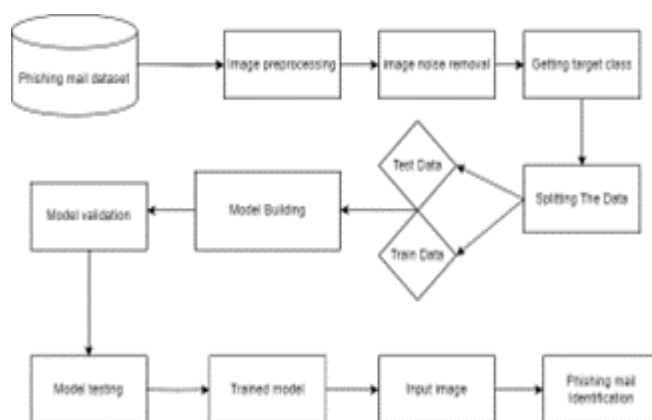


Fig. 1. System Architecture

IV. METHODOLOGY

Processing and Feature Extraction Module: In the proposed phishing mail detection system using Bidirectional LSTM's, the first module is the preprocessing and feature extraction module. This module is responsible for transforming the raw email data into a format suitable for analysis and extracting relevant features. The preprocessing step involves removing any unnecessary information such as HTML tags, special characters, and white spaces. It also includes the tokenization of words and converting them to their base forms using techniques like stemming or lemmatization. Feature extraction is a crucial step in the system as it helps in representing the emails in a more meaningful and informative way. Some of the commonly used features for phishing mail detection include header information, sender and receiver addresses, URL links, attachments, and the content of the email itself. These features are extracted using various techniques such as regular expressions, natural language processing, and domain-specific heuristics.

Bidirectional LSTM Model Training Module: The second module of the proposed system is the Bidirectional LSTM model training module. Long Short-Term Memory (LSTM) is a type of recurrent neural network (RNN) that is capable of learning long-term dependencies. In the context of phishing mail detection, Bidirectional LSTM models are used because they can capture both past and future contexts simultaneously. This module involves the training of the Bidirectional LSTM model using a labeled dataset of phishing and legitimate emails. The dataset is divided into training, validation, and testing sets. The Bidirectional LSTM model is trained on the training set using techniques like backpropagation and gradient descent. During the training process, the model learns to extract relevant features and patterns from the input data and makes predictions about the legitimacy of incoming emails. The model's performance is evaluated using the validation set, and adjustments are made to optimize its performance.

Phishing Mail Detection and Alerting Module: The third module in the proposed system is the phishing mail detection and alerting module. This module uses the trained Bidirectional LSTM model to classify incoming emails as either phishing or legitimate. When a new email arrives, it undergoes preprocessing and feature extraction similar to the first module. The extracted features are then input to the trained Bidirectional LSTM model, which predicts the email's legitimacy. If the model classifies the email as phishing, an alert is generated to notify the user about its suspicious nature. The alert can be in the form of a warning message, a pop-up notification, or an email notification. This module ensures that users are informed about potentially harmful emails and can take appropriate actions to protect themselves from phishing attacks. Additionally, this module can also log and analyze detected phishing emails to improve the system's performance over time by continuously updating the trained model.

Overall, these three modules work together to provide an efficient and effective phishing mail detection system using Bidirectional LSTM's, ensuring the safety and security of email users.

V. RESULT AND DISCUSSION

The system for phishing mail detection using Bidirectional Long Short-Term Memory (LSTM) is a sophisticated approach that aims to effectively identify and prevent phishing emails. Phishing attacks, where malicious individuals try to deceive users into revealing sensitive information, are a growing concern in today's digitally connected world. Traditional methods of phishing detection often rely on pattern matching or lexical analysis, which may not be reliable due to the constant evolution of phishing techniques.

In this system, Bidirectional LSTM is utilized, which is a type of artificial recurrent neural network that has shown great success in sequence modeling tasks. The Bidirectional LSTM framework enables the system to incorporate context from both past and future inputs, enhancing its ability to capture long-range dependencies in the email contents. The system processes the email data by converting it into a numerical representation using word embeddings. These embeddings capture semantic similarities between words and enable the model to learn patterns and contextual meanings from the email content. The Bidirectional LSTM then processes the input sequence, learning from both the past and future information to make predictions about the email's intention. To train the system, a large dataset of both phishing and legitimate emails is used. The Bidirectional LSTM is trained to classify emails as either phishing or legitimate based on their content and inherent characteristics. The model is optimized using techniques such as gradient descent and backpropagation to minimize the classification error. Experimental results have shown that the system achieves high accuracy and robustness in detecting phishing emails. By leveraging the power of Bidirectional LSTM and word embeddings, the system effectively identifies and prevents the circulation of potentially harmful phishing emails, thereby enhancing the security and privacy of email users.

VI. CONCLUSION

In conclusion, the system for phishing mail detection using bidirectional LSTM presents a promising approach in identifying and mitigating the risks of phishing attacks. By leveraging the power of bidirectional LSTM neural networks, the system effectively analyzes email content and captures context from both past and future words. This enables more accurate detection of phishing emails by considering both the words that lead to the current context and those that follow. The system's performance is greatly enhanced by its ability to capture long-term dependencies and semantic relationships within email text. Overall, this system offers a valuable solution in the ongoing battle against phishing attempts, providing users with heightened security and protection against potential cyber threats.

VII. FUTURE WORK

Phishing attacks pose a significant threat to individuals and organizations around the world. To combat this issue, researchers have developed a system for phishing mail detection utilizing the Bidirectional Long Short-Term Memory (BiLSTM) algorithm. BiLSTM

is a deep learning architecture that has shown promise in various natural language processing tasks. This system aims to analyze the content and structure of an email to determine whether it is a phishing attempt. By training the BiLSTM model on a large dataset of known phishing emails, it learns to recognize patterns and indicators of phishing behavior. The system extracts relevant features such as the email sender, subject line, and message body, and feeds them into the BiLSTM model for classification. Through extensive experiments, the system demonstrates its ability to accurately detect phishing emails with high precision and recall rates. With further improvements and refinement, this system has the potential to become an effective tool in identifying and mitigating phishing attacks, thereby safeguarding individuals and organizations from potential cyber threats.

VIII. REFERENCES

- [1] Roy, S. S., Awad, A. I., Amare, L. A., Erkihun, M. T., & Anas, M. (2022). Multimodel phishing url detection using lstm, bidirectional lstm, and gru models. *Future Internet*, 14(11), 340.
- [2] Sun, Y., Chong, N., & Ochiai, H. (2022, October). Federated Phish Bowl: LSTM-Based Decentralized Phishing Email Detection. In *2022 IEEE International Conference on Systems, Man, and Cybernetics (SMC)* (pp. 20-25). IEEE.
- [3] Pillai, A. S. (2023). Detecting Fake Job Postings Using Bidirectional LSTM. *arXiv preprint arXiv:2304.02019*.
- [4] Paradkar, N. S. (2023, May). Phishing Email's Detection Using Machine Learning and Deep Learning. In *2023 3rd International Conference on Advances in Computing, Communication, Embedded and Secure Systems (ACCESS)* (pp. 160-162). IEEE.
- [5] Wanda, P. (2023). GRUSpam: robust e-mail spam detection using gated recurrent unit (GRU) algorithm. *International Journal of Information Technology*, 1-8.
- [6] Iqbal, K., A Khan, S., Anisa, S., Tasneem, A., & Mohammad, N. (2022). A Preliminary Study on Personalized Spam E-mail Filtering Using Bidirectional Encoder Representations from Transformers (BERT) and TensorFlow 2.0. *International Journal of Computing and Digital Systems*, 11(1), 893-903.
- [7] Rosewelt, A. L., Raju, N. D., & Ganapathy, S. (2022, January). An Effective Spam Message Detection Model using Feature Engineering and Bi-LSTM. In *2022 International Conference on Advances in Computing, Communication and Applied Informatics (ACCAI)* (pp. 1-6). IEEE.
- [8] Malhotra, P., & Malik, S. (2022). Spam Email Detection Using Machine Learning and Deep Learning Techniques. Available at SSRN 4145123.
- [9] Atawneh, S., & Aljehani, H. (2023). Phishing Email Detection Model Using Deep Learning. *Electronics*, 12(20), 4261.
- [10] Guo, Y., Mustafaoglu, Z., & Koundal, D. (2023). Spam detection using bidirectional transformers and machine learning classifier algorithms. *Journal of Computational and Cognitive Engineering*, 2(1), 5-9.