

Paper Title:

Detection and Analyzing Phishing Emails Using NLP Techniques

Paper Link:

<https://ieeexplore.ieee.org/document/10156738>

1 Summary**1.1 Motivation**

The motivation of this paper is to address the increasing threat of phishing attacks, particularly in the Arabic language, and the lack of research and solutions available in this area. The paper aims to propose a new model that utilizes natural language processing (NLP) techniques to extract and analyze the content of Arabic email messages. The goal is to differentiate between legitimate emails and phishing emails by using three determinants. By combining multiple methods, the paper demonstrates the effectiveness of using a comprehensive approach to achieve successful results in phishing email detection.

1.2 Contribution

The contribution of this paper is to propose a new model for detecting phishing emails in Arabic by extracting and analyzing their content. The three determinants used are a blacklist of commonly used Arabic phishing words, the roots of those words, and a list of commonly used Arabic phishing sentences. The paper presents the results with accuracy rates of 99% for identifying legitimate emails and 96% for identifying phishing emails. The paper also explores anti-phishing solutions, including phishing blacklists and NLP and ML-based detection of phishing emails.

1.3 Methodology

The methodology of this paper involves proposing a new model to extract Arabic email content. The process starts by emailing a specific account with all the acquired electronic messages. The messages are then checked to see if they are in Arabic and are analyzed based on their content using three separate methods:

- Comparing the phrases in the email content with a blacklist of common Arabic phishing words
- Obtaining the root of each word in the email content and comparing it to a blacklist of Arabic common phishing words.
- Comparing the phrases in the email text with a list of common Arabic phishing sentences The paper also explores various anti-phishing solutions, including phishing blacklists and NLP and ML-based detection of phishing emails.

1.4 Conclusion

In conclusion, phishing is a pervasive cyber attack that exploits the gullibility of internet users to deceive them into revealing personal information or downloading malicious software. The determinants in the document include a blacklist of common phishing words, the roots of these words, and a list of phishing sentences. The results obtained using these determinants where some show high accuracy in identifying legitimate emails and phishing emails. This research contributes to addressing the scarcity of Arabic-language-based phishing email detection and provides insights for future studies in content-based phishing email detection.

2 Limitations

2.1 First Limitation

One of the limitations of this paper is that using only one method/determinant gives very low accuracy. In one experiment, the accuracy for recognizing legitimate emails was 98.5 percent, while the accuracy for phishing emails was only 42.0 percent when using the root of words in the blacklist as a determinant. The document does not provide a specific explanation for the low accuracy of using the root of words in the blacklist. However, based on the results mentioned, it can be inferred that the root of words alone might not be sufficient to accurately detect phishing emails.

2.2 Second Limitation

The document does not provide detailed information about the size and diversity of the dataset used for training and testing the model, which could potentially impact the generalizability and reliability of the results.

3 Synthesis

This paper proposes a new model for detecting phishing emails in Arabic by utilizing neural language programming (NLP) techniques. The model extracts the Arabic email content and analyzes it using three determinants: a blacklist of Arabic common phishing words, the roots of a blacklist of Arabic common phishing words, and a list of Arabic common phishing sentences. At first, the researchers provided individual methods to work with the analysis of detecting phishing and legal emails. This shows very low accuracy, which led to the combination of two different methods to work with. Furthermore, they proposed a complex hybrid solution using three conditions (list of words, roots, and sentences) to get a higher accuracy of 99% for identifying legitimate emails and 96% for identifying phishing emails. The document also explores other anti-phishing solutions, such as phishing blacklists, NLP, and machine learning-based detection methods. The paper also explains the future work and improvements in detecting phishing emails.