

An Optimized Ensemble Learning Framework for Phishing and Malicious URL Detection

Mrs. J. Hima Bindu¹, Kota Reethi Chandrika², Paidi Ashritha³

¹Assistant Professor, Mahatma Gandhi Institute of Technology

^{2,3}UG Student, Mahatma Gandhi Institute of Technology

Abstract - Malicious URL attacks frequently utilize counterfeit websites and legitimate logos to mislead users, presenting a considerable risk in the digital landscape. These attacks can result in the exposure of sensitive information, such as banking credentials and passwords. Although there are existing anti-phishing strategies, cybercriminals continually adapt their methods, underscoring the necessity for a reliable prediction mechanism to safeguard users effectively. Classification serves as one of the techniques for identifying phishing websites. This paper introduces a model designed to detect phishing attacks through the application of various Machine Learning (ML) classifiers. By analysing the "phishing_legitimate_full.csv" dataset sourced from Kaggle, our proposed system incorporates ML models, including Random Forest, CATB, and Naive Bayes, to identify malicious URLs. The dataset comprises 48 features derived from phishing sites and legitimate sites. The methodology encompasses dataset partitioning, model training, prediction, and evaluation to improve accuracy and minimize latency, employing effective feature selection methods to extract relevant features from raw URLs. Ultimately, our system aspires to advance the development of automated solutions capable of effectively classifying and mitigating phishing threats, thereby enhancing user security in the online environment.

Keywords: Malicious URL detection, phishing attacks, machine learning, classification, Random Forest, CATB, Naive Bayes, anti-phishing, cyber security, feature selection, URL analysis, automated threat detection, digital security.

I. INTRODUCTION

The swift growth of online services has resulted in a notable increase in cyber threats, particularly phishing attacks that manipulate users' trust to acquire sensitive information. Malicious URLs are a primary means of these attacks, as they imitate legitimate websites to trick users into revealing their credentials, financial information, and other confidential data. Although traditional security measures, such as blacklisting and heuristic detection, have been implemented, phishing tactics continue to advance, rendering these conventional methods less effective. In response to these challenges, machine learning-based strategies have gained traction due to their capacity to analyze URL patterns and classify them as either safe or harmful. Among these strategies, ensemble

learning techniques, which integrate multiple models, have shown improved detection accuracy by minimizing biases and capitalizing on the strengths of individual classifiers. This study investigates the efficacy of ensemble models, including Random Forest, CatBoost, and Naïve Bayes, in detecting phishing attempts through a systematic evaluation of their classification performance. The methodology encompasses data preprocessing, feature selection, model training, and evaluation of predictive capabilities based on essential performance metrics. The comparative analysis of these models seeks to determine the most effective and accurate method for identifying phishing URLs, taking into account factors such as scalability, real-time applicability, and computational efficiency. By utilizing machine learning in this domain, this research aims to enhance cybersecurity solutions, improving defences against evolving online threats and ensuring safer digital interactions for users.

A. Problem Statement.

Phishing attacks have emerged as one of the most significant threats in the realm of cybersecurity, utilizing social engineering tactics to trick individuals into disclosing sensitive information, including login credentials, financial data, and personal details. Traditional methods for detecting phishing, such as blacklisting and rule-based systems, often fall short in adapting to the ever-changing strategies employed by cybercriminals. These conventional approaches frequently experience high rates of false positives, delayed identification of threats, and limited flexibility in responding to new phishing schemes. With the increasing complexity of phishing websites and the heightened dependence on online transactions, there is an urgent requirement for a sophisticated detection system capable of accurately and promptly identifying malicious URLs in real-time. This research intends to tackle this issue by employing machine learning-based ensemble methods, specifically Random Forest, CatBoost, and Naïve Bayes, to improve the precision and resilience of phishing detection systems. Through a comprehensive analysis and comparison of these models, the study aims to create an optimized strategy that effectively reduces phishing risks while ensuring scalability and practical applicability in rapidly changing online environments.

B. Existing System

Current phishing detection systems predominantly utilize conventional machine learning algorithms, including K-Nearest Neighbors (KNN), Support Vector Machines (SVM), Random Forest, and Logistic Regression. These algorithms are generally

trained on datasets characterized by a static set of features derived from URLs, such as domain age, the presence of suspicious keywords, and URL length. Although Random Forest has demonstrated commendable performance in terms of precision, recall, and F1-score, other models like SVM and Logistic Regression have exhibited inconsistent results, particularly when confronted with a variety of evolving phishing strategies. Furthermore, the exploration of ensemble methods that integrate classifiers—such as Random Forest and CatBoost—has been undertaken; however, these approaches incur significant computational expenses, which restrict their scalability for real-time applications. Additionally, these systems frequently encounter challenges related to generalizability, often struggling to maintain accuracy on datasets that differ from those utilized during training. Problems such as potential dataset bias, susceptibility to adversarial attacks, and difficulties in managing imbalanced data further compromise the robustness of these systems, underscoring the necessity for more flexible and efficient solutions.

Disadvantages of Existing System:

- *Limited Adaptability to Evolving Phishing Tactics*
- *Dataset Bias and Generalizability Issues*
- *High Computational Cost and Scalability Concerns*
- *Challenges in Handling Imbalanced Data*
- *Limited Feature Set and Feature Extraction Limitations*

PROPOSED SYSTEM

A. Architecture of Proposed System.

The proposed system adopts a systematic methodology to effectively identify phishing URLs through the application of machine learning techniques. It initiates with the collection and preprocessing of data, wherein labeled URLs are assembled and refined to address missing values, standardize features, and extract pertinent attributes. Techniques such as Recursive Feature Elimination (RFE) are employed for feature selection and engineering, enabling the identification of essential attributes, including domain-based, lexical, and content-based features, which contribute to enhanced classification accuracy. Subsequently, the system trains three distinct machine learning models—Random Forest, CatBoost, and Naïve Bayes—on the preprocessed dataset, with each model learning to distinguish between legitimate and phishing URLs. An ensemble learning strategy is implemented to amalgamate predictions, capitalizing on the strengths of individual classifiers to bolster overall detection accuracy. Upon completion of the training phase, the system classifies new, previously unseen URLs by assigning a probability score that indicates whether a URL is

safe or malicious. The performance of the models is assessed using metrics such as precision, recall, F1-score, and accuracy, facilitating a comparative analysis to identify the most effective model or combination of models. Ultimately, the optimized model is deployed within a real-time phishing detection framework, ensuring both scalability and adaptability to the continuously evolving landscape of phishing threats.

Advantages of Proposed System.

- *Enhanced Detection Accuracy*
- *Real-time Threat Analysis*
- *Adaptive Learning for Evolving Threats*
- *Scalability for Large Networks*
- *Improved Efficiency and Speed*

II. LITERATURE SURVEY

This research examines the application of machine learning methodologies for the identification of phishing websites, with an emphasis on the performance and efficiency of various models. The authors evaluate a range of classification techniques, including Random Forest, CatBoost, and Naïve Bayes, utilizing a dataset comprising 48 extracted features. The findings underscore the benefits of ensemble methods in enhancing accuracy and minimizing false positive rates. Nonetheless, issues such as computational complexity, optimization of feature selection, and the need for real-time adaptability persist as significant challenges.[1]

This article introduces an ensemble learning strategy that merges Random Forest and Gradient Boosting for the purpose of phishing detection. The investigation reveals that the combination of multiple classifiers leads to improved detection accuracy. The authors observe that while ensemble techniques bolster classification performance, they also impose substantial computational requirements, complicating real-time detection efforts. There is a call for further research aimed at enhancing efficiency and scalability.[2]

A framework for phishing detection based on deep learning is discussed in this paper, which utilizes neural networks to scrutinize URL structures and webpage content. The model demonstrates superior accuracy compared to conventional machine learning approaches. However, it necessitates extensive datasets and considerable computational power, and it is susceptible to adversarial attacks. The study advocates for additional research into lightweight deep learning models to improve their applicability in real-world scenarios.[3]

This review assesses the efficacy of support vector machines (SVM) in the detection of phishing URLs. The findings suggest that SVM is effective in differentiating phishing sites from legitimate ones through feature-based analysis. However, the research recognizes the dependency on specific datasets and the necessity for enhanced generalization across various phishing attack methodologies. The paper proposes hybrid strategies that

integrate SVM with other classifiers to achieve greater robustness.[4]

A comparative analysis of various machine learning models for phishing detection is conducted, focusing on the effectiveness of Random Forest, Naïve Bayes, and Logistic Regression. The findings indicate that Random Forest consistently outperforms the other models in terms of accuracy, precision, and recall. However, the study also highlights ongoing challenges, including the management of imbalanced datasets and the optimization of feature engineering. This research emphasizes the necessity for adaptive machine learning strategies to address the continuously evolving tactics employed in phishing attacks. [5].

CONCLUSION

This research provides an in-depth examination of machine learning methodologies employed for phishing detection, with particular emphasis on the efficacy of ensemble techniques including Random Forest, CatBoost, and Naïve Bayes. By leveraging a dataset that encompasses a variety of URL characteristics, the proposed framework significantly improves the accuracy of phishing detection while minimizing both false positives and false negatives. The results indicate that CatBoost surpasses other classifiers, attributed to its superior management of categorical data and effective feature selection processes.

The study emphasizes the importance of combining multiple machine learning models to enhance the resilience of phishing detection systems. Nonetheless, issues such as computational demands, the need for real-time responsiveness, and potential biases within datasets present ongoing challenges that warrant further investigation. Future research should prioritize the optimization of feature engineering, the reduction of model latency, and the integration of deep learning approaches to fortify phishing detection strategies. Through the advancement of these techniques, this study aspires to contribute to the evolution of more secure and efficient cybersecurity frameworks.

REFERENCES

[1] F. Ali, M. Unnisa Begum, "A Machine Learning Approach to Detect Malicious URLs for Phishing Attack Prevention," in *Journal of Science Technology and Education*, vol. 12, no. 2, pp. 336–345, 2023.

[2] H. Liu, Y. Zhang, X. Wang, "Single and Hybrid-Ensemble Learning-Based Phishing Website Detection," in *IEEE*

Transactions on Information Forensics and Security, vol. 18, no. 3, pp. 678–685, 2023, doi: 10.1109.

[3] J. Wang, X. Li, Z. Chen, "Phishing Detection Using Neural Networks: A Comparative Analysis," in *Sensors*, vol. 23, no. 7, p. 3467, 2023, doi: 10.3390/s23073467.

[4] S. Zhang, L. Huang, Q. Feng, "Deep Learning for Malicious URL Detection: Challenges and Opportunities," in *IEEE Access*, vol. 10, pp. 36505-36517, 2022, doi: 10.1109.

[5] A. Ahmed, P. Singh, R. Verma, "Comparative Study of Phishing Detection Algorithms Using Machine Learning Techniques," in *Journal of Ambient Intelligence and Humanized Computing*, vol. 13, pp. 1–14, 2022, doi: 10.1007.