ABSTRACT

Phishing attacks have become a significant and growing cybersecurity threat. To address this challenge, this study proposes the development of a phishing detection model based on an ensemble learning approach using a Soft Voting Classifier. The research aims to enhance accuracy and reduce the false positive rate by leveraging the predictive power of multiple classification models. The methodology involved collecting 10,000 phishing URLs from PhishTank.org and 5,000 legitimate URLs from Cisco Umbrella, which were then extracted into 15 key features. Single classification models such as Logistic Regression, K-Nearest Neighbor, Random Forest, and Naive Bayes were implemented and trained, along with various Soft Voting combinations of these models. The evaluation results show that the Random Forest model achieved the highest accuracy among single models at 87.60%. Meanwhile, the Soft Voting combination of Logistic Regression, K-Nearest Neighbor, and Random Forest demonstrated a very competitive performance with an accuracy of 87.47%, the highest Recall of 0.9175. This study concludes that while the Soft Voting Classifier did not significantly improve performance compared to the best single model, it successfully achieved comparable results and provided a stable performance profile. The resulting balance between precision and recall makes it a reliable and relevant solution for mitigating complex phishing threats.

Keywords: Phishing Detection, Machine Learning, Soft Voting Classifier, Ensemble Learning, Cybersecurity