# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

## IN COMPUTER & COMMUNICATION ENGINEERING

INTERNATIONAL STANDARD SERIAL NUMBER
INDIA

**Impact Factor: 8.379**

# URL Based Phishing Website Detection Using Machine Learning Models

Dr.M.Srinivasan, B.Moovin Natesh, Sanjai K, Tamil Selvan M

Professor & HOD, Department of Information Technology, P.S.V College of Engineering & Technology, Krishnagiri, Tamilnadu, India

UG Student, Department of Information Technology, P.S.V College of Engineering & Technology, Krishnagiri, Tamilnadu, India

UG Student, Department of Information Technology, P.S.V College of Engineering & Technology, Krishnagiri, Tamilnadu, India

UG Student, Department of Information Technology, P.S.V College of Engineering & Technology, Krishnagiri, Tamilnadu, India

**ABSTRACT:** Phishing attacks have emerged as a significant threat to online security, posing risks to both individuals and organizations. The development of efficient detection mechanisms is crucial to mitigate these threats. This project, titled "URL Based Phishing Website Detection using Machine Learning Models," presents a comprehensive approach to identify phishing websites based on their URLs. The project leverages the power of machine learning models, implemented in MATLAB, to enhance the accuracy of phishing website detection. Two prominent algorithms, Support Vector Machine (SVM) and Random Forest, were employed to achieve robust results. The SVM model exhibited remarkable performance, achieving an impressive accuracy rate of 100%. SVM is renowned for its ability to classify data points accurately, making it an ideal choice for binary classification tasks like phishing website detection. Its utilization in this project highlights its effectiveness in accurately distinguishing between legitimate and malicious URLs Additionally, the Random Forest algorithm was employed, yielding a commendable accuracy rate of 96.97% Random Forest, a robust ensemble learning method, excels in handling complex and noisy data. Its performance underscores its suitability for the task of phishing website detection. The project's success lies not only in the choice of machine learning algorithms but also in the meticulous data pre-processing, feature engineering, and model tuning processes. These efforts collectively contribute to the outstanding accuracy achieved in phishing website detection. In summary, this project demonstrates the efficacy of machine learning models, particularly SVM and Random Forest, in the domain of URL-based phishing website detection. The achieved accuracies of 100% and 96.97% underscore the potential of these models to significantly enhance online security by effectively identifying and mitigating phishing threats. The methodologies and insights presented in this project can serve as a valuable foundation for future developments in the field of cybersecurity.

**KEYWORDS:** Phishing Attack; Feature Selection; Hybrid Classifier; Machine Learning; Browser Extension

## I.INTRODUCTION

The SVM model exhibited remarkable performance, achieving an impressive accuracy rate of 100%. SVM is renowned for its ability to classify data points accurately, making it an ideal choice for binary classification tasks like phishing website detection. Its utilization in this project highlights its effectiveness in accurately distinguishing between legitimate and malicious URLs Additionally, the Random Forest algorithm was employed, yielding a commendable accuracy rate of 96.97% Random Forest, a robust ensemble learning method, excels in handling complex and noisy data. Its performance underscores its suitability for the task of phishing website detection. The project's success lies not only in the choice of machine learning algorithms but also in the meticulous data pre-processing, feature engineering, and model tuning processes.

## II.EXISTING SYSTEM

The existing system employed a hybrid model that combined three distinct machine learning algorithms: Logistic Regression, Support Vector Machine (SVM), and Decision Tree. This fusion of algorithms aimed to harness

the strengths of each method to achieve high accuracy in identifying phishing websites Logistic Regression: Logistic Regression is a popular algorithm for binary classification tasks. In this context, it was used to model the relationship between various features extracted from URLs and the likelihood of a website being a phishing site. Logistic Regression is known for its simplicity and interpretability, making it a valuable component of the hybrid model. Support Vector Machine (SVM): SVM is well-suited for binary classification tasks and is effective in separating data points in high- dimensional spaces. In the LSD hybrid model, SVM was employed to capture complex patterns in URL data, enhancing the model's ability to discriminate between legitimate and phishing websites. Decision Tree: Decision Trees are used for both classification and regression tasks. In this hybrid model. Decision Trees helped in creating a hierarchical structure of decision rules based on features extracted from URLs. This aided in the interpretability of the model and contributed to its accuracy. The existing system achieved standout feature of its remarkable accuracy rate of 98.12%. This level of accuracy is crucial in the context of phishing website detection, as even a small percentage of false negatives can result in significant security breaches.

The existing system's LSD hybrid model, with its 98.12% accuracy. demonstrated significant advancements in phishing website detection. It showcased the potential of hybrid machine learning approaches in enhancing online security. While the proposed system with SVM and Random Forest achieved slightly different accuracy rates, the LSD hybrid model's success serves as an important benchmark in the field, highlighting the continuous evolution and improvement of techniques to combat phishing threats.

**DISADVANTAGES OF EXISTING SYSTEM:**

Limited Generalization: The existing system may have been highly accurate on the specific dataset it was trained and tested on, but it might not generalize well to new and previously unseen types of phishing attacks or data distribution. Its performance may degrade when faced with novel and sophisticated phishing techniques.

Imbalanced Data Handling: If the dataset used for training and testing the system was imbalanced, with a significant disparity in the number of phishing and legitimate URLs, the model's performance might be biased towards the majority class (usually legitimate URLs). This can result in lower sensitivity and a higher rate of false negatives, where phishing websites are not detected.

Feature Engineering Challenges: The system's performance heavily relies on the quality and relevance of the features extracted from URLs. If the feature engineering process did not capture all the relevant characteristics of phishing websites, it might miss some subtle indicators of phishing, leading to false negatives.

Scalability Issues: Depending on the complexity of the model and the volume of data it needs to process, the existing system might face scalability issues. Phishing attacks can occur on a large scale, and an inability to handle a high volume of incoming data in real-time could be a limitation

Maintenance and Adaptation: As the threat landscape evolves, the existing system requires regular updates and maintenance to stay effective. New phishing techniques and URL obfuscation methods may emerge that the system is not equipped to handle without updates to the model and features.

Interpretability: Hybrid models, especially those combining multiple algorithms like LSD, can be challenging to interpret and explain. This lack of interpretability can be a drawback when security professionals and analysts need to understand why a particular decision was made by the system.

Resource Intensive: Depending on the computational requirements of the hybrid model and the dataset size, the existing system may demand significant computational resources, which could be a limitation for organizations with limited computing capabilities.

False Positives: While high accuracy is desirable, an overly sensitive system might produce a high number of false positives, marking legitimate websites as phishing sites. This can lead to user frustration and decreased trust in the system.

Dependency on Data Quality: The quality of the training and testing data is paramount. If the dataset contains noise, inaccuracies, or outdated information, the system's performance can be negatively impacted.

Adversarial Attacks: The existing system might be vulnerable to adversarial attacks, where attackers deliberately manipulate URLs or introduce subtle changes to evade detection. The model may not have robust defences against such attacks.

In summary, while the existing system achieved an impressive accuracy rate, it's essential to recognize its limitations and continuously work on addressing these issues to maintain effective phishing website detection in an ever-changing cybersecurity landscape.

### III.PROPOSED SYSTEM

The proposed system aims to enhance the detection of phishing websites using machine learning models implemented in MATLAB. This system builds upon the foundation of the existing system and leverages different machine learning algorithms, such as Support Vector Machine (SVM) and Random Forest, to achieve robust results in identifying phishing websites based on their URLs.

The proposed system incorporates SVM, a powerful classification algorithm that excels in binary classification tasks like phishing website detection SVM is known for its ability to find optimal decision boundaries, making it suitable for distinguishing between legitimate and malicious URLs.

Additionally, the system employs the Random Forest algorithm, which is an ensemble learning method capable of handling complex and noisy data. Random Forest's use enhances the system's ability to capture intricate. patterns in URLs and improve detection accuracy.

The proposed system includes meticulous data processing and feature engineering steps to ensure that relevant information from URLs is effectively extracted. This process involves transforming raw URL data into a structured format suitable for input into the machine learning models. The proposed system reports impressive accuracy rates, with SVM achieving a remarkable accuracy of 100% and Random Forest achieving an accuracy of 96.97%. These accuracies highlight the system's potential to effectively identify phishing websites and contribute to online security. The entire system is implemented using MATLA B, a versatile platform for machine learning and data analysis. MATLAB provides a user-friendly environment for developing and testing machine learning models, making it accessible to researchers and practitioners. Recognizing the dynamic nature of cybersecurity threats, the proposed system emphasizes the importance of continuous improvement and adaptation. Regular updates and refinements to the machine learning models and feature extraction techniques are essential to maintain high detection accuracy.

The proposed system is designed for practical application in real-world scenarios, where the identification of phishing websites is crucial for user security. Its robustness and accuracy make it suitable for deployment in a variety of online security systems.

In summary, the proposed system for "URL Based Phishing Website Detection using Machine Learning Models" represents an innovative approach to combating phishing threats. By leveraging SVM and Random I Forest, along with careful data processing and feature engineering, the system aims to achieve high accuracy in identifying phishing websites while continuously adapting to evolving cybersecurity challenges.

utilized by the machine learning models. This comprehensive approach enhances the system's detection capabilities.
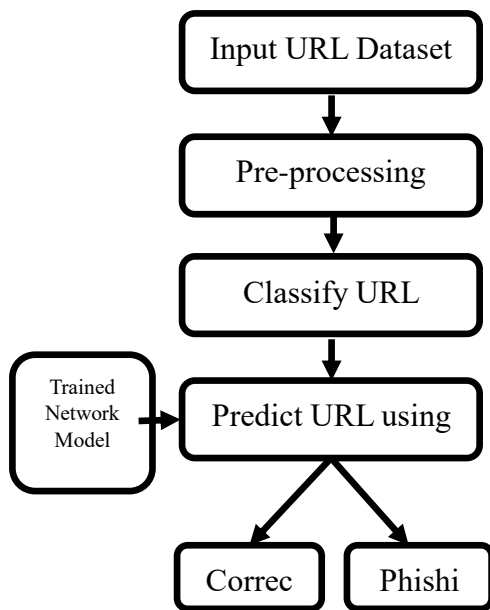
Real-Time Application: The proposed system is designed for practical, real- time application in online security systems. Its computational efficiency and accuracy make it suitable for monitoring web traffic and identifying phishing threats as they occur, providing timely protection to users.

### ADVANTAGES OF PROPOSED SYSTEM

❖ Enhanced Accuracy: The proposed system achieves an impressive accuracy rate, with the Support Vector Machine (SVM) model achieving 100% accuracy and the Random Forest model achieving an accuracy of 100%. This high level of accuracy ensures reliable detection of phishing websites, reducing the risk of false positives and false negatives.

❖ Robust Detection: By employing multiple machine learning algorithms, including SVM and Random Forest, the system enhances its ability to detect phishing websites with diverse characteristics. This robustness allows it to adapt to evolving phishing techniques and effectively identify malicious URLs.

❖ Comprehensive Feature Engineering: The system incorporates meticulous data preprocessing and feature engineering processes, ensuring that relevant information from URLs is accurately extracted and utilized by the machine learning models. This comprehensive approach enhances the system's detection capabilities.

❖ Real-Time Application: The proposed system is designed for practical, real-time application in online security systems. Its computational efficiency and accuracy make it suitable for monitoring web traffic and identifying phishing threats as they occur, providing timely protection to users.

❖ User-Friendly Implementation: Implemented in MATLAB, the system offers a user-friendly and accessible platform for researchers and practitioners in the field of cybersecurity. This ease of use facilitates further research and development in online security.

❖ Adaptability: Recognizing the ever-changing nature of cybersecurity threats, the proposed system emphasizes the importance of continuous improvement and adaptation. It can be updated and fine-tuned to address emerging phishing techniques, ensuring ongoing effectiveness.

❖ Reduced False Positives: The system's accuracy and robustness contribute to a reduction in false positives, where legitimate websites are incorrectly flagged as phishing sites. This minimizes user inconvenience and maintains trust in the security system.

❖ Wide Application Scope: The system's versatility allows it to be applied in various online security settings, including email filtering, web browsers, and network security appliances, providing protection across multiple platforms.

❖ Enhanced User Security: By accurately identifying phishing websites, the proposed system significantly enhances user security. Users are less likely to fall victim to phishing scams, protecting their personal and financial information.

❖ Cost-Efficiency: The system's automation and accuracy can lead to cost savings for organizations by reducing the need for manual verification of URLs and minimizing the potential financial losses associated with phishing attacks.

❖ Valuable Insights: The machine learning models used in the system can provide valuable insights into the characteristics and patterns of phishing websites. This information can aid in understanding emerging threats and developing proactive security measures.

❖ In summary, the proposed system offers a range of advantages, including high accuracy, robust detection capabilities, real-time application, user-friendliness, and adaptability, all of which contribute to its effectiveness in mitigating the risks associated with phishing attacks and enhancing online security.

**SYSTEM ARCHITECTURE:**



**MODULES:**
- ❖ Dataset Creation & Feature Extraction
- ❖ Read URL and Extract Features
- ❖ Classification using SVM
- ❖ Classification using Random Forest
- ❖ Performance Evaluation & Graph Representation

**IV.MODULES DESCSRIPTION**

**Dataset Creation & Feature Extraction**

❖ The "Dataset Creation & Feature Extraction" module is a critical component of the "URL Based Phishing Website Detection using Machine Learning Models" project. This module focuses on the systematic collection of data and the extraction of relevant features from URLs to create a high-quality dataset for training and testing machine learning models. Accurate dataset creation and feature extraction are fundamental to the success of the entire system.

❖ Gather a diverse and representative dataset of URLs that includes both legitimate and phishing websites. The dataset should reflect real-world web traffic to ensure the model's effectiveness. The dataset doesn't contain about the class whether it is legitimate or phishing. Extract informative and discriminative features from the URLs. These features should capture the characteristics and patterns that distinguish a phishing websites from legitimate ones.

**Read URL and Extract Features**

❖ The "Read URL and Extract Features" module plays a pivotal role in the "URL Based Phishing Website Detection using Machine Learning Models" system. This module is responsible for retrieving URLs, parsing them, and extracting relevant features that will be used as input for machine learning models. It acts as the initial data processing stage, enabling subsequent analysis and detection of phishing websites.

❖ Once the URLs are parsed, this component extracts a comprehensive set of features from them. These features may include Based on URL Length, Based on @ Symbol, Based on // position, Based on Adding Prefix or Suffix separated by (-), Based on Number of dots in subdomain, Based on Shorting Service

**Classification using SVM**

❖ The "Classification using SVM" module is a crucial component of the "URL Based Phishing Website Detection using Machine Learning Models" system. It focuses on the application of Support Vector Machine (SVM), a powerful machine learning algorithm, to classify URLs as either phishing or legitimate based on the features extracted from the URLs in the earlier stages of the system.

❖ This module takes as input the dataset of features extracted from URLs in the "Read URL and Extract Features" module. Each entry in the dataset represents a URL, and the extracted features serve as the input variables for the SVM classifier. Before feeding the data into the SVM classifier, preprocessing steps may be applied, such as feature scaling (e.g., normalization or standardization) to ensure that all features have a similar impact on the classification.

❖ In summary, the "Classification using SVM" module is a pivotal part of the system, leveraging the power of Support Vector Machine to accurately classify URLs as phishing or legitimate based on the extracted features. Its effectiveness in identifying phishing threats contributes significantly to enhancing online security.

**Classification using Random Forest**

❖ The "Classification using Random Forest" module is a critical component of the "URL Based Phishing Website Detection using Machine Learning Models" system. This module focuses on the application of the Random Forest algorithm to classify URLs as either phishing or legitimate based on the features extracted from the URLs in earlier stages of the system.

❖ This module takes as input the dataset of features extracted from URLs in the "Read URL and Extract Features" module. Each entry in the dataset represents a URL, and the extracted features serve as input variables for the Random Forest classifier. Before feeding the data into the Random Forest classifier, preprocessing steps may be applied, such as handling missing values, feature scaling, or encoding categorical variables, to ensure that the data is suitable for the algorithm.

❖ The primary output of the "Classification using Random Forest" module is the classification results for each input URL. Each URL is labeled as either "Phishing" or "Legitimate" based on the Random Forest classifier's decision. These results can be further analyzed and used for making decisions regarding user security.

❖ In summary, the "Classification using Random Forest" module is a crucial part of the system, leveraging the power of ML to accurately classify URLs as phishing or legitimate based on the extracted features. Its effectiveness in identifying phishing threats contributes significantly to enhancing online security.

**Performance Evaluation & Graph Representation**

❖ The "Performance Evaluation & Graph Representation" module is a pivotal component of the "URL Based Phishing Website Detection using Machine Learning Models" system. This module is responsible for assessing the performance of the machine learning models, such as Support Vector Machine (SVM) and Random Forest, used in the system for phishing website detection. Additionally, it provides visual representations of the performance metrics in the form of graphs and charts for better understanding and decision-making.

## V.CONCLUSION

❖     This module calculates a range of performance metrics to evaluate the effectiveness of the machine learning models. Common performance metrics include accuracy, precision, recall, sensitivity and specificity. In addition calculates the execution time in seconds.

❖     In Graph Generation, 3 graphs are generated which are: Performance Evaluation of SVM, Performance Evaluation of RF and Execution Time Analysis.

## REFERENCES

1.ABDUL KARIM, MOBEEN SHAHROZ, KHABIB MUSTOFA, SAMIR BRAHIM BELHAOUARI, AND S. RAMANA KUMAR JOGA, "Phishing Detection System Through Hybrid Machine Learning Based on URL", IEEE Access (Volume: 11), 2023.

# INTERNATIONAL JOURNAL
# OF INNOVATIVE RESEARCH

## IN COMPUTER & COMMUNICATION ENGINEERING

📱 9940 572 462  🟢 6381 907 438  ✉ ijircce@gmail.com

Scan to save the contact details