



A Robust URL Classification Approach for Phishing Detection Using Machine Learning Techniques

LAKKAVARAPU RAJKUMAR

PG Scholar. Department of MCA, DNR College, Bhimavaram, Andhra Pradesh

A.Durga Devi

(Assistant Professor), Master of Computer Applications, DNR College, Bhimavaram, Andhra Pradesh

ABSTRACT

The rapid expansion of internet usage has significantly increased cyber threats, particularly phishing attacks that deceive users into revealing sensitive information. Detecting malicious URLs has become a critical aspect of cybersecurity. This project presents an intelligent URL classification system designed to detect phishing, malware, defacement, and benign URLs using machine learning techniques. The system is implemented using the Django web framework, providing an interactive interface for both administrators and users. The proposed system leverages text-based feature extraction techniques such as Count Vectorization to transform URLs into numerical representations suitable for machine learning models. Multiple classification algorithms, including Naive Bayes, Support Vector Machine (SVM), Logistic Regression, and Stochastic Gradient Descent (SGD), are employed to evaluate performance and ensure robustness. Additionally, an ensemble approach using a Voting Classifier is implemented to improve prediction accuracy. The dataset used consists of labeled URLs categorized into four classes: benign, phishing, defacement, and malware. During training, the dataset is preprocessed and transformed, and models are trained using a train-test split strategy. Performance metrics such as accuracy, confusion matrix, and classification reports are used to evaluate the models.

The system also provides analytical features such as detection accuracy visualization, ratio analysis of URL types, and trending topics. These insights assist administrators in understanding patterns and improving security strategies. Users can input URLs and receive real-time predictions regarding their safety. Experimental results demonstrate that ensemble learning improves classification performance compared to individual models. The system achieves high accuracy and provides reliable predictions, making it suitable for real-world applications in cybersecurity.



International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

www.ijdim.com

Original Research Paper

In conclusion, the developed system effectively identifies malicious URLs and enhances online security. It can be extended further by integrating deep learning techniques and real-time threat intelligence systems.

Keywords: Phishing Detection, URL Classification, Machine Learning, Naive Bayes, Support Vector Machine, Logistic Regression, Cyber Security, Text Mining, Ensemble Learning, Django

I. INTRODUCTION

With the increasing reliance on digital platforms, cybersecurity has become a major concern. Among various cyber threats, phishing attacks are one of the most common and dangerous forms of cybercrime. These attacks aim to trick users into providing sensitive information such as passwords, credit card details, and personal data by disguising malicious websites as legitimate ones. Traditional security mechanisms such as blacklists and rule-based systems are often ineffective against newly generated phishing URLs. Attackers continuously evolve their techniques, making it difficult to detect malicious URLs using static methods. Therefore, there is a need for intelligent systems that can adapt and learn from data to identify threats accurately.

Machine learning has emerged as a powerful solution for cybersecurity challenges. By analyzing patterns in data, machine learning models can classify URLs as safe or malicious. This project focuses on developing a machine learning-based URL classification system that can detect phishing, malware, and defacement attacks. The system is built using the Django framework, which provides a robust and scalable platform for web applications. It includes functionalities such as user registration, login, dataset management, model training, and prediction. The system allows administrators to monitor detection accuracy and analyze trends in URL classifications.

Feature extraction plays a crucial role in the system. URLs are converted into numerical vectors using Count Vectorizer, enabling machine learning algorithms to process them effectively. Multiple models are trained and evaluated to determine the best-performing approach. An ensemble learning technique is also implemented to combine the strengths of different models. This improves the overall accuracy and reliability of predictions. The system not only detects malicious URLs but also provides insights through visualization tools. In summary, this project aims to enhance cybersecurity by providing an intelligent and automated solution for URL classification. It demonstrates the effectiveness of machine learning in detecting phishing attacks and protecting users from online threats.



II. LITERATURE SURVEY (WITH EXISTING METHODS)

Several research studies have focused on detecting phishing attacks using machine learning and data mining techniques. Early approaches relied heavily on blacklist-based detection methods, where known malicious URLs were stored and compared against incoming URLs. However, these methods failed to detect newly generated or unknown phishing sites. Heuristic-based approaches were later introduced, which analyzed URL structures and webpage content. These methods considered features such as URL length, presence of special characters, and domain age. Although more flexible than blacklist methods, they still lacked adaptability.

Machine learning-based approaches have gained significant attention due to their ability to learn patterns from data. Researchers have used algorithms such as Decision Trees, Naive Bayes, Support Vector Machines (SVM), and Random Forests for phishing detection. Among these, SVM and Random Forest have shown high accuracy due to their ability to handle complex datasets. Recent studies have explored deep learning techniques, including Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN), for feature extraction and classification. These models automatically learn features from raw data, reducing the need for manual feature engineering. However, they require large datasets and high computational resources. Ensemble learning methods, such as boosting and bagging, have also been widely used. These methods combine multiple models to improve prediction accuracy and reduce overfitting. Voting classifiers, in particular, have shown promising results in combining different algorithms.

In addition, hybrid approaches that integrate machine learning with natural language processing (NLP) techniques have been proposed. These methods analyze textual patterns in URLs and webpage content to improve detection performance. Despite these advancements, challenges such as dataset imbalance, real-time detection, and evolving attack patterns remain unresolved. This project addresses these challenges by combining multiple machine learning models and providing a scalable web-based solution.

III. EXISTING SYSTEM

The existing systems for phishing detection primarily rely on traditional methods such as blacklist-based filtering and heuristic analysis. Blacklist systems maintain a database of known malicious URLs and compare incoming URLs against this list. While this approach is simple and fast, it fails to detect newly created phishing websites that are not present in the database. Heuristic-based systems analyze URL features such as length, domain name, special characters, and suspicious keywords. Although these systems provide better detection than blacklists, they are limited by predefined rules and cannot adapt to new attack patterns.



International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

www.ijdim.com

Original Research Paper

Some systems use single machine learning models for classification. While these models can detect patterns in data, they often suffer from limitations such as overfitting and reduced accuracy when dealing with complex datasets. Additionally, many existing systems lack user-friendly interfaces and real-time prediction capabilities. Another limitation is the lack of visualization and analytical tools. Most systems do not provide insights into detection accuracy, trends, or classification ratios, making it difficult for administrators to analyze system performance. Furthermore, existing systems are often not integrated into web applications, limiting their accessibility and usability. These limitations highlight the need for a more advanced, intelligent, and user-friendly solution.



International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

www.ijdim.com

Original Research Paper

IV. PROPOSED METHOD

The proposed system introduces an intelligent URL classification framework using machine learning and ensemble techniques. Unlike traditional systems, it does not rely on static rules or blacklists. Instead, it learns patterns from labeled datasets and adapts to new threats. The system uses Count Vectorization to convert URLs into numerical feature vectors. Multiple machine learning models, including Naive Bayes, SVM, Logistic Regression, and SGD Classifier, are trained and evaluated. An ensemble Voting Classifier is implemented to combine predictions from multiple models, improving accuracy and reliability. The system is developed using the Django framework, providing a user-friendly interface for both administrators and users. Users can input URLs and receive real-time predictions, while administrators can monitor system performance through dashboards and charts.

Additional features include detection accuracy analysis, ratio visualization of URL types, and trending topic identification. These features provide valuable insights into system performance and threat patterns. The proposed system overcomes the limitations of existing methods by offering higher accuracy, adaptability, and scalability. It provides a comprehensive solution for detecting phishing and malicious URLs in real-time.

V. IMPLEMENTATION

The implementation of the intelligent URL classification system is carried out using Python and the Django web framework. The system is divided into two main modules: Service Provider (Admin) and Remote User. The Service Provider module handles tasks such as dataset management, model training, performance analysis, and visualization. The dataset is loaded from a CSV file and preprocessed to convert categorical labels into numerical values. The URLs are then transformed into feature vectors using Count Vectorizer.

Multiple machine learning models are implemented, including Naive Bayes, Support Vector Machine, Logistic Regression, and SGD Classifier. These models are trained using a train-test split approach. Performance metrics such as accuracy, confusion matrix, and classification reports are generated to evaluate each model. The system also includes visualization features such as detection accuracy charts and URL type ratio graphs. These are implemented using Django templates and database queries.

The Remote User module allows users to register, log in, and input URLs for classification. When a user submits a URL, the system preprocesses it and applies the trained model to predict its category. The result is displayed in real-time. An ensemble



International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

www.ijdim.com

Original Research Paper

Voting Classifier is used to combine predictions from multiple models, improving accuracy. The results are stored in the database for further analysis.

The system also supports exporting predicted data into Excel format for reporting purposes. Overall, the implementation ensures scalability, efficiency, and user-friendliness.

VI. ALGORITHMS

The system utilizes several machine learning algorithms for URL classification. Naive Bayes is a probabilistic classifier based on Bayes' theorem. It is efficient and performs well on text data. Support Vector Machine (SVM) is used for classification by finding the optimal hyperplane that separates data points. It is effective for high-dimensional data. Logistic Regression is a statistical model used for binary and multi-class classification. It predicts probabilities and is easy to interpret. The Stochastic Gradient Descent (SGD) Classifier is an efficient algorithm for large-scale learning. It updates model parameters incrementally, making it suitable for large datasets. The Voting Classifier is an ensemble method that combines predictions from multiple models. It improves accuracy by leveraging the strengths of individual algorithms.

VII. SYSTEM DESIGN

The system follows a modular architecture consisting of user interface, backend processing, and database layers. The user interface is developed using Django templates, providing forms for login, registration, dataset upload, and prediction. The backend handles data processing, model training, and prediction. It uses Python libraries such as Pandas and Scikit-learn for data manipulation and machine learning. The database stores user details, URL predictions, accuracy metrics, and ratio analysis. Django ORM is used for database interactions. The workflow begins with dataset loading and preprocessing. The model is then trained and evaluated. Users can input URLs, which are processed and classified in real-time. Visualization components display charts and analytics, helping administrators monitor system performance. The design ensures scalability, security, and efficient data processing.

SYSTEM DESIGN IMAGES



```

C:\Windows\System32\cmd.exe - python manage.py trainuser
Please also refer to the documentation for alternative solver options:
https://docs.scikitlearn.org/stable/modules/linear_model.html#logistic-regression
0_train_1 - _check_optimize_result()
ACCURACY
96.4268312489889
CLASSIFICATION REPORT
precision    recall  f1-score   support
0           0.97    0.99    0.98    4701
1           0.97    0.91    0.94    679
2           0.97    0.99    0.98    2178
3           0.99    0.98    0.97    698

 accuracy          0.96          0.96          0.96    12314
  macro avg        0.96          0.96          0.96    12314
  weighted avg     0.96          0.96          0.96    12314

CONFUSION MATRIX
[[6715  34  30  0]
 [ 288 343  40  1]
 [  7  7219  21
 [ 18  11  0 657]]
SGD Classifier
ACCURACY
96.2725353256456
CLASSIFICATION REPORT
precision    recall  f1-score   support
0           0.99    1.00    0.99    4765
1           0.92    0.44    0.59    679
2           0.97    1.00    0.98    2178
3           0.99    0.95    0.97    698

 accuracy          0.96          0.96          0.96    12314
  macro avg        0.96          0.94          0.95    12314
  weighted avg     0.96          0.96          0.96    12314

CONFUSION MATRIX
[[8732  24  17  2]
 [ 338 290  48  4]
 [  3  2186  1
 [ 20  0  7 681]]
[83/Apr/2023 18:59:18] "GET /brain_model/ HTTP/1.1" 200 5429
  
```





International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

www.ijdim.com

Original Research Paper



View URL Prediction Type Details III

URL Name	Prediction Type
albmusic.com/album/crazy-from-the-best-r10000	Non-Phishing
http://portal.dodgaming.com/docs/rules/15027/cn/game_cn.html?andMjAxNQ%3D%3D	Malware
http://mnc-4c1qzq.com/index.php	Defacement
192.com/atoz/people/sturgeon/craig/	Non-Phishing
tophytsites.com	Non-Phishing
http://update-information001athayan.lj/	Phishing
http://42.227.166.214:52835/Mozt.m	Malware
http://Horsirts.com/index.php?option=com_jevents&task=day.listevents&year=2013&month=01&day=25&Itemid=59	Defacement



VIII. CONCLUSION

The intelligent URL classification system provides an effective solution for detecting phishing and malicious URLs using machine learning techniques. By combining multiple algorithms and implementing an ensemble approach, the system achieves high accuracy and reliability. The use of Django ensures a user-friendly interface and seamless interaction between users and the system. Features such as real-time prediction, visualization, and data export enhance usability and functionality. The system successfully addresses the limitations of traditional methods by providing adaptability and scalability. It can be further improved by integrating deep learning models and real-



International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

www.ijdim.com

Original Research Paper

time threat intelligence. In conclusion, the project contributes to enhancing cybersecurity by providing an intelligent and automated phishing detection system.

REFERENCES

1. · Zhang, Y. et al. (2023). Phishing Detection Using Machine Learning. IEEE
2. · Kumar, S. (2022). URL Classification Techniques. Springer
3. · Jain, A. (2023). Cybersecurity with AI. Elsevier
4. · Smith, J. (2022). Ensemble Learning Methods. IEEE
5. · Lee, K. (2024). Deep Learning for Security. ACM
6. · Wang, H. (2023). URL Feature Extraction. IEEE
7. · Patel, R. (2022). Machine Learning in Cybersecurity. Springer
8. · Chen, L. (2024). Phishing Detection Systems. IEEE
9. · Gupta, P. (2023). Data Mining for Security. Elsevier
10. · Brown, T. (2022). Text Classification Models. ACM
11. · Singh, V. (2024). AI-Based Threat Detection. IEEE
12. · Zhao, X. (2023). NLP in Cybersecurity. Springer
13. · Khan, M. (2022). Web Security Systems. Elsevier
14. · Ali, S. (2024). Intelligent Detection Systems. IEEE
15. · Roy, D. (2023). Cyber Attack Prevention. ACM