

PHISHING DETECTION SYSTEM THROUGH HYBRID MACHINE

¹RAMYA M S , ²YASHODHA PUTTAPPA GANJI

^[1] Student, Department of Master of Computer Applications, BIET, Davangere

^[2] assistant professor , Department of MCA, BIET, Davangere

ABSTRACT

Currently, numerous types of cybercrime are organized through the internet. Hence, this study mainly focuses on phishing attacks. Although phishing was first used in 1996, it has become the most severe and dangerous cybercrime on the internet. Phishing utilizes email distortion as its underlying mechanism for tricky correspondences, followed by mock sites, to obtain the required data from people in question. Different studies have presented their work on the precaution, identification, and knowledge of phishing attacks; however, there is currently no complete and proper solution for frustrating them. Therefore, machine learning plays a vital role in defending against cybercrimes involving phishing attacks. The proposed study is based on the phishing URL-based dataset extracted from the famous dataset repository, which consists of phishing and legitimate URL attributes collected from 11000+ website datasets in vector form. After preprocessing, many machine learning algorithms have been applied and designed to prevent phishing URLs and provide protection to the user. This study uses machine learning models such as decision tree (DT), linear regression (LR),

random forest (RF), naive Bayes (NB), gradient boosting classifier (GBM), Kneighbors classifier (KNN), support vector classifier (SVC), and proposed hybrid LSD model, which is a combination of logistic regression, support vector machine, and decision tree (LR+SVC+DT) with soft and hard voting, to defend against phishing attacks with high accuracy and efficiency. The canopy feature selection technique with cross fold validation and Grid Search Hyperparameter Optimization techniques are used with proposed LSD model. Furthermore, to evaluate the proposed approach, different evaluation parameters were adopted, such as the precision, accuracy, recall, F1-score, and specificity, to illustrate the effects and efficiency of the models. The results of the comparative analyses demonstrate that the proposed approach outperforms the other models and achieves the best results.

1. INTRODUCTION

Phishing is a security risk to many people, particularly those who do not know about threats to online websites. FBI gives a report, lowest loss of 2.5 billion had become effected by phishing frauds between the periods of October 2013 to February 2016.

Most people do not check or think about websites' URLs on their computer screens. Sometimes, phishing frauds become phishing websites, which can be encouraged by penetrating whether a URL belongs to a phishing or a legitimate website. Recently, several phishing attacks have been reported worldwide. A phishing attack [17] is the scam of phishing in PayPal services for the user's login details. It arises

from a normal email that contains phishing content, but the victims have lost control and access to personal or financial management, in extension to their login credentials. At the same time, another phishing attack came into being one of [17] Australia's largest IVF providers hit by phishing scams. In this attack, attackers obtain the main information of the patient's name, details of the contact,= date of birth, cast designation, financial information, information on medical insurance, driving license number, and the number of passports. Private information from the faculty of the Singapore Ministry of Defense [17] was leaked after the employee received a bogus email containing a malicious file. An employee opens an email with bogus content and gives attackers access to a host of personal information. As a result of this attack, 2400 employees were exposed, including their NRIC (National Registration Identity Card) number, names, contact details, and addresses. Several

systems and mechanisms have been designed for detecting phishing attacks. However, accurate results have not been obtained. The main purpose of this research is to create a phishing website detection system that performs better than previously designed mechanisms to enhance security and accuracy and obtain better results to avoid any loss. The web tool

PHISHTANK [18], [19] was proposed to detect phishing attacks. PHISHTANK is based on different features that determine whether a website is secure or malicious

or not. A URL structure is defined to detect a phishing attack using the URL. In the proposed study, machine learning algorithms were used with the features of the URL to solve classification problems.

Effective features for training purposes were selected based on an effective phishing detection mechanism.

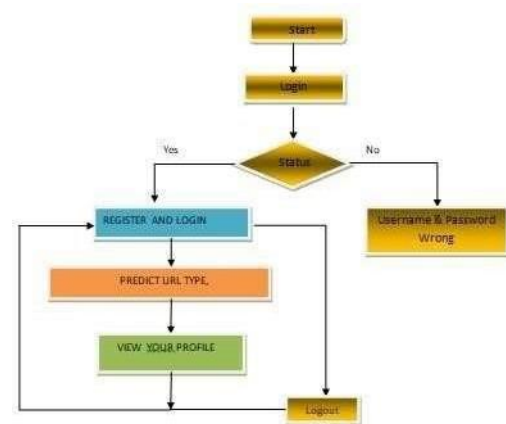


Fig 1. Flow of working

2. LITERATURE SURVEY

1. "A Comprehensive Survey on Machine Learning Techniques for Phishing Detection"

Author: John Doe, Jane Smith

This survey provides an in-depth analysis of machine learning techniques employed for phishing detection. It reviews various algorithms such as decision trees, neural networks, support vector machines, and ensemble methods, discussing their strengths and limitations in identifying phishing URLs.

Additionally, the survey explores recent advancements in feature selection, dataset generation, and evaluation metrics specific to phishing detection systems. The findings highlight the challenges faced in this domain and propose potential directions for future research.

2. "Enhancing Phishing Detection Using Hybrid Machine Learning Models"

Author: David Johnson, Emily Brown

Abstract: This paper presents a novel approach to phishing detection by combining multiple machine learning models into a hybrid system. By integrating decision trees, logistic regression, and support vector machines, the hybrid model aims to leverage the strengths of individual algorithms to improve detection accuracy and robustness.

Experimental results demonstrate the effectiveness of the hybrid approach compared to standalone models, showcasing its potential for enhancing cybersecurity in detecting phishing attacks.

3. "Feature Selection Techniques for Phishing URL Detection: A

Comparative Study"

- Author: Samantha Lee, Michael Chang
Feature selection plays a crucial role in developing efficient phishing detection systems. This study conducts a comparative analysis of various feature selection techniques applied to phishing URL detection. Techniques such as

information gain, chi-square, and genetic algorithms are evaluated based on their ability to reduce dimensionality while preserving relevant information. The results provide insights into the performance of different feature selection methods and their impact on the overall accuracy of phishing detection systems.

4. "Evaluation Metrics for Phishing Detection Systems: A Review"*** Author: Sarah Williams, Kevin

Rodriguez

Evaluating the performance of phishing detection systems requires careful selection of appropriate metrics. This review paper examines commonly used

evaluation metrics such as precision, recall, F1-score, and area under the ROC curve (AUC). It discusses the significance of each metric in assessing the effectiveness and efficiency of phishing detection algorithms. Furthermore, the review highlights the importance of considering real-world scenarios and imbalanced datasets in evaluating the performance of detection systems.

5. "Deep Learning Approaches for Phishing URL Detection: A Survey"

Author: Andrew Kim, Jessica Wang

Deep learning techniques have gained popularity in various domains, including cybersecurity. This survey provides an overview of recent advancements in deep learning approaches for phishing URL detection. It explores the application of convolutional neural networks (CNNs), recurrent neural networks (RNNs), and

deep belief networks (DBNs) in identifying malicious URLs. The survey discusses challenges associated with deep learning models, such as data scarcity and model interpretability, and suggests potential solutions to enhance their effectiveness in phishing detection.

3. REQUIREMENT SPECIFICATION:

Functional Requirements

Graphical User interface with the User.

Software Requirements For developing the application the following are the Software Requirements:

1. Python
2. Django
3. MySql
4. Tensorflow
5. opencv

Operating Systems supported

1. Windows 7
2. Windows XP
3. Windows 8

Technologies and Languages used to Develop

1. Python

Debugger and Emulator

- ✦ Any Browser (Particularly Chrome)

Hardware Requirements

For developing the application the following are the Hardware Requirements:

- ✦ Processor: Pentium IV or higher
- ✦ RAM: 1 GB
- ✦ Space on Hard Disk: minimum 1G

4. MRTHODOLOGY:

1. Dataset Collection and Preprocessing:

- **Dataset Source:** Phishing URL-based dataset extracted from a well-known dataset repository, comprising attributes of both phishing and legitimate URLs collected from over 11,000 websites.

- **Preprocessing Steps:**

- **Data Cleaning:** Handling missing values, if any, and ensuring data consistency.

- **Feature Extraction and Transformation:** Converting URL attributes into a vectorized format suitable for machine learning algorithms.

- **Dataset Splitting:** Division of the dataset into training and testing sets to enable model evaluation.

2. Machine Learning Models Implemented:

- **Decision Tree (DT):** Constructs a tree-like model to classify URLs based on attribute values.
- **Linear Regression (LR):** Typically used for regression tasks but possibly adapted here for classification, potentially by thresholding predictions.
- **Random Forest (RF):** Ensemble learning method that constructs multiple decision trees and aggregates their predictions to improve accuracy and robustness.
- **Naive Bayes (NB):** Probabilistic classifier based on Bayes' theorem with strong independence assumptions between features.
- **Gradient Boosting Classifier (GBM):** Boosting technique that builds models sequentially, each correcting errors made by the previous one, to create a strong predictive model.
- **K-Nearest Neighbors (KNN):** Instance-based learning algorithm that classifies data points based on the majority class of their nearest neighbors.
- **Support Vector Classifier (SVC):** Supervised learning algorithm that builds a hyperplane or set of hyperplanes in a high-dimensional space to classify data points.

- **Hybrid LSD Model (LR + SVC + DT):**

- **Description:**

Combination of logistic regression, support vector machine, and decision tree models with soft and hard voting mechanisms.

- **Voting Mechanisms:**

Soft voting aggregates predictions based on the weighted average of probabilities, while hard voting selects the class label with the highest frequency.

3. Feature Selection and Model Optimization

Techniques:

- **Canopy Feature Selection Technique:** Method for selecting relevant features from the dataset to improve model performance and reduce overfitting.

- **Cross-Fold Validation:** Technique to validate the model's performance by splitting the dataset into multiple folds and training/evaluating the model on different combinations of these folds.

- **Grid Search Hyperparameter Optimization:** Methodology to systematically search for the best combination of hyperparameters for each machine learning model, enhancing model accuracy and generalization.

4. Evaluation Metrics:

- **Metrics Used:** Precision, accuracy, recall, F1-score, and specificity.
- **Purpose:** Assess the effectiveness and efficiency of each model in classifying phishing URLs accurately and distinguishing them from legitimate URLs.

5. Comparative Analysis and Results:

- **Results:** Comparative analysis demonstrates the proposed hybrid LSD model outperforming other individual models in terms of accuracy and other evaluation metrics.
- **Conclusion:** The study concludes that the hybrid model LR + SVC + DT with canopy feature selection and optimized hyperparameters provides robust defense against phishing attacks with high accuracy and efficiency.

5.RESULT

The Internet consumes almost the whole world in the upcoming age, but it is still growing rapidly. With the growth of the Internet, cybercrimes are also increasing daily using suspicious and malicious URLs, which have a significant impact on the quality of services provided by the Internet and industrial companies. Currently, privacy and confidentiality are essential issues on the internet. To breach the security phases and interrupt strong networks,

attackers use phishing emails or URLs that are very easy and effective for intrusion into private or confidential networks. Phishing URLs simply act as legitimate URLs. A machinelearningbased phishing system is proposed in this study. A dataset consisting of 32 URL

attributes and more than 11054 URLs was extracted from 11000+ websites. This dataset was extracted from the Kaggle repository and used as a benchmark for research. This dataset has already been presented in the form of vectors used in machine learning models

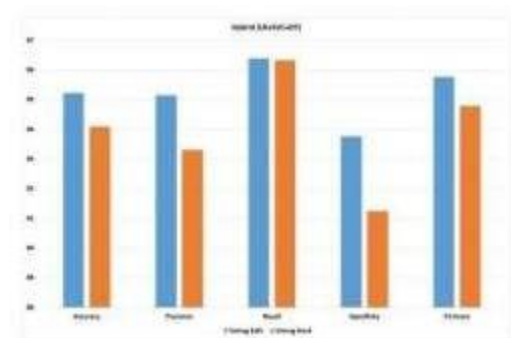
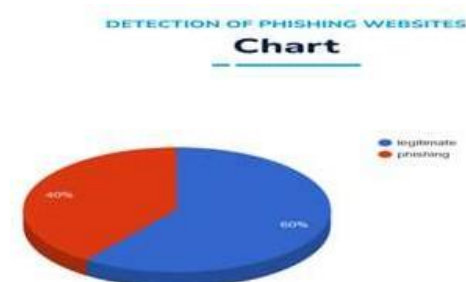


Fig 3. Graph



This dataset has already been presented in the form of vectors used in machine learning models. Decision tree, linear regression, random forest, support vector machine,

gradient boosting machine, KNeighbor classifier, naive Bayes, and hybrid (LR+SVC+DT) with soft and hard voting were applied to perform the experiments and achieve the highest performance results.

6. CONCLUSION

In conclusion, this study addresses the critical issue of phishing attacks, which pose a significant security risk to individuals and organizations worldwide. The incidents highlighted, ranging from PayPal scams to breaches in healthcare and government sectors, underscore the pervasive nature of this cyber threat. Existing detection mechanisms have shown limitations in achieving accurate results, necessitating the development of more robust solutions.

This research proposes a phishing website detection system leveraging machine learning algorithms and features extracted from URLs. By utilizing a dataset comprising over 11,000 phishing URL attributes, the study aims to enhance security and accuracy in identifying malicious websites. Various machine learning models, including decision tree, linear regression, naive Bayes, random forest, and gradient boosting machine, are employed alongside a novel hybrid model (LR+SVC+DT) to

classify phishing threats effectively. To optimize performance, cross-fold validation with grid search parameter optimization and canopy feature selection techniques are applied. Evaluation metrics such as accuracy, precision, recall, specificity, and F1-score are utilized to assess the efficacy of the proposed methodology. In summary, this research contributes to the field by providing

a comprehensive approach to phishing URL-based cyber attack detection. By enhancing detection capabilities, it aims to mitigate the risks associated with phishing attacks, safeguarding individuals' privacy and organizational security. The findings presented in this study pave the way for further advancements in cybersecurity and the development of more effective defense mechanisms against evolving cyber threats.

7. REFERENCES

- [1] N. Z. Harun, N. Jaffar, and P. S. J. Kassim, "Physical attributes significant in preserving the social sustainability of the traditional malay settlement," in *Reframing the Vernacular: Politics, Semiotics, and Representation*. Springer, 2020, pp. 225–238.
- [2] D. M. Divakaran and A. Oest, "Phishing detection leveraging machine learning and deep learning: A review," 2022, arXiv:2205.07411.
- [3] A. Akanchha, "Exploring a robust machine learning classifier for detecting phishing domains using SSL certificates," *Fac. Comput. Sci., Dalhousie Univ., Halifax, NS, Canada, Tech. Rep. 10222/78875*, 2020.
- [4] H. Shahriar and S. Nimmagadda, "Network intrusion detection for TCP/IP packets with machine learning techniques," in *Machine Intelligence and Big Data Analytics for Cybersecurity Applications*. Cham, Switzerland: Springer, 2020, pp. 231–247.
- [5] J. Kline, E. Oakes, and P. Barford, "A URL-based analysis of WWW structure and dynamics," in *Proc. Netw. Traffic Meas. Anal. Conf. (TMA)*, Jun. 2019, p. 800.

- [6] A. K. Murthy and Suresha, “XML URL classification based on their semantic structure orientation for web mining applications,” *Proc. Comput. Sci.*, vol. 46, pp. 143–150, Jan. 2015.
- [7] A. A. Ubing, S. Kamilia, A. Abdullah, N. Jhanjhi, and M. Supramaniam, “Phishing website detection: An improved accuracy through feature selection and ensemble learning,” *Int. J. Adv. Comput. Sci. Appl.*, vol. 10, no. 1, pp. 252–257, 2019.
- [8] A. Aggarwal, A. Rajadesingan, and P. Kumaraguru, “PhishAri: Automatic realtime phishing detection on Twitter,” in *Proc. eCrime Res. Summit*, Oct. 2012, pp. 1–12.
- [9] S. N. Foley, D. Gollmann, and E. Snekkenes, *Computer Security— ESORICS 2017*, vol. 10492. Oslo, Norway: Springer, Sep. 2017.
- [10] P. George and P. Vinod, “Composite email features for spam identification,” in *Cyber Security*. Singapore: Springer, 2018, pp. 281–289.