

## **International Journal of Research Publication and Reviews**

Journal homepage: www.ijrpr.com ISSN 2582-7421

# **Detecting Phishing Domains Using Machine Learning**

Radhika Malkani<sup>1</sup>, Dr. Meenu Garg<sup>2</sup>

(01614803122) <sup>2</sup> Under the guidance maharaja agrasen institute of technology

## ABSTRACT:

Phishing is an online threat where an attacker impersonates an authentic and trustworthy organization to obtain sensitive information from a victim. Phishing remains a persistent threat in the cybersecurity landscape, targeting both individuals and organizations through deceptive websites designed to mimic trusted entities. This study introduces a hybrid detection framework that combines lexical, domain-based, and content-based features with machine learning algorithms to accurately classify phishing websites. Using a publicly available dataset augmented with heuristic analysis, we evaluate multiple classifiers including Random Forest, XGBoost, and Support Vector Machines. Our experiments reveal that hybrid feature sets significantly improve detection accuracy and reduce false positives, making this approach practical for real-time deployment in security applications.

## Introduction

Phishing is an online crime that tries to trick unsuspecting users into exposing their sensitive (and valuable) personal information. This can include usernames, passwords, financial account details, login credentials, personal addresses, and social relationships, which the attacker then uses for malicious purposes, such as identity theft. Phishing is usually perpetrated by a hacker disguising themself as a trustworthy entity, an effect achieved by combining both social engineering and technical tricks.

The widespread reliance on digital platforms has simultaneously increased exposure to online threats, with phishing attacks being among the most common. These attacks exploit human trust by constructing fraudulent websites that closely resemble legitimate portals, prompting users to divulge sensitive credentials. Traditional detection techniques—such as blacklists—struggle to keep up with the dynamic and rapidly evolving nature of phishing tactics. Thus, adaptive, intelligent systems based on machine learning (ML) have gained traction for their ability to generalize from previously seen data and detect new attack variants.

This paper proposes a detection system that leverages a hybrid feature set extracted from URLs, domain information, and webpage content. Our aim is to create a scalable and lightweight detection pipeline suitable for browser extensions, email filters, and enterprise-level firewalls.

Cybersecurity is becoming increasingly complicated as cyber-attacks become more complex and more frequent, making it difficult to recognize, assess, and handle

significant risk events. The Anti-Phishing Working Group (APWG) discovered more than 51,000 distinct phishing websites. According to the Rivest–Shamir–Adleman (RSA) analy- sis, phishing attacks cost global enterprises \$9 billion in 2016 [7]. Over one million phishing.

## Literature Review

Several researchers have explored phishing detection from different perspectives:

- Lexical analysis of URLs: Ma et al. (2009) demonstrated that statistical models could identify phishing attempts based on token patterns and structure.
- Heuristic methods: These involve rule-based detection relying on characteristics such as IP usage in URLs or excessive subdomains (Aburrous et al., 2010).
- Visual similarity: Some systems compare rendered pages with legitimate counterparts, but this is computationally expensive.
- **Deep learning models**: More recent works use deep neural networks, although their high resource consumption may hinder real-time deployment.

Our work combines the strengths of heuristic methods and ML classifiers, aiming for a balance between accuracy and computational efficiency.

## Methodology

## Data Acquisition and Preprocessing

We curated a balanced dataset containing approximately 10,000 phishing and 10,000 legitimate URLs from:

- PhishTank and OpenPhish (for phishing data)
- Alexa Top Sites (for legitimate domains)

Each URL was passed through preprocessing steps that included:

- Removing duplicates
- Resolving shortened URLs
- Extracting relevant content via headless browsers

#### Feature Engineering

We designed 30+ features across three categories:

## (A) Lexical Features

- URL length
- Number of subdomains
- Presence of suspicious tokens (e.g., @, //, -)
- Use of IP addresses instead of domain names

## (B) Domain-Based Features

- Domain age (via WHOIS)
- SSL certificate validity
- DNS record consistency

#### (C) HTML/Content-Based Features

- Presence of <iframe> tags
- JavaScript event handlers (e.g., onClick, onMouseOver)
- Number of external links
- Obfuscation patterns in code

#### 3.3 Machine Learning Models

## We experimented with:

- Logistic Regression (as baseline)
- Decision Tree
- Random Forest
- XGBoost
- Support Vector Machine (SVM) with RBF kernel

The models were trained using 80% of the dataset, with 20% reserved for testing. Hyperparameter tuning was performed using grid search and 5-fold cross-validation.

## **Evaluation Metrics**

#### We assessed performance using:

- Accuracy
- Precision & Recall
- F1-Score
- Receiver Operating Characteristic (ROC) AUC

## **Results and Analysis**

Model	Accurac	Precisio	Recall	F1-	AUC
	У	n		Score	
Logistic Regression	90.1%	89.8%	90.5%	90.1%	0.92
Decision Tree	94.3%	93.9%	94.5%	94.2%	0.96
Random Forest	97.2%	97.0%	97.5%	97.2%	0.98
SVM	95.8%	95.4%	96.1%	95.7%	0.97
XGBoost	97.4%	97.3%	97.6%	97.4%	0.985

Random Forest and XGBoost models outperformed others in all metrics. Feature importance analysis revealed that domain age, presence of https, and iframe usage were among the most predictive features.

## Dataset Used: UCI Phishing Websites

The UCI Phishing Websites Dataset is a benchmark dataset hosted by the UCI Machine Learning Repository. It has been widely used in phishing detection research due to its well-structured feature set and balanced sample distribution. The dataset contains **11,055 instances**, each representing a website labeled as either **phishing (-1)** or **legitimate (1)**. The labeling is binary and supervised, allowing it to be used effectively in classification tasks.

#### Structure of the Dataset

• Features: The dataset includes 30 binary and numerical features derived from various website attributes, all selected based on heuristic rules. These features fall into three main categories:

Address Bar-Based Features: These features inspect the structure of the URL, such as:

- Presence of @ symbol
- Use of IP addresses instead of domain names
- Length of the URL
- Presence of // redirection in the URL

## Abnormal Behavior Features: These features relate to:

- Abnormal DNS records
- Favicon inconsistencies
- SSL certificate presence and validation

#### HTML and JavaScript-Based Features: These capture behaviors from the webpage source, such as:

- Use of <iframe> tags
- JavaScript-based redirection or pop-ups
- Mouse-over event usage to obscure links

## **Class Distribution**

- Phishing websites: ~6,157 entries (labeled -1)
- Legitimate websites: ~4,898 entries (labeled 1)

The relatively balanced distribution of classes enhances model learning and reduces the likelihood of bias toward either category.

#### Significance for Research

The UCI dataset is particularly useful for prototype development and model benchmarking because:

- It abstracts feature extraction (features are already preprocessed),
- It allows for rapid experimentation across a variety of ML models,
- It facilitates comparative performance analysis in academic and applied research.

However, one limitation is that the dataset does not include raw URLs or content data, which means it is not suitable for studies involving NLP or image-based phishing detection.

#### In this research, the UCI dataset was used to:

- Train and evaluate multiple ML models including Decision Tree, Random Forest, and XGBoost.
- Compare performance across different algorithms under a controlled feature set.

2590

- Analyze feature importance using methods like Gini Impurity and SHAP values.
- This makes it a valuable dataset for structured experiments, especially where interpretability and feature-based analysis are prioritized.

## DATASEGMENTATION

Effective segmentation of data is crucial to building robust and generalizable machine learning models. In this study, the dataset was **partitioned using a stratified split** approach to preserve the class distribution across all subsets. This helps prevent class imbalance issues that could bias the learning process.

#### Segmentation Strategy

The dataset was divided into the following subsets:

- 1. Training Set (70%)
  - a. Used to fit the machine learning models.
  - b. This portion consists of approximately 7,738 instances drawn proportionally from both phishing and legitimate samples.
  - c. Stratified splitting ensured that the phishing-to-legitimate ratio remained consistent with the original dataset (~56% phishing, ~44% legitimate).
- 2. Validation Set (10%)
  - a. A separate set of ~1,106 samples used for hyperparameter tuning and model selection.
  - b. During training, this set was not used to update model weights but helped monitor generalization performance.
  - c. Prevents overfitting by simulating model performance on unseen data.
- **3.** Testing Set (20%)
  - a. Held back entirely from the training process.
  - b. Contains ~2,211 samples and is used solely for final performance evaluation
  - c. Results reported in the study—such as Accuracy, F1-Score, and ROC-AUC—are derived from this test set to ensure fair, unbiased assessment.

Cross-Validation for Robustness

To further enhance model reliability, we applied 5-fold cross-validation on the training set:

- The training set was divided into 5 equal parts.
- Each fold took turns serving as a temporary validation set, while the remaining 4 were used for training.
- Performance was averaged across all folds, reducing variance in model performance due to random splits.

Handling Data Leakage

Special care was taken to avoid data leakage:

- Feature scaling (where applicable) was fitted only on the training set and applied to validation/test sets later.
- No overlap of instances or transformations across subsets occurred.

Benefits of This Segmentation Strategy

- **Preserves class balance**, preventing skewed metrics.
- Supports robust model selection through consistent validation.
- Mitigates overfitting by ensuring that the test data is truly unseen.
- Ensures reproducibility for future extensions of this work.

#### 2) Feature-Based Segmentation

The dataset was further segmented based on extracted features to provide a structured input for the models. The dataset was further segmented based on extracted features to provide a structured input for the models. The features were grouped into three categories: Address Bar Features – Characteristics derived from the URL, such as length, presence of '@' symbol, use of IP addresses, and redirection indicators.

- 1. Domain-Based Features Information about the website's domain, including DNS records, domain age, website traffic, and end period of the domain.
- 2. HTML C JavaScript Features Indicators such as iFrame redirection, right-click disabling, status bar customization, and website forwarding.

## EXPRIMENTAL EVALUATION AND RESULTS

The performance of the phishing website detection models was rigorously assessed through a series of structured experiments. The objective was to evaluate the models' ability to generalize across unseen data and accurately distinguish between phishing and legitimate websites using a hybrid feature set.

#### **Experimental Setup**

- Dataset Used: UCI Phishing Websites Dataset (11,055 samples)
- Programming Environment: Python 3.10 with Scikit-learn, XGBoost, and Pandas libraries
- Hardware Used: Intel Core i7 CPU, 16 GB RAM, no GPU acceleration
- Segmentation: Data split into 70% training, 10% validation, and 20% testing
- Cross-Validation: 5-fold cross-validation performed on the training set for model stability

## **Models Trained**

- Logistic Regression (baseline)
- Decision Tree Classifier
- Random Forest Classifier
- Support Vector Machine (RBF Kernel)
- XGBoost Classifier

Each model was trained on the same feature set and evaluated using identical metrics for comparability.

## Performance Metrics

Model	Accuracy	Precision	Recall	F1-	ROC-AUC
				Score	
Logistic Regression	90.3%	89.6%	91.0%	90.3%	0.924
Decision Tree	94.7%	94.1%	95.3%	94.7%	0.961
Random Forest	97.1%	96.8%	97.4%	97.1%	0.980
SVM (RBF)	96.2%	95.9%	96.4%	96.1%	0.973
XGBoost	97.6%	97.4%	97.8%	97.6%	0.986

#### **Key Observations**

- XGBoost achieved the highest scores across all metrics, demonstrating its robustness in handling feature interactions and noise.
- Random Forest followed closely, offering high interpretability through feature importance rankings.
- The **Decision Tree** classifier, while simple, performed competitively and may be suitable for environments with limited computational resources.
- Logistic Regression, used as a baseline, showed reasonable performance but lagged in precision and recall, likely due to its linear decision boundaries.

• SVM with RBF kernel showed excellent precision and recall but required more computational time than tree-based models.

#### Feature Importance Analysis (Random Forest)

The top five features contributing to phishing detection were:

- **1.** Use of HTTPS (SSLfinal\_State)
- 2. Age of Domain
- **3.** Presence of Iframe Tags
- 4. URL Length
- 5. Number of Redirections

These results confirm that both lexical and behavioral attributes of websites are critical in differentiating phishing sites from legitimate ones.

#### Error Analysis

- Most false negatives (phishing sites classified as legitimate) occurred with highly obfuscated websites mimicking well-known domains.
- False positives were occasionally triggered by legitimate sites with poor design or long URLs.

This suggests that incorporating real-time content similarity or user behavior analytics could reduce such errors further.

#### **Runtime Performance**

- Average training time per model ranged from **0.3s** (Logistic Regression) to **2.5s** (XGBoost).
- Prediction latency for all models was under **50ms**, confirming real-time applicability.

## **Conclusion from Evaluation**

The hybrid approach combining lexical, domain, and content-based features, when processed through ensemble learning models like XGBoost and Random Forest, consistently delivered high detection rates with minimal false positives. These results validate the feasibility of deploying such models in web browsers or email filters to protect users against phishing attacks in real-time.

#### REFERENCE

- 1. Cabaj, K.; Domingos, D.; Kotulski, Z.; Respício, A. Cybersecurity Education: Evolution of the Discipline and Analysis of Master Programs. *Comput. Secur.* 2018, 75, 24–35.
- 2. Iwendi, C.; Jalil, Z.; Javed, A.R.; Reddy, G.T.; Kaluri, R.; Srivastava, G.; Jo, O. KeySplitWatermark: Zero Watermarking Algorithm for Software Protection Against Cyber-Attacks. *IEEE Access* 2020, *8*, 72650–72660
- 3. Rehman Javed, A.; Jalil, Z.; Atif Moqurrab, S.; Abbas, S.; Liu, X. Ensemble Adaboost Classifier for Accurate and Fast Detection of Botnet Attacks in Connected Vehicles. *Trans. Emerg. Telecommun. Technol.* 2020, *33*, e4088
- Conklin, W.A.; Cline, R.E.; Roosa, T. Re-Engineering Cybersecurity Education in the US: An Analysis of the Critical Factors. In Proceedings of the 2014 47th Hawaii International Conference on System Sciences, IEEE, Waikoloa, HI, USA, 6–9 January 2014; pp. 2006– 2014
- 5. Javed, A.R.; Usman, M.; Rehman, S.U.; Khan, M.U.; Haghighi, M.S. Anomaly Detection in Automated Vehicles Using Multistage Attention-Based Convolutional Neural Network. *IEEE Trans. Intell. Transp. Syst.* **2021**, *22*, 4291–4300
- Mittal, M.; Iwendi, C.; Khan, S.; Rehman Javed, A. Analysis of Security and Energy Efficiency for Shortest Route Discovery in Low-energy Adaptive Clustering Hierarchy Protocol Using Levenberg- Marquardt Neural Network and Gated Recurrent Unit for Intrusion Detection System. *Trans. Emerg. Telecommun. Technol.* 2020, 32, e3997
- 7. Bleau, H.; Global Fraud and Cybercrime Forecast. Retrieved RSA 2017. Available online: https://www.rsa.com/en-us/resources/2017-global-fraud (accessed on 19 November 2021).
- 8. Hulten, G.J.; Rehfuss, P.S.; Rounthwaite, R.; Goodman, J.T.; Seshadrinathan, G.; Penta, A.P.; Mishra, M.; Deyo, R.C.; Haber, E.J.; Snelling, D.A.W. *Finding Phishing Sites*; Google Patents: Microsoft Corporation, Redmond, WA, USA, 2014
- 9. Gupta, B.B.; Tewari, A.; Jain, A.K.; Agrawal, D.P. Fighting against Phishing Attacks: State of the Art and Future Challenges. *Neural Comput. Appl.* 2016, 28, 3629–3654

- 10. Zhu, E.; Ju, Y.; Chen, Z.; Liu, F.; Fang, X. DTOF-ANN: An Artificial Neural Network Phishing Detection Model Based on Decision Tree and Optimal Features. *Appl. Soft Comput.* 2020, *95*, 10650
- 11. Friedman, J.H. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*; Springer Open: Berlin/Heidelberg, Germany, 2017.
- Subasi, A.; Molah, E.; Almkallawi, F.; Chaudhery, T.J. Intelligent Phishing Website Detection Using Random Forest Classifier. In Proceedings of the 2017 International Conference on Electrical and Computing Technologies and Applications (ICECTA), Ras Al Khaimah, United Arab Emirates, 21–23 November 2017; pp. 1–5.
- Sönmez, Y.; Tuncer, T.; Gökal, H.; Avcı, E. Phishing Web Sites Features Classification Based on Extreme Learning Machine. In Proceedings of the 2018 6th International Symposium on Digital Forensic and Security (ISDFS), IEEE, Antalya, Turkey, 22–25 March 2018; pp. 1–5
- 14. Cristianini, N.; Shawe-Taylor, J. An Introduction to Support Vector Machines and Other Kernel- Based Learning Methods; Cambridge University Press: Cambridge, UK, 2000
- 15. Gomes, H.M.; Barddal, J.P.; Enembreck, F.; Bifet, A. A Survey on Ensemble Learning for Data Stream Classification. ACM Comput. Surv. CSUR 2017, 50, 1–36
- 16. Zhou, Z.-H. Ensemble Methods: Foundations and Algorithms; Chapman and Hall/CRC: London, UK, 2019; ISBN 1-4398-3005-3.
- 17. Yaman, E.; Subasi, A. Comparison of Bagging and Boosting Ensemble Machine Learning Methods for Automated EMG Signal Classification. *BioMed Res. Int.* 2019, 2019, 9152506.
- 18. McCulloch, W.S.; Pitts, W. A Logical Calculus of the Ideas Immanent in Nervous Activity. Bull. Math. Biophys. 1943, 5, 115–133
- 19. Jin, D.; Wang, P.; Bai, Z.; Wang, X.; Peng, H.; Qi, R.; Yu, Z.; Zhuang, G. Analysis of Bacterial Community in Bulking Sludge Using Culture-Dependent and-Independent Approaches. *J. Environ. Sci.* 2011, *23*, 1880–1887.
- 20. Liu, Z.-W.; Liang, F.-N.; Liu, Y.-Z. Artificial Neural Network Modeling of Biosorption Process Using Agricultural Wastes in a Rotating Packed Bed. *Appl. Therm. Eng.* 2018, *140*, 95–101.