Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

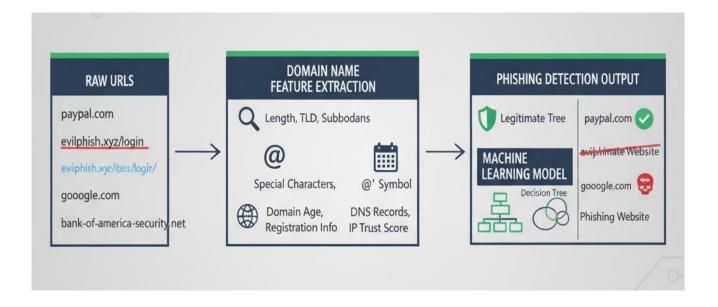
ML-Based Phishing Website Detection Using Domain Name Features

Harnessing Domain Name Intelligence for Phishing Attack Prevention

Aaftab Khan MBA (Data Science), Amity University 7absec@gmail.com

Abstract

Phishing attacks remain one of the most pervasive cybersecurity threats, exploiting human vulnerabilities through deceptive websites. Traditional blacklist and heuristic methods fail to detect newly created and zero-day phishing domains, necessitating proactive approaches. This study proposes a machine learning-based detection system leveraging lexical and domain-level features extracted directly from URLs. Two classification algorithms—Decision Tree (DT) and Support Vector Machine (SVM)—were implemented and evaluated using a balanced dataset of 50,000 phishing and legitimate URLs. Results demonstrate that the Decision Tree outperformed SVM across accuracy (97.11% vs. 92.72%), precision (0.9702 vs. 0.9239), and recall (0.9781 vs. 0.9472). The findings highlight the effectiveness of lightweight ML models in real-time phishing prevention, with practical implications for browser integration, enterprise gateways, and security training.



Keywords — Phishing detection, machine learning, domain name analysis, Decision Tree, SVM, cybersecurity.

1. Introduction

ISSN:2394-2231

Phishing, a form of social engineering, manipulates victims into disclosing sensitive information by imitating trusted organizations. Despite advancements in defense mechanisms, phishing remains the most common entry vector in data breaches [Verizon DBIR 2024]. Traditional defenses—such as blacklists and heuristic rules—struggle with rapidly evolving phishing techniques. Modern phishing campaigns have become increasingly sophisticated, ranging from

International Journal of Computer Techniques – IJCT Volume 12 Issue 5, October 2025

Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

targeted spear phishing and whaling to smishing and voice-based attacks, often exploiting real-world events to increase credibility. Reports from the Anti-Phishing Working Group (2024) indicate over 1.3 million phishing websites were detected in a single quarter, with financial services, SaaS/webmail, and e-commerce among the most targeted sectors. Attackers frequently use tactics such as typosquatting, homograph attacks, and URL obfuscation to evade detection, making manual or rule-based defenses insufficient. In this context, machine learning (ML) offers a proactive solution by analyzing lexical and domain-level features of URLs to identify phishing attempts in real time.

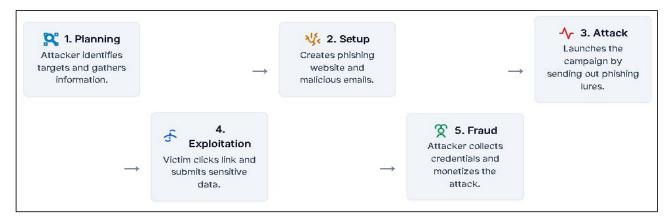


Figure 1: The Phishing Attack Lifecycle

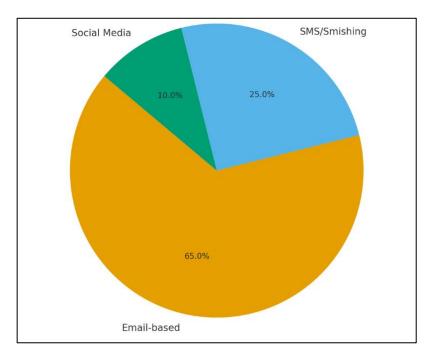


Figure 2: Distribution of Phishing Attack Types (2024)

Machine Learning (ML) introduces predictive capabilities by analyzing structural and lexical characteristics of URLs to classify websites as phishing or legitimate. This paper presents a comparative evaluation of two ML models—Decision Tree and Support Vector Machine—using URL-based features.

Contributions:

ISSN:2394-2231



<u>International Journal of Computer Techniques – IJCT Volume 12 Issue 5, October 2025</u>

Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

- Development of a phishing detection system using domain-level and lexical features.
- 2. Comparative analysis of DT and SVM on a large balanced dataset.
- 3. Discussion of practical deployment scenarios and limitations.

2. Related Work

Early phishing detection relied on blacklists and heuristics, effective only against known threats. URL-based ML approaches (Mohammad et al., 2015) demonstrated improvements by incorporating features such as URL length, domain age, and keyword presence. Deep learning methods (Bahnsen et al., 2017; Li et al., 2023) achieve higher accuracy but require significant computational resources.

Comparative studies (Basit et al., 2021; Zhang et al., 2022) highlight the effectiveness of ensemble models, while adversarial studies reveal vulnerabilities to evasion techniques. The research gap lies in lightweight, interpretable models capable of real-time deployment without extensive resource requirements.

3. Methodology

3.1 Dataset

- **Phishing URLs:** 27,500 from PhishTank archives (2024–2025).
- Legitimate URLs: 27,500 from Alexa Top Sites and Common Crawl.
- **Final Dataset:** ~50,000 balanced instances.

3.2 Feature Engineering

Thirty lexical and domain-based features were used (e.g., URL length, subdomain count, SSL state, domain age). Features were encoded into {-1, 0, 1} categories representing legitimate, suspicious, or phishing indicators.

3.3 Model Development

- Decision Tree (DT): Implemented with scikit-learn, using information gain for feature splits.
- Support Vector Machine (SVM): Implemented with linear kernel. Dataset split: 80% training, 20% testing with stratification.

3.4 Evaluation Metrics

Accuracy, Precision, Recall, F1-score, and Confusion Matrix were used for assessment.

4. Results

4.1 Decision Tree - The Decision Tree classifier was trained using the extracted URL-based features, and its performance was evaluated on the test dataset. The model demonstrated consistently strong results across all evaluation metrics, as summarized below:

Accuracy: 97.11%Precision: 0.9702

Recall: 0.9781F1-score: 0.9741

ISSN:2394-2231

- E1 N (22 ()

• False Negatives: 32 (out of 2,211 test samples).

International Journal of Computer Techniques – IJCT Volume 12 Issue 5, October 2025

Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

Figure 3: Training and evaluating Decision Tree model

4.2 Support Vector Machine - The Support Vector Machine (SVM) classifier with a linear kernel was implemented to provide a comparative baseline against the Decision Tree. The performance of the SVM model on the test dataset is presented below:

Accuracy: 92.72%
Precision: 0.9239
Recall: 0.9472
F1-score: 0.9354
False Negatives: 65.

Figure 4: Training and evaluating SVM model

4.3 Comparative Analysis

The DT outperformed SVM across all metrics. Its ability to handle categorical, discrete features proved advantageous. SVM's linear kernel failed to capture complex feature interactions, leading to higher false negatives.

Metric	Decision Tree	Support Vector Machine (SVM)
Accuracy	0.9711	0.9272
Precision	0.9702	0.9239
Recall	0.9781	0.9472
F1-Score	0.9741	0.9354

Table 1: Comparative Performance of Models

5. Discussion

ISSN:2394-2231

International Journal of Computer Techniques – IJCT Volume 12 Issue 5, October 2025

Open Access and Peer Review Journal ISSN 2394-2231

https://ijctjournal.org/

The study validates the utility of ML for phishing detection. While DT offers superior interpretability and recall, it risks overfitting if not pruned. SVM performed moderately but required longer training times.

Practical Implications:

- Lightweight Deployment: DT models can be integrated into browsers, mobile apps, or email clients.
- Enterprise Gateways: Ensembles (e.g., Random Forest, Gradient Boosting) can further improve accuracy for server-side defenses.
- Awareness Training: Technical tools must be paired with phishing simulations and employee education.

Limitations:

- Dataset primarily English-based, limiting generalizability to multilingual phishing.
- Exclusion of content-based and dynamic features (e.g., HTML scripts, DNS lookups).
- Vulnerable to adversarial attacks and evolving phishing strategies.

6. Conclusion

This research demonstrates that machine learning models using lexical and domain-based features effectively detect phishing websites. The Decision Tree classifier achieved a superior balance of accuracy and interpretability compared to SVM. Future research should explore ensemble learning, adversarial robustness, and hybrid approaches combining lexical, host-based, and content features.

References

ISSN:2394-2231

- Mohammad, R.M. et al. (2015). "Phishing Detection Based on URL Features." Applied Computing and Informatics.
- 2. Bahnsen, A.C. et al. (2017). "DeepPhish: Detecting Phishing URLs using Deep Learning." *eCrime Researchers Summit*
- 3. Basit, A. et al. (2021). "Comparative Study of ML Models for Phishing Detection." *IEEE Access*.
- Odonnat, Ambroise (2024). phishing.arff. figshare. Dataset. https://doi.org/10.6084/m9.figshare.26232710.v1
- 5. Li, X. et al. (2023). "Transformer Models for Phishing Email Detection." ACM TDSC.
- 6. Verizon. (2024). Data Breach Investigations Report.