

Abstract

Phishing attacks pose a significant threat to online security, deceiving users into disclosing sensitive information through fraudulent websites. Traditional detection mechanisms, such as blacklists and heuristics, are increasingly ineffective against sophisticated phishing tactics like domain obfuscation and cloaking. To address this challenge, our project develops a machine learning-based phishing detection system that analyzes website URLs and content to classify them as legitimate or phishing. Using models like **XGBoost**, **Random Forest**, and **Neural Networks**, the system identifies key patterns and features to ensure high accuracy in detection. It is implemented as a scalable and user-friendly web application with browser extension integration, providing real-time alerts to protect users from phishing attacks. This innovative approach bridges advanced machine learning techniques and practical cybersecurity solutions, aiming to create a safer online environment for individuals and organizations.

Introduction

Phishing is one of the most persistent cybersecurity threats, targeting users with fraudulent websites to steal sensitive information like passwords and financial data. Traditional detection methods, such as blacklists and user reports, often fail against evolving phishing tactics like domain obfuscation or cloaking.

To combat these challenges, our project leverages **machine learning** for accurate phishing detection. By analyzing URL and content features, models like XGBoost and Random Forest classify websites in real-time. The system integrates into a web application and browser extension, offering instant alerts and enhancing online security.

Existing System

Several existing applications focus on phishing detection, primarily utilizing a mix of URL analysis, content evaluation, and machine learning techniques. Content-based detection systems analyze the textual content on websites, while URL-based systems look for suspicious patterns, such as abnormal domain names or the presence of special characters.

Machine learning-based applications, such as those employing Random Forest or Support Vector Machines, classify websites by analyzing a range of features extracted from URLs. Despite their utility, these tools have limitations: URL-based systems may fail to detect zero-day attacks, content-based systems are vulnerable to fake copies of legitimate content, and machine learning models require constant retraining and large datasets. Therefore, an integrated, real-time phishing detection solution is needed, which can be deployed efficiently on users' devices.

Proposed System

The proposed system is a real-time phishing detection tool that allows users to input URLs and receive feedback on whether the site is phishing or legitimate. This application is built around a machine learning model trained on a dataset of phishing and legitimate URLs. Features such as domain name length, presence of special characters, and domain reputation are analyzed to classify websites.

The core model uses XGBoost for its ability to process and classify large, complex datasets with high accuracy. The application not only provides a user-friendly interface for URL input and classification but also offers integration with a browser extension for real-time protection. By incorporating machine learning, this system ensures scalability and adaptability against the evolving nature of phishing tactics.

Literature Review:

S. No	Author(s)	Title	Year	Journal/ Source	Methodology	Metris	Demerits
1	Leon Reznik	Computer Security with Artificial Intelligence, Machine Learning, and Data Science Combination	2024	IEEE	Combines AI, ML, and Data Science for computer security.	Comprehensive approach to enhance security systems.	Complex integration of multiple technologies.
2	L. Tang and Q. H. Mahmoud	A Deep Learning-Based Framework for Phishing Website Detection	2023	IEEE Access, vol. 10	Deep learning framework for phishing detection.	High accuracy in detecting phishing websites.	Requires significant computational resources.
3	SatheeshKumar, M., Srinivasagan, K.G.	A lightweight and proactive rule-based incremental construction approach to detect phishing scam	2022	Inf Technol Manag	Rule-based incremental approach for phishing detection.	Lightweight and proactive, suitable for real-time applications.	Limited adaptability to evolving phishing techniques.

S.No	Author(s)	Title	Year	Journal/ Source	Methodology	Metris	Demerits
4	Anti-Phishing Working Group	Phishing Activity Trends Report-1Q	2022	APWG Report [Online]	Analyzes phishing trends across industries.	Comprehensive report on the latest phishing trends.	Doesn't provide specific solutions or methodologies.
5	P. Zhang, A. Oest, H. Cho et al.	CrawlPhish: Large-scale analysis of client-side cloaking techniques in phishing	2021	IEEE Symp. Secur. Privacy (SP)	Analyzes client-side cloaking in phishing websites.	Identifies hidden techniques used by phishing sites.	Focuses only on client-side cloaking.
6	O. K. Sahingoz, U. Cekmez and A. Buldu	Internet of Things (IoTs) Security: Intrusion Detection using Deep Learning	2021	Journal of Web Engineering	Deep learning for intrusion detection in IoT networks.	Enhances security in IoT devices with automated threat detection.	High computational overhead for IoT networks.
7	A. Awasthi and N. Goel	Generating Rules to Detect Phishing Websites Using URL Features	2021	ODICON Conference	Rule generation based on URL features to detect phishing.	Efficient detection using URL characteristics.	Not suitable for complex phishing attempts involving obfuscation.

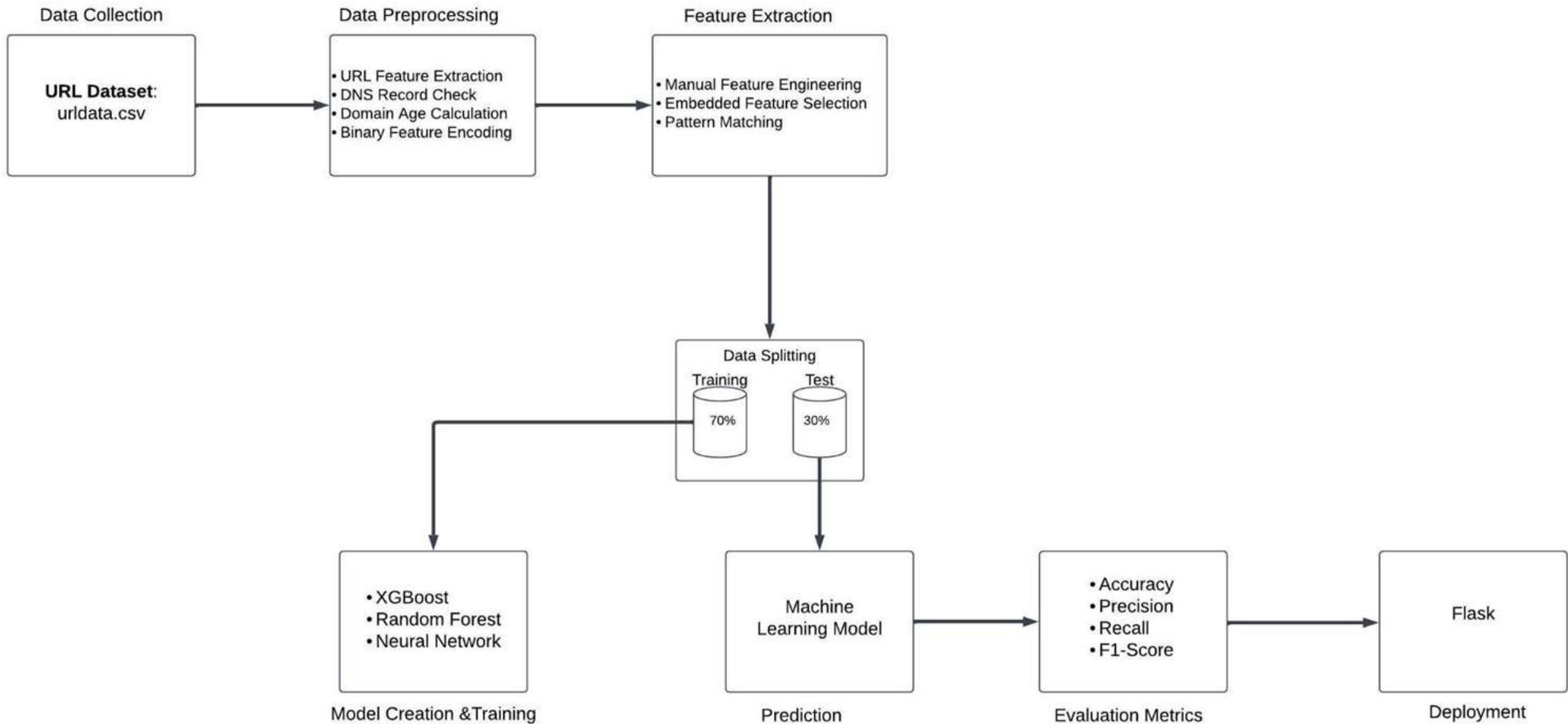


FIGURE-1: BLOCK DIAGRAM OF PHISHING WEBSITE DETECTION

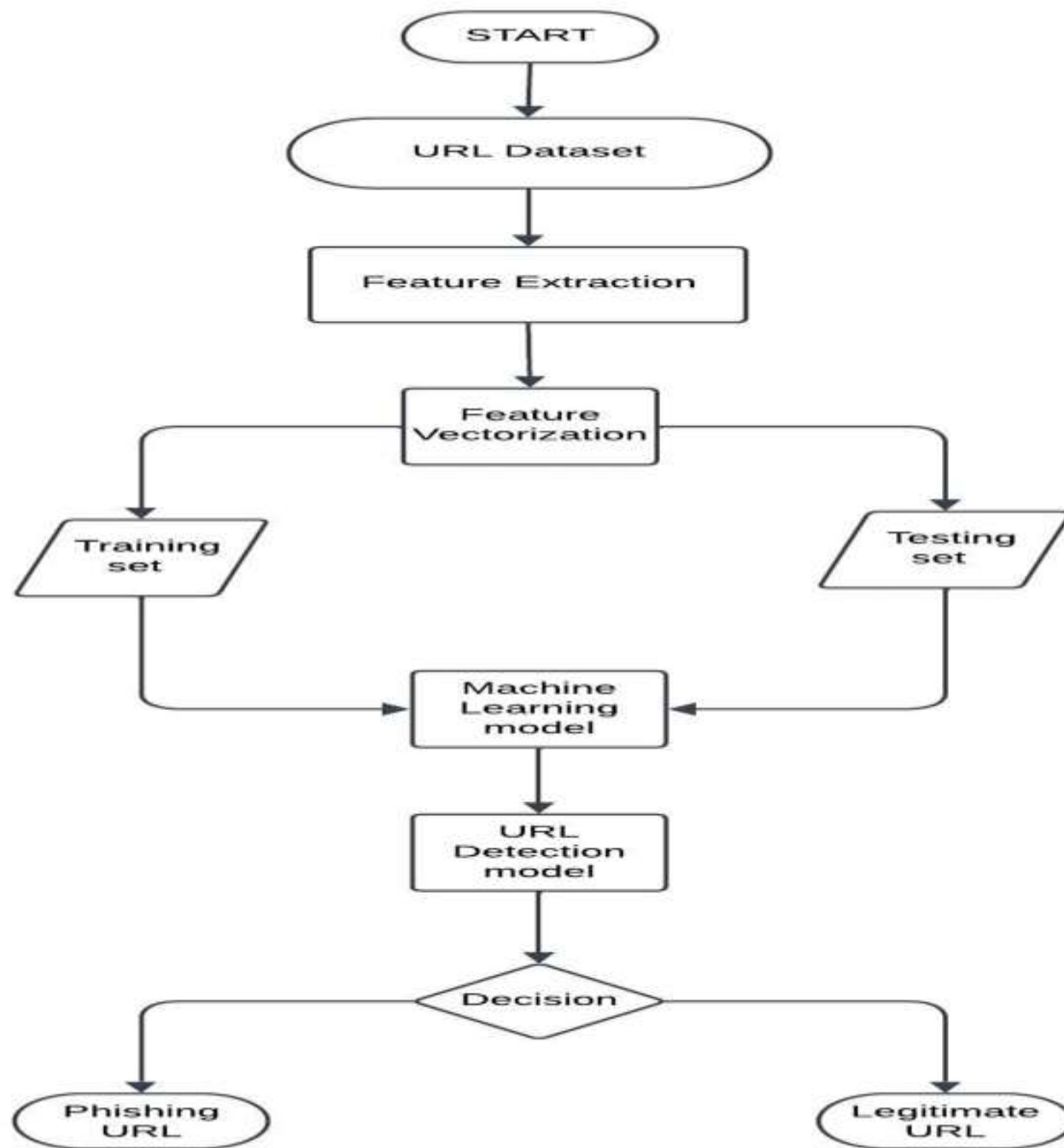


FIGURE-2: FLOWCHART DIAGRAM OF PHISHING WEBSITE DETECTION

Technical Specification

- **Software Requirements**

- ***Text Processing:*** NLTK
- ***Machine Learning:*** Scikit-learn
- ***Data Manipulation:*** Pandas and NumPy
- **IDE:** Jupyter Notebook
- **Datasets:** Kaggle Phishing Dataset

- **Hardware Requirements**

- **Processor:** Modern multi-core CPU for efficient processing and analysis.
- **RAM:** 8 GB to handle large datasets and model computations.
- **Storage:** SSD with at least 256 GB of space for data and software installation.
- **Network:** Stable internet connection for accessing datasets and cloud-based resources.

Conclusion

The phishing detection web application developed in this project leverages machine learning techniques, particularly the XGBoost model, to effectively identify phishing websites. By analyzing key features of URLs, the system can accurately classify websites as phishing or legitimate, providing an essential tool for users to protect themselves against phishing attacks. With real-time URL analysis, helping users make informed decisions about whether to trust a website or not. In addition to the robust backend functionality powered by machine learning, the frontend of the application provides an intuitive and user-friendly interface that simplifies the detection process for all users. This project demonstrates the potential of machine learning offering a practical solution for detecting phishing websites. It successfully bridges the gap between sophisticated algorithms and user-friendly applications, giving individuals a powerful tool to safeguard their online activities.

Bibliography

- [1] Leon Reznik," Computer Security with Artificial Intelligence, Machine Learning, and Data Science Combination," in Intelligent Security Systems: How Artificial Intelligence, Machine Learning and Data Science Work For and Against Computer Security, IEEE, 2024
- [2] L. Tang and Q. H. Mahmoud, "A Deep Learning-Based Framework for Phishing Website Detection," in IEEE Access, vol. 10, pp. 1509-1521, 2023
- [3] SatheeshKumar, M., Srinivasagan, K.G. UnniKrishnan, G. A lightweight and proactive rule-based incremental construction approach to detect phishing scam. Inf Technol Manag (2022)
- [4] Anti-Phishing Working Group—APWG. (2022). Phishing Activity Trends Report-1Q. [Online]. Available: https://docs.apwg.org/reports/apwg_trends_report_q1_2022.pdf
- [5] P. Zhang, A. Oest, H. Cho, Z. Sun, R. Johnson, B. Wardman, S. Sarker, A. Kapravelos, T. Bao, R. Wang, Y. Shoshitaishvili, A. Doupé, and G.-J. Ahn, "CrawlPhish: Large-scale analysis of client-side cloaking techniques in phishing," in Proc. IEEE Symp. Secur. Privacy (SP), May 2021
- [6] O. K. Sahingoz, U. Cekmez and A. Buldu, "Internet of Things (IoT)s Security: Intrusion Detection using Deep Learning" 2021, Journal of Web Engineering, 2021
- [7] A. Awasthi and N. Goel, "Generating Rules to Detect Phishing Websites Using URL Features," 2021 1st Odisha International Conference on Electrical Power Engineering, Communication and Computing Technology(ODICON), 2021
- [8] E. Kocyigit, M. Korkmaz, O.K. Sahingoz, B. Diri, "Real-Time ContentBased Cyber Threat Detection with Machine Learning". In: Abraham, A., Piuri, V., Gandhi, N., Siarry, P., Kaklauskas, A., Madureira, A. (eds) Intelligent Systems Design and Application, 2021
- [9] Abdullateef O. et al., "Improving the phishing website detection using empirical analysis of Function Tree and its variants", Heliyon, vol 7, Issue 7, 2021
- [10] A.V. Ramana, Rao, K.L. Rao, R.S. Stop-Phish: an intelligent phishing detection method using feature selection ensemble. Soc. Netw. Anal. Min. 11, 110 (2021)