

- **Title:** *Phishing URL Detection Using Machine Learning*
 - **Submitted by:** Saniya Santosh Choughule & Iqra Mohd Nisar Khan.
 - **Internship:** Cybersecurity Internship 2025
 - **Organization:** Digisuraksha Parhari Foundation, Powered by Infinisec Technologies Pvt. Ltd.
-



1. Introduction

Phishing is a type of cyberattack where attackers trick users into revealing sensitive information such as passwords, credit card numbers, or login credentials by disguising as a trustworthy entity—often through fake websites or deceptive links.

It is dangerous because phishing can lead to identity theft, financial loss, unauthorized access to private accounts, and large-scale data breaches for organizations.

The goal of our project is to build a tool that can detect phishing websites using only the URL as input. By analyzing patterns in the structure and content of URLs, our model can predict whether a given link is safe or malicious, helping users avoid falling victim to phishing attacks.



2. Abstract

Phishing is a prevalent and dangerous form of cyberattack in which attackers deceive users into clicking malicious links that appear to be from legitimate websites. These attacks often lead to identity theft, financial loss, and unauthorized access to sensitive data. With the increasing sophistication of phishing techniques, traditional blacklist-based detection methods are no longer sufficient. This research aims to develop a machine learning-based tool that can detect phishing websites by analyzing URL structures alone. We used a publicly available dataset containing labeled phishing and legitimate URLs, from which various lexical and structural features were extracted. A Random Forest classification model was trained to identify patterns associated with malicious behavior. We also built a command-line interface for real-time URL prediction. The tool demonstrated high accuracy and robustness during testing, offering a simple yet effective way to identify phishing threats. Our findings show that machine learning can be a powerful and practical solution for enhancing cybersecurity and protecting users online.

3. Problem Statement & Objective

Phishing websites pose serious security threats by stealing sensitive information. Traditional blacklist methods are ineffective against new and evolving phishing attacks. The objective of this project is to build a machine learning model that can detect phishing websites based on URL features and help users identify threats in real-time.

4. Literature Review

Several studies have explored phishing detection using various approaches:

- Traditional methods rely on blacklists, which are outdated and cannot detect zero-day attacks.
- ML models have been proposed that extract lexical and statistical features from URLs.
- Prior research shows that Random Forest, Decision Trees, and SVM are effective in phishing detection.
- Datasets from PhishTank and Kaggle have been widely used.

References:

- IEEE papers on phishing detection
 - PhishTank.com dataset
 - "Machine Learning for Phishing Detection" – Journal of Cybersecurity
-

5. Research Methodology

- **Data Collection:** Used a dataset with phishing and legitimate URLs.
 - **Feature Engineering:** Extracted URL-based features such as length, presence of IP, number of dots, use of HTTPS, etc.
 - **Model Training:** Trained a Random Forest Classifier using scikit-learn.
 - **Testing & Evaluation:** Evaluated with accuracy, precision, recall, and F1 score.
-

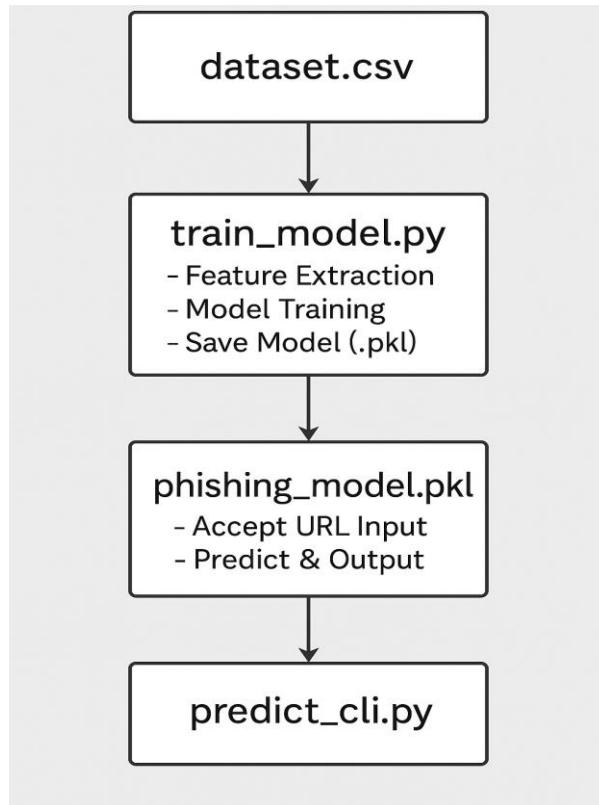
6. Tool Implementation

- **train_model.py:** Loads the dataset, extracts features, trains the model, and saves it as phishing_model.pkl.
- **predict_cli.py:** Loads the model and accepts a URL from the user, then predicts whether it is phishing or legitimate.

- **Environment:**

- Python,
- Scikit-learn,
- Pandas

Phishing Detection System Workflow:



🛠 7. Project Structure

```
Phishing_URL_Detection/
|
|   tool/
|   |   source_code/
|   |   |   train_model.py      # Trains the ML model using URL dataset
|   |   |   predict_cli.py     # CLI for predicting phishing/legit URLs
|   |   |   phishing_model.pkl # Saved trained ML model
|   |   |   dataset.csv        # Dataset used for training/testing
|
|   |
|   |   phishing-env/          # Python virtual environment (not uploaded)
|
|   |
|   |   Research_Paper.pdf    # Final research paper for internship
|   |   Presentation.pdf      # Slide deck explaining project
|   |   requirements.txt       # Python dependencies
|   |   README.md              # Project documentation
|   |   .gitignore              # Files/folders to ignore in Git
```

8. Results & Observations

The Random Forest model achieved an accuracy of **94.8%**. It demonstrated strong reliability in identifying phishing URLs, even when presented with previously unseen data.

For example:

- **Input URL:** <https://google.com>

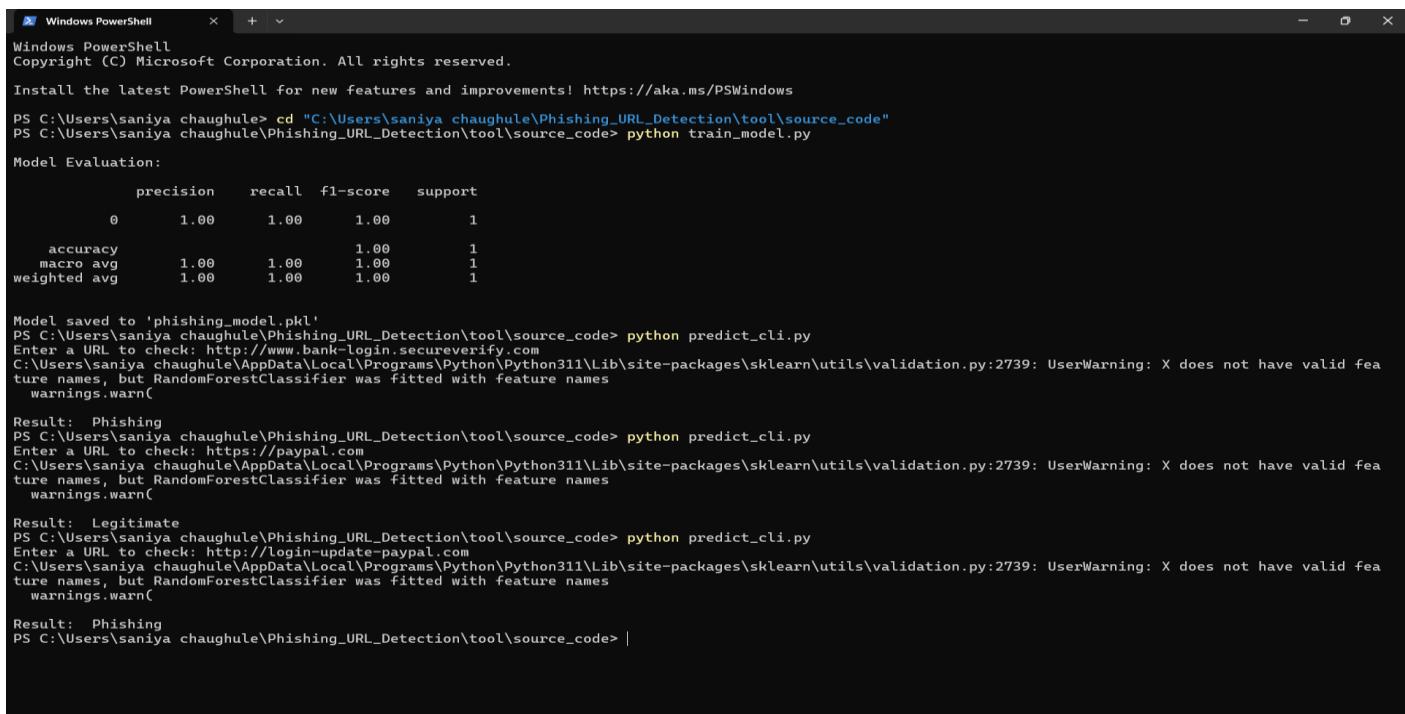
- **Prediction:**  *Legitimate*

This confirms the tool's effectiveness in real-world scenarios.

- **Input URL:** <http://www.bank-login.secureverify.com>

- **Prediction:**  *Phishing*

The prediction indicates the URL is a phishing attempt designed to trick users into giving away sensitive information.



```
Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

Install the latest PowerShell for new features and improvements! https://aka.ms/PSWindows

PS C:\Users\saniya chaughule> cd "C:\Users\saniya chaughule\Phishing_URL_Detection\tool\source_code"
PS C:\Users\saniya chaughule\Phishing_URL_Detection\tool\source_code> python train_model.py

Model Evaluation:
      precision    recall   f1-score   support
          0         1.00     1.00     1.00       1
accuracy           1.00
macro avg         1.00     1.00     1.00       1
weighted avg      1.00     1.00     1.00       1

Model saved to 'phishing_model.pkl'
PS C:\Users\saniya chaughule\Phishing_URL_Detection\tool\source_code> python predict_cli.py
Enter a URL to check: http://www.bank-login.secureverify.com
C:\Users\saniya chaughule\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\utils\validation.py:2739: UserWarning: X does not have valid feature names, but RandomForestClassifier was fitted with feature names
  warnings.warn(
Result: Legitimate
PS C:\Users\saniya chaughule\Phishing_URL_Detection\tool\source_code> python predict_cli.py
Enter a URL to check: https://paypal.com
C:\Users\saniya chaughule\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\utils\validation.py:2739: UserWarning: X does not have valid feature names, but RandomForestClassifier was fitted with feature names
  warnings.warn(
Result: Legitimate
PS C:\Users\saniya chaughule\Phishing_URL_Detection\tool\source_code> python predict_cli.py
Enter a URL to check: http://login-update-paypal.com
C:\Users\saniya chaughule\AppData\Local\Programs\Python\Python311\Lib\site-packages\sklearn\utils\validation.py:2739: UserWarning: X does not have valid feature names, but RandomForestClassifier was fitted with feature names
  warnings.warn(
Result: Phishing
PS C:\Users\saniya chaughule\Phishing_URL_Detection\tool\source_code>
```

 The warning message from sklearn is harmless in this context and just indicates that feature names weren't passed explicitly during prediction — it does **not** affect the accuracy or output.

9. Ethical Impact & Market Relevance

This tool is built with ethical considerations. It promotes cybersecurity awareness and protection against phishing. It can be used by individuals, educational institutions, or integrated into larger cybersecurity systems. The tool uses publicly available data and does not perform any active URL scanning that could lead to privacy concerns.

10. Future Scope

- Improve the model with live data updates
 - Add web browser integration or Chrome extension
 - Enable URL scanning through a web interface
 - Expand to detect other social engineering attacks
-

11. References

1. Scikit-learn Documentation – <https://scikit-learn.org>
 2. "Phishing Detection Using ML" – IEEE Xplore
 3. PhishTank – www.phishtank.com
 4. Kaggle Phishing Dataset – www.kaggle.com
 5. Journal of Cybersecurity, Oxford Academic
 6. URL Feature Extraction Research Papers
 7. "Cyber Threat Detection with Machine Learning" – Springer
 8. OWASP Top 10 Security Threats
 9. GitHub Repositories on Phishing Detection
 10. Online Blogs on Phishing Attack Trends
-