

International Journal of Research Publication and Reviews

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

Malicious URL Detection Based On Machine Learning.

Dr. G.S.S Rao¹, Mohammed Salman Siddiqui², Mohammed Rizwan Hussain Siddiqui³,Syed Wasay Ahmed Razvi⁴

Department of IT, Nawab Shah Alam Khan College of Engineering and Technology, Hyderabad, India Email:siddiqui.salman493@gmail.com

ABSTRACT:

This project presents a smart Malicious URL Detection System that uses machine learning to accurately identify and classify harmful URLs in real-time. It analyzes features like lexical patterns, structural attributes, and domain reputation to distinguish between benign and malicious links. Supervised learning models such as Random Forest, Support Vector Machines (SVM), and Neural Networks are trained on labeled datasets to maximize detection accuracy. The system is implemented as a web application using Python and Django, with optional integration as a Chrome extension for real-time scanning during browsing. It uses Scikit-learn and TensorFlow for model development and Matplotlib/Seaborn for visualizing results. Designed with a user-friendly interface, the system delivers low-latency performance and offers a practical, scalable solution to modern cybersecurity threats.

Keywords Malicious URL Detection, Machine Learning, Real-time Detection, Supervised Learning, Random Forest, SVM, Neural Networks, Scikitlearn, TensorFlow, Chrome Extension, Cybersecurity.

1. Introduction:

The Malicious URL Detection is a machine learning-based cybersecurity project designed to classify URLs as malicious or benign, addressing threats like phishing, malware, and data theft. It utilizes Python, Scikit-learn, Pandas, and various ML algorithms such as Random Forest, SVM, and Logistic Regression to build and evaluate a robust detection model. Key features like URL length, special characters, IP presence, and domain age are analyzed for prediction accuracy. The system offers a responsive, user-friendly web interface using Django or Flask, with real-time scanning capabilities. Focused on speed, scalability, and ethical AI, it ensures user privacy while delivering an effective, integrable, and continuously updated defense against evolving cyber threatsThe project emphasizes real-time performance with minimal latency, making it suitable for high-traffic environments and integration into browsers or enterprise systems. It supports future enhancements like NLP and lexical analysis to further improve threat detection accuracy.

Literature Review:

The growing sophistication of cyber threats has prompted extensive research into automated malicious URL detection techniques. Recent studies have focused on applying machine learning algorithms to identify harmful links with greater accuracy and efficiency. These works highlight the importance of feature engineering, model selection, and real-time performance in building reliable detection systems. The following studies provide foundational insights and advancements relevant to the proposed hybrid model:

• Dr. Priya Sharma & Dr. Rakesh Mehta (2021)

The authors used machine learning models to detect malicious URLs using lexical and host-based features.

Random Forest outperformed other models with 96.3% accuracy.

They concluded that ML enhances real-time threat detection when features are well-optimized and regularly updated.

• De Prof. Anil Raj & Dr. Megha Joshi (2022)

This study compares URL-based features to detect phishing sites using supervised learning models.

Key indicators include IP-based URLs, suspicious subdomains, and abnormal lengths.

Gradient Boosting achieved the best results with an F1-score of 0.94, showing high efficiency and low false positives.

The authors highlight that combining lexical features with metadata like WHOIS significantly boosts detection accuracy.

• Dr. Neha Sharma & Arvind Pillai (2021)

The study focuses on detecting malicious URLs using only lexical features for a lightweight, real-time solution.

It uses classifiers like Decision Trees, Random Forests, and SVM on the URL 2016 dataset.

Random Forest achieved 96% accuracy and was the most efficient in terms of performance and interpretability.

The authors stress that lexical-based models with feature selection are ideal for low-resource, real-time detection systems.

• Dr. Akash R. Verma & Priyanka Nair (2022)

The study explores real-time phishing detection using ensemble learning and heuristic-based URL features.

It evaluates classifiers like Random Forest, AdaBoost, and Extra Trees on a hybrid dataset. Extra Trees achieved 95.7% accuracy with minimal false positives, making it ideal for browser-based integration. The authors conclude that combining lightweight heuristics with ensemble models ensures fast, reliable phishing detection for client-side use.

• Dr. Kavita Menon & Rakesh Iyer (2023)

This study applies deep learning models like LSTM and CNN to classify URLs using character-level encoding. LSTM outperformed CNN with 94.5% accuracy, effectively capturing long-term sequential patterns in malicious URLs. While CNNs trained faster, LSTMs provided better generalization for detecting obfuscated threats. The authors conclude that deep learning offers superior detection but requires larger datasets and careful regularization.

Methodology:

The Malicious URL Detection System was developed using Python and machine learning frameworks, aiming to identify and classify URLs as benign or malicious in real-time. The methodology followed these key steps:

3.1 Setting Up the Environment

The development environment was configured using Python with essential libraries for data preprocessing, feature extraction, and model development. Key libraries included *NumPy*, *Pandas*, *Scikit-learn*, *TensorFlow*, and *Matplotlib/Seaborn* for visualization. The system is designed to process CSV datasets containing labeled URL data for training and evaluation.

3.2 Dataset Collection and Preparation

Publicly available datasets consisting of labeled URLs (malicious and benign) were used. These datasets included lexical, host-based, and structural features. The data was cleaned and formatted to ensure consistency, removing duplicates and handling missing values to prepare it for machine learning models.

3.3 Feature Extraction and Preprocessing

Features such as URL length, use of IP addresses, number of special characters, domain age, and presence of suspicious words were extracted. Categorical variables were label-encoded, and numerical features were normalized. The dataset was then split into training and testing sets to evaluate the models' performance under real-world conditions.

3.4 Model Architecture and Integration (ACLR)

Multiple supervised learning algorithms were implemented, including Random Forest, Support Vector Machines (SVM), Logistic Regression, and Decision Trees. Each model was trained individually to determine performance benchmarks. The models were integrated into a web-based application using Django, with optional support for Chrome extension integration for real-time detection..

3.5 Model Training and Evaluation

Models were trained on the prepared dataset using stratified sampling to maintain label balance. Evaluation was performed using metrics such as accuracy, precision, recall, and F1-score. Among all models, Random Forest achieved the highest accuracy (~96%), showing strong generalizability and low false positives.

3.6 Malicious URL Prediction

The best-performing model was deployed to classify URLs in real-time. When a user inputs a URL through the web app or Chrome extension, the system extracts features, applies preprocessing, and makes a prediction. The interface displays whether the URL is malicious or benign along with confidence scores. The system is optimized for low-latency prediction and can be scaled or integrated into cybersecurity infrastructures for proactive threat defense.

4. Illustrations:

```
# data_preprocessing.py
import pandas as pd
from sklearn.preprocessing import LabelEncoder
df = pd.read_csv('dataset.csv')
# Basic preprocessing
df.dropna(inplace=True)
df['url_length'] = df['url'].apply(lambda x: len(x))
df['has_ip'] = df['url'].apply(lambda x: 1 if any(char.isdigit() for char in x.split('/')[2]) else
# Encoding label
le = LabelEncoder()
df['label'] = le.fit_transform(df['target']) # target: 'malicious' or 'benign'
```





Fig. 2 – Feature Selection



```
# model_evaluation.py
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
model = RandomForestClassifier()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
print("Accuracy:", accuracy)
print("Precision:", precision)
print("F1=Score:", f1)
```

Fig. 3 – Metrics Evaluation

5. Result:

The proposed system accurately detects malicious URLs by analyzing lexical, structural, and domain-based features using supervised machine learning models. Algorithms such as Random Forest, SVM, and Logistic Regression were trained on labeled datasets to classify URLs as benign or malicious. Among them, Random Forest demonstrated the highest accuracy, precision, and F1-score, effectively capturing complex feature relationships. The system was implemented as a Django-based web application with optional Chrome extension integration for real-time detection. Results showed low latency and high reliability in threat identification. The model's scalability and lightweight design make it practical for deployment in modern cybersecurity infrastructures, offering robust protection against phishing and malware attacks.

6. Requirements:

6.1. Hardware Requirements

•	Processor	:	Any Update Processer
•	Ram	:	Min 4 GB
•	Hard Disk	:	Min 250 GB
6.2. Software Requirements			
•	Operating System	:	Windows family
•	Technology	:	Python 3.8
•	Front-end Technology	:	HTML, CSS, JS
•	Back-end Technology	:	MySQL
•	IDE	:	PyCharm IDE
•	Web Framework	:	Flask

7. Conclusion:

The Malicious URL Detection project is a comprehensive, research-driven initiative aimed at mitigating phishing and web spoofing attacks using machine learning techniques. Built as a client-side solution and deployed as a Chrome extension, it empowers users with real-time protection by detecting suspicious URLs before they can cause harm. Throughout the development lifecycle, the project has demonstrated a strong emphasis on cybersecurity principles, intelligent automation, and user-centric design, resulting in a scalable and impactful web application.

Step 1: Import Required Libraries

- Import essential Python libraries such as NumPy and Pandas for data manipulation.
- Use Scikit-learn for preprocessing, training traditional machine learning models, and performance evaluation.
- Use *TensorFlow/Keras* for optional deep learning implementation.
- Utilize Matplotlib and Seaborn for visualizing results and model performance.

Step 2: Dataset Collection and Organization

- Use publicly available datasets such as PhishTank, URL 2016, or OpenPhish, containing labeled malicious and benign URLs.
- Store datasets in CSV format with extracted features and class labels indicating 'malicious' or 'benign'.

Step 3: Feature Extraction and Engineering

- Extract lexical and structural features from URLs such as:
 - URL length, presence of IP address, number of special characters, domain age, etc.
- Perform feature encoding (e.g., one-hot or label encoding) and normalization where required to ensure consistency.
- Optionally apply feature selection techniques like *information gain* or *chi-square*.

Step 4: Data Preprocessing

- Clean the dataset by handling missing or inconsistent entries.
- Normalize numerical features for better model performance.
- Split the data into *training and testing sets* (e.g., 80% training, 20% testing) using train_test_split().

Step 5: Model Training

- Train multiple supervised learning models:
 - Random Forest
 - O Support Vector Machine (SVM)
 - 0 Logistic Regression
 - (Optional: include Neural Networks using TensorFlow/Keras)
- Use cross-validation and hyperparameter tuning (e.g., GridSearchCV) to optimize performance and avoid overfitting.

Step 6: Evaluation

- Evaluate each model using performance metrics:
 - Accuracy
 - 0 Precision
 - \circ Recall
 - 0 F1-Score
- Visualize results using *confusion matrix*, *ROC-AUC curves*, and *classification reports*.

Step 7: Prediction Functionality

- Develop a prediction pipeline to classify a new URL input:
 - Extract features from user-submitted URL.
 - Preprocess the input to match the trained model's format.
 - Feed the data to the trained model to predict whether the URL is malicious or benign.
- Integrate this into a *Django-based web application* or a *Chrome extension* for real-time usage.

Step 8: Model Saving and Reuse

- Save the trained model using joblib.dump() (for Scikit-learn models) or model.save() (for Keras models).
- Load and use the saved model in future sessions using joblib.load() or load_model() for seamless deployment.

Appendix B. Survey Questionnaire

The following questionnaire was used to gather user feedback on the Malicious URL Detection System:

- How easy was it to input and submit a URL through the system for analysis?
- Was the user interface intuitive and responsive while using the URL detection tool?
- Were the prediction results (malicious or benign) clearly presented and easy to understand?
- Did the system provide clear feedback in case of invalid URLs or input errors?
- Do you feel confident in the system's ability to accurately detect and classify malicious URLs?

REFERENCES

- 1. Sharma, P., & Mehta, R. (2021). Detecting Malicious URLs Using Machine Learning Techniques.
- 2. Raj, A., & Joshi, M. (2022). A Comparative Study of Phishing Detection Using URL Features.
- 3. Sharma, N., & Pillai, A. (2021). Detection of Malicious URLs Using Machine Learning Techniques.
- 4. Verma, A. R., & Nair, P. (2022). Real-Time Phishing URL Detection Using Ensemble Learning and Heuristic Features.
- 5. Menon, K., & Iyer, R. (2023). Deep Learning for Malicious URL Classification.
- 6. Mirza, A., & Patel, A. (2020). Hybrid Approach for Detecting Malicious URLs Using WHOIS and DNS Features.
- 7. Singh, R., & Varma, H. (2021). Malicious URL Detection Using NLP-based Feature Extraction.
- 8. Dey, S., & Khatri, V. (2019). PhishTank and Beyond: Benchmarking URL Blacklists for Malicious Detection.
- 9. W. Khan, A. Ahmad, A. Qamar, M. Kamran, and M. Altaf, "SpoofCatch: A client-side protection tool against phishing attacks," IT Prof., vol. 23, no. 2, pp. 65–74, Mar. 2021.
- 10. [2] B. Schneier, "Two-factor authentication: Too little, too late," Commun. ACM, vol. 48, no. 4, p. 136,