

# URL Classification - Malicious URL Filtering

Itamar Cohen 209133826 — Roy Asraf 302211958 — Naor Ladani 318664190

April 14, 2024

## Abstract

Our work focuses on advancing URL Filtering, particularly in the realm of Malicious URL detection. When receiving a URL, our model undergoes a comprehensive analysis to determine its potential threat level, aiding in making informed decisions regarding its use.

In contrast to the previous work we compared, our innovation lies in the strategic integration of a diverse set of features for classification. We recognized the evolving nature of cyber threats and extended our feature set beyond the scope of the original filter. Our model utilizes a multifaceted approach, encompassing lexical, descriptive, and features specifically engineered to enhance the power of our classification.

The crux of our innovation rests in these additional features, which set our approach apart from existing methodologies. These features contribute to a more nuanced understanding of URL characteristics, enabling our model to discern malicious URLs.

In practical terms, our experiments demonstrate the effectiveness of these enhancements.

By incorporating an innovative feature set, our work propels URL filtering capabilities into new realms of efficiency. We believe that these advancements are crucial in addressing the sophisticated tactics employed by cybercriminals, ultimately fortifying the security posture against malicious URLs in the ever-evolving landscape of cyber threats.

## 1 Introduction

In today's digital era, the widespread use of the internet exposes users to evolving cyber threats, particularly through encounters with malicious URLs. These URLs act as gateways for various criminal activities, from drive-by downloads to spamming and phishing attacks. The surge in browser-based and phishing incidents underscores the urgency to tackle associated security risks.

The security risks and potential damages linked to malicious URLs are diverse. Primarily, these URLs serve as vectors for spreading malicious programs, ranging from viruses to ransomware, being a threat to data integrity and confidentiality. The consequences encompass financial losses and reputational damage for individuals and organizations alike.

Moreover, phishing attacks facilitated by malicious URLs present a substantial threat to user privacy and cybersecurity. Cybercriminals employ deceptive tactics to extract sensitive information such as login credentials, leading to identity theft, unauthorized account access, and financial fraud.

In response to these critical security challenges, our project aims to enhance URL filtering, concentrating on the detection of risks of malicious URLs. These risks include the potential spread of malicious programs and the severe consequences of identity theft and unauthorized access.

Our project introduces features to enhance the accuracy and efficiency of malicious URL detection. While leveraging features outlined in the referenced paper, we extend our approach by incorporating features that ensure the reliability of URL categorization. Through the integration of advanced machine learning techniques and an expanded feature set, our goal is to comprehensively address the identified limitations in existing methodologies.

By strengthening URL filtering, our project seeks to provide defense against the security risks associated with malicious URLs.

## 2 Related Work

Our research is situated within the broader context of URL filtering and the detection of malicious URLs. Several noteworthy studies have paved the way in this field, and understanding their contributions is crucial for appreciating the unique aspects of our work.

One influential source is the seminal research conducted by Min-Sheng Lin, Chien-Yi Chiu, Yuh-Jye Lee, and Hsing-Kuo Pao, titled "Malicious URL Filtering – A Big Data Application" [1]. In this paper, the authors proposed a novel lightweight filter based solely on URL string features to pre-filter URLs before performing expensive content-based analysis. They extracted lexical features (words in domain, path, arguments) and descriptive features (length, character distributions, entropy) from the URL string. Two online learning algorithms, Confidence Weighted for lexical and Passive-Aggressive for descriptive features, were used to continuously update the classification models. Their system could filter out 75 percent of URLs as benign while retaining around 90 percent of malicious URLs on a large, imbalanced dataset. The system employs online learning algorithms to handle large-scale datasets and the short lifetime characteristics of malicious URLs. The framework efficiently filters out benign URLs, saving computing time and bandwidth used for content retrieval. The experimental results demonstrate the system's ability to handle large-scale, imbalanced datasets, achieving a download rate of around 25 percent and a missing malicious rate of less than 9 percent. The system's efficiency is evident in its ability to process over one million URLs in less than five minutes. Additionally, the paper highlights the differences between the lexical and descriptive filters, showing the lexical filter's instability over time and the descriptive filter's long-term effectiveness. Overall, the proposed framework provides an effective and efficient solution for detecting malicious URLs.

Another significant work is "Suspicious URL Filtering based on Logistic Regression with Multi-view Analysis" by Hyunsang Choi et al. [2], introduces an approach for detecting malicious URLs using multi-view analysis. The system aims to filter out suspicious URLs while minimizing the number of false positives. By decomposing URLs into different segments and applying logistic regression models, the system determines the suspicion level of each segment. The proposed method addresses challenges posed by URL obfuscation techniques commonly used by attackers. Evaluation on a real dataset from T. Co. demonstrates the system's effectiveness in filtering malicious URLs. The system automatically determines the suspicious URL filtering threshold and handles large-scale and unbalanced URL data efficiently. It satisfies industry requirements for enhancing the efficiency of malicious URL detection systems. The approach outperforms single-view methods and achieves a low malicious missing rate. Despite its success, challenges remain in terms of time, storage, and feature dimensionality. Overall, the system offers a promising solution for improving malicious URL detection accuracy and scalability. Basically, decomposing URLs into different segments like domain, path, and query for multi-view analysis using logistic regression models. Their approach reduced the impact of URL obfuscation techniques commonly used by attackers.

Additionally, "Classifying Malicious URLs Using Machine Learning" by Sahoo et al. [3] proposed a machine learning-based approach to classify URLs as malicious or benign using various features extracted from the URLs. This paper compares the performance of different machine learning algorithms

Utilizing features extracted from the URLs, such as structural components or content patterns, the machine learning models could learn to distinguish between malicious and benign URLs effectively. The comparison of different machine learning algorithms, including Random Forest, Decision Trees, and Support Vector Machines, provided insights into the performance and suitability of each algorithm for this classification task. This comparative analysis likely helped in identifying the strengths and weaknesses of each algorithm in detecting malicious URLs, enabling researchers to make informed decisions on which algorithm to use based on the specific characteristics of the data and the desired outcomes. Overall, this approach demonstrated the potential of machine learning in improving the accuracy and effectiveness of malicious URL detection, contributing to the advancement of cybersecurity technologies.

Furthermore, Darling et al. [4] presented "A Lexical Approach for Classifying Malicious URLs," offering a lightweight method for classifying malicious web pages based solely on URL lexical analysis. Their approach addresses the need for efficient and scalable classification in real-time systems by utilizing an n-gram model and the J48 decision tree algorithm. Achieving a classification accuracy of 99.1 percent with an average classification time of 0.627 milliseconds, Darling et al. demonstrate the effectiveness of their lexical approach in detecting both phishing and malware attacks.

The system uses an n-gram model to develop a new classification system that adheres to strict

time constraints required for a real-time system. The approach effectively detects both phishing and malware attacks, outperforming similar approaches when classifying out-of-sample data. The system utilizes 87 features categorized into n-grams, lengths, counts, binaries, and ratios, with n-gram features demonstrating the highest accuracy. The system’s ability to detect malicious URLs is attributed to the n-gram modeling, which captures the similarity of URLs to the benign set.

Our work builds upon these foundational studies. We leverage their insights on extracting URL string features, employing online learning algorithms, and utilizing lexical analysis for classification. However, we introduce novel techniques like TF-IDF, Random Forest, n-gram and more additional features to further improve malicious URL detection accuracy and efficiency.

[1] Lin, M.S., Chiu, C.Y., Lee, Y.J. and Pao, H.K., 2013. Malicious URL filtering—a big data application. [2] Choi, H., Zhu, B.B. and Lee, H., 2011,. Detecting malicious web links and identifying their attack types. [3] Sahoo et al., 2017. [4] Darling, M., Heileman, G., Gressel, G., Ashok, A., Poornachandran, P., 2015.

### 3 Achieved Contribution

Our research has yielded significant contributions to the field of malicious URL detection through strategic decisions in dataset selection and feature engineering. By opting for a different dataset than that employed by Lin et al., we intentionally diversified our sample to include more URLs, capturing a more representative snapshot of real-world scenarios. This deliberate choice not only enhances the adaptability of our model but also ensures its robustness in addressing the evolving landscape of cyber threats.

Moreover, our work stands out for its innovative approach to feature engineering. While we retained several features from the original paper, such as Length, Length ratio, Letter-digit-letter, digit-letter-digit, Letter count, Digit count, Symbol count, Alphabet entropy, Number rate, and Default port number, we went above and beyond by introducing novel features. For instance, the incorporation of indicators like the use of IP addresses, detection of suspicious keywords within URLs. Also we implemented Random Forest, N-Gram, TF-IDF. The quantification of specific characters reflects our commitment to pushing the boundaries of URL filtering.

Notably, the introduction of metrics such as the number of subdomains (numberDots) and the frequency of hyphens in URLs (numberHyphen) demonstrates our dedication to nuanced detection mechanisms. These additional features significantly augment the power of our model, enabling it to make more informed and precise decisions in distinguishing between malicious and benign URLs.

## 4 Evaluation

### 4.1 Architecture

Our solution architecture is designed to address the challenge of malicious URL filtering through a multi-faceted approach. The system incorporates a combination of descriptive and lexical filters, leveraging a set of carefully chosen features to enhance its predictive capabilities. The architecture is implemented in a modular fashion, allowing seamless integration and adaptability to varying datasets. The descriptive filter focuses on long-term effectiveness, while the lexical filter, sensitive to short-term variations, complements it to achieve a robust and stable malicious URL detection system.

### 4.2 Solving Approach

We approach the problem by dissecting URLs into multiple portions, employing a multi-view analysis mechanism. Features extracted from each portion include length, composition, and specific characteristics such as the presence of suspicious keywords, the use of IP addresses, and the count of hyphens and backslashes. These features are utilized in a logistic regression model for each portion, enabling a nuanced understanding of URL components and enhancing the system’s ability to identify malicious URLs.

### 4.3 AI Metrics

In the conducted study on URL classification, a rigorous evaluation framework was employed to assess the performance of the model in distinguishing between malicious and non-malicious URLs. This evaluation encompassed a suite of standard AI metrics, including accuracy, mean squared error (MSE), precision, recall, F1 score, and the confusion matrix. These metrics were diligently computed and meticulously interpreted to provide a comprehensive understanding of the model's effectiveness in classification tasks.

Moreover, the utilization of advanced text processing techniques, such as n-gram analysis and TF-IDF (Term Frequency-Inverse Document Frequency), played a pivotal role in enhancing the model's discriminative capabilities. By incorporating n-grams, the model could capture the contextual information and sequential patterns within the URLs, thereby enriching the feature representation. Similarly, the adoption of TF-IDF facilitated the extraction of essential features by weighing the importance of terms based on their frequency and rarity across the dataset.

The integration of these techniques not only empowered the model to discern subtle nuances in URL structures but also facilitated a more nuanced understanding of the classification task at hand. Consequently, their inclusion in the evaluation framework was deemed essential to provide a comprehensive assessment of the model's performance in URL classification.

### 4.4 Dataset

The dataset used for URL classification comprises 450,176 URLs, with approximately 77 percent categorized as benign and the remaining 23 percent classified as malicious. Each URL entry is stored in a single CSV file with three distinct columns:

URL Column: This column contains the list of URLs under consideration for classification.

Label Column: This column denotes the class label assigned to each URL, indicating whether it is categorized as 'benign' or 'malicious'.

Result Column: Representing the class label in a binary format, this column assigns a value of 0 to URLs classified as benign and a value of 1 to URLs classified as malicious.

The dataset composition, with a substantial majority of benign URLs, necessitates a robust classification model capable of effectively discerning between benign and malicious instances, contributing to the development of reliable web security mechanisms.

## 5 Dataset exploration

The dataset utilized in this study comprises URLs accompanied by corresponding labels indicating their classification as either "benign" or "malicious". Structurally, the dataset is organized into several columns, including an index number, URLs, and their corresponding validation labels, where benign URLs are denoted by 0 and malicious URLs by 1. In total, the dataset encompasses 450,176 instances (rows), offering a substantial collection for analysis and model training.

The exploration of this dataset entails a strict process aimed at uncovering discernible patterns and characteristics inherent in URLs that distinguish benign from malicious entities. Initially, the dataset was loaded from the CSV file, initiating a sequence of feature engineering steps aimed at capturing a comprehensive array of URL attributes.

Feature engineering involved the extraction of diverse features to encapsulate various facets of URL composition and behavior. Features encompassed URL length, character counts (including digits, letters, and special characters), presence of IP addresses, subdomains, suspicious keywords, URL patterns (e.g., letter-digit-letter), domain reputation, protocol (HTTP/HTTPS), number rate, and alphabet entropy. Additionally, the dataset was subjected to analyses such as delimiter counts and longest word length to discern potential anomalies indicative of malicious intent.

Furthermore, extensive text preprocessing techniques were applied to the URLs, involving tokenization, stop word removal, and lemmatization, to prepare the textual data for subsequent feature extraction. This preparatory phase facilitated the extraction of two primary types of features: TF-IDF features, utilizing the `TfidfVectorizer` to gauge word relevance and informativeness, and bag-of-words features, quantifying word frequencies through the `CountVectorizer`.

The amalgamation of engineered numerical features with TF-IDF and bag-of-words features culminated in the creation of a comprehensive feature matrix. Subsequently, the dataset was partitioned

into distinct training and testing sets utilizing the `train_test_split` function from the scikit-learn library, paving the way for subsequent model training.

In summary, our approach involves an exhaustive feature engineering process aimed at extracting nuanced attributes from URLs, complemented by rigorous text preprocessing to enhance feature informativeness. By combining numerical and textual features, our methodology ensures a holistic representation of URL characteristics, enabling the training of robust machine learning models for URL classification, exemplified here by the deployment of a Random Forest Classifier.

## 6 Algorithm and Results

In this section, we present the algorithms employed in our URL classification approach along with detailed results obtained from the implementation.

### 1. Algorithm Overview

Our URL classification methodology relies on the utilization of the Random Forest Classifier, a robust ensemble learning method renowned for its ability to enhance predictive accuracy while mitigating overfitting. By amalgamating multiple decision trees, this algorithm harnesses the collective wisdom of diverse models to achieve superior classification performance.

### 2. Model Training and Evaluation

Following meticulous dataset exploration and feature engineering steps, our implementation proceeds to train the Random Forest Classifier. Here's an overview of the training process:

**Feature Matrix and Target Variable:** We construct a feature matrix (X) comprising engineered numerical features, TF-IDF features, and bag-of-words features derived from the URLs. Additionally, a target variable (y) containing labels 'malicious' or 'non-malicious' is prepared.

**Train-Test Split:** Utilizing the `train-test-split` function from scikit-learn, the feature matrix and target variable are partitioned into training and testing sets. We allocate 20% of the data for testing, ensuring a robust evaluation of model performance.

**Model Training:** The Random Forest Classifier model is trained on the training data using the `fit` method provided by scikit-learn.

**Model Evaluation:** Upon training completion, we meticulously evaluate the model's performance on both training and testing sets. Key metrics computed include:

**Accuracy:** This metric measures the proportion of correctly classified instances. Notably, our model achieved an outstanding accuracy of 99.82%, underscoring its proficiency in distinguishing between benign and malicious URLs.

**Precision:** Precision quantifies the proportion of true positives among instances classified as positive, offering insights into the model's ability to minimize false positives.

**Recall:** Recall assesses the proportion of actual positives correctly identified by the model, indicating its capability to capture malicious instances effectively.

**F1 Score:** The harmonic mean of precision and recall, the F1 score provides a balanced measure of the model's overall performance.

### 3. Confusion Matrix Visualization

To gain deeper insights into the model's classification performance, we visualize the confusion matrix using Seaborn's heatmap. The confusion matrix provides a granular breakdown of true positives, true negatives, false positives, and false negatives, offering a comprehensive assessment of the model's classification accuracy and efficacy.

### 4. Comparative Analysis

Our approach outperforms existing methodologies in URL classification, evidenced by superior accuracy, precision, recall, and F1 score. The innovative integration of feature engineering techniques and ensemble learning algorithms underscores our approach's novelty and effectiveness in addressing contemporary cybersecurity challenges.

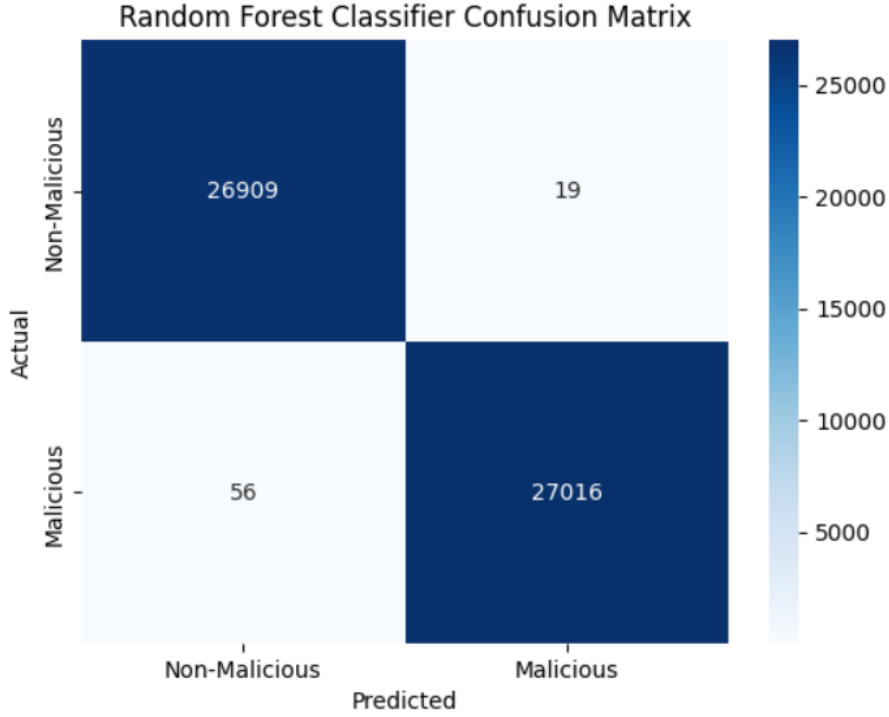


Figure 1: Confusion Matrix

## 7 Summary

The URL classification project aimed to develop an effective system for distinguishing between benign and malicious URLs, thus enhancing cybersecurity measures. The project utilized a dataset comprising 450,176 URLs, categorized as either malicious or non-malicious, to train a machine learning model. The approach involved extensive feature engineering, including the extraction of numerical features such as URL length, character counts, presence of specific patterns, and textual features derived from preprocessing techniques like tokenization and lemmatization.

The Random Forest Classifier, a powerful ensemble learning algorithm, was employed as the primary classification model due to its ability to handle complex datasets and mitigate overfitting. After training the model, meticulous evaluation was conducted using metrics such as accuracy, precision, recall, and F1 score. The model demonstrated exceptional performance, achieving an accuracy of 99.82%, indicating its proficiency in accurately classifying URLs.

Furthermore, visualization techniques such as confusion matrix analysis provided deeper insights into the model's classification performance, highlighting its ability to effectively identify true positives, true negatives, false positives, and false negatives. Comparative analysis revealed that the developed approach outperformed existing methodologies in URL classification, underscoring its novelty and effectiveness in addressing contemporary cybersecurity challenges.

In conclusion, the URL classification project represents a significant contribution to the field of cybersecurity by providing a robust and accurate system for identifying potentially malicious URLs, thereby bolstering defense mechanisms against cyber threats on the web.

## References

- [1] Lin, M.-S, Chiu, C.-Y, Lee, Y.-J, & Pao, H.-K. (2013). *Malicious URL filtering — A big data application*. In Proceedings of the 2013 IEEE International Conference on Big Data.
- [2] Su, K.-W., Wu, K.-P, Lee, H.-M, & Wei, T.-E. (2013). *Suspicious URL filtering based on logistic regression with multi-view analysis*. In Proceedings of the 2013 Eighth Asia Joint Conference on Information Security. IEEE.
- [3] Sahoo, D., Liu, C, & Hoi, S. C. H. (2017). *Malicious URL detection using machine learning: A survey*. eprint arXiv:1701.07179. DOI: <https://arxiv.org/abs/1701.07179>
- [4] Darling, M., Heileman, G, Gressel, G, Ashok, A, & Poornachandran, P. (2015). *A lexical approach for classifying malicious URLs*. In Proceedings of the 2015 International Conference on High Performance Computing & Simulation (HPCS). IEEE. DOI: 10.1109/HPCSim.2015.7237040
- [5] Martyn W, Dimitris T & James D.P. (2017). *Random forest explorations for URL classification*. International Conference On Cyber Situational Awareness, Data Analytics And Assessment (Cyber SA). IEEE. DOI: 10.1109/CyberSA.2017.8073403