

Malicious URL Attack Type Detection using Multiclass Classification

Improving Model Robustness

Orel Zamler ✉

School of Computer science
Ariel University

Anna Pinchuk ✉

School of Computer science
Ariel University

Eilon Barashi ✉

School of Computer science
Ariel University

Aviel Edri ✉

School of Computer science
Ariel University

Abstract

The persistent threat of malicious URLs in the digital landscape necessitates advanced detection strategies beyond traditional blacklists, which falter against newly minted threats. Building on the foundation of using machine learning for the classification of malicious URLs, this study introduces an innovative, integrated approach. We enhance the multiclass classification framework with adversarial training techniques and refined feature engineering, targeting phishing, malware, and spam URLs. Our methodology advances the detection capabilities by incorporating a novel domain similarity score alongside established features like Kullback-Leibler Divergence and bag-of-words segmentation. By deploying a diverse ensemble of classifiers, including Extreme Gradient Boosting (XGBoost), Adaptive Boosting (AdaBoost), Light Gradient Boosting (LightGBM), Categorical Boosting (CatBoost), and introducing Random Forest and adversarial training with ZOO attacks, we demonstrate improvements in model robustness and accuracy. The application of our model within a user-friendly Flask web application illustrates the practicality and accessibility of our approach for real-time malicious URL detection. Evaluated on a dataset of [107,114 URLs], our enhanced model achieves superior performance, markedly improving upon existing methodologies with overall accuracy rates round [0.95]. This work represents a substantial step forward in the application of machine learning to cybersecurity, offering a robust tool for the enhancement of existing anti-phishing, anti-spam, and anti-malware platforms.

1. Introduction

In an era where the Internet has become the backbone of both business and personal communication, the threat posed by malicious URLs to cybersecurity has never been more pronounced. Malicious URLs, acting as conduits for phishing,

malware, and spam attacks, undermine the integrity, confidentiality, and availability of digital resources. Recent reports underscore the escalating trend in cyberattacks, with phishing attempts reaching unprecedented levels and the sophistication of attacks increasing. The Semantic Internet Security Threat Report reveals a surge in web attacks and underscores the vulnerability of small organizations to such threats. Despite efforts to mitigate these risks, traditional approaches like blacklists remain fundamentally flawed, primarily due to their inability to preemptively recognize newly generated malicious URLs. The reliance on blacklists, which catalog known malicious URLs, is increasingly untenable in the face of rapidly evolving cyber threats. Over 90% of malicious URL clicks occur before these URLs are appended to blacklists, highlighting the reactive nature of this approach. Furthermore, the dependency on human intervention for blacklist updates introduces delays and potential inaccuracies. Recognizing these limitations, the cybersecurity research community has pivoted towards machine learning as a proactive and dynamic solution for malicious URL detection. By extracting and analyzing URL features, machine learning models offer a promising avenue for distinguishing between benign and malicious URLs, contingent on the quality and breadth of the data processed. Building on this foundation, our study elevates the detection of malicious URLs to a sophisticated multiclass classification challenge, scrutinizing the efficacy of advanced ensemble learning techniques. We extend the conventional machine learning paradigm by incorporating adversarial training and a refined feature engineering process, focusing on enhancing the model's accuracy and robustness against novel threats. Our approach introduces a novel domain similarity score, augmenting traditional features like Kullback-Leibler Divergence and bag-of-words segmentation.

This study meticulously evaluates the performance of cutting-edge ensemble learners, including Extreme Gradient Boosting (XGBoost), Adaptive Boosting (AdaBoost), Light Gradient Boosting (LightGBM), Categorical Boosting (CatBoost), and Random Forest - with a particular emphasis on the resilience to adversarial attacks facilitated by ZOO attack simulations on XGBoost. We demonstrate the practical application of our findings through the deployment of a Flask-based web application, offering an accessible platform for real-time malicious URL detection. This application showcases the applicability of our research.

2. Motivation

Our research is primarily motivated by the pursuit of a more dynamic, efficient, and forward-looking strategy in the detection of malicious URLs. Drawing inspiration from the foundational work presented in the IEEE research paper “Towards Fighting Cybercrime: Malicious URL Attack Type Detection using Multiclass Classification” we aim to extend the boundaries of malicious URL detection by incorporating advanced machine learning techniques and novel feature engineering methodologies. Recognizing the potential limitations of traditional URL features and the static nature of blacklists, our study introduces an integrated approach that leverages both adversarial training and an optimized set of URL-based features. Adversarial training, a concept not extensively explored in the context of malicious URL detection in previous studies, including the IEEE research paper, offers a promising avenue for enhancing model robustness. By simulating attacks during the training phase, our model is designed to anticipate and counteract evasion techniques employed by attackers, thereby improving the resilience against novel threats. This approach aligns with our goal to not only match but surpass the detection capabilities outlined in existing literature. Furthermore, our research places a significant emphasis on refining the feature selection process. Beyond the conventional use of bag-of-words segmentation and KL divergence, we introduce TF-IDF and a novel domain similarity score, aiming to capture the nuanced distinctions between benign and malicious URLs more effectively. This enhancement is predicated on the hypothesis that a more discerning feature set will lead to improved classification accuracy, particularly in a multi-class setting that encompasses a broader spectrum of URL attack types.

3. Related Work

The landscape of cybersecurity has been significantly shaped by the advent of machine learning algorithms tailored for the detection of malicious URLs. This body of work spans across various attack vectors such as phishing, spam, and malware, each presenting unique challenges and requiring specialized detection methodologies. **Phishing URL Detection:** Phishing attacks, aiming to illicitly acquire personal information, have been combated using machine learning techniques with notable success. Research has leveraged algorithms like Random Forest and SVM, achieving high precision and recall. These studies have focused on the nuances of URL features, domain squatting, and innovative methods like PhishDef for real-time detection, highlighting the critical role of lexical features and the potential of synthetic URL generation for enhancing model training. **Spam URL Detection:** Beyond phishing, spam URL detection has been a focal point, with efforts to differentiate between benign and spam URLs through classifiers like SVM and exploring the potential of NLP and language models. Real-time filtering solutions such as Monarch have demonstrated the feasibility and efficiency of managing large-scale URL traffic with high accuracy. The misuse of short URLs for spam activities has also been a significant area of investigation, revealing the importance of behavioral signals in detection methodologies. **Malware URL Detection:** The detection of malware URLs has seen the application of SVM classifiers using discriminative lexical features, underscoring a consistent pattern in the construction of malware URLs. Innovative approaches, including browser emulation techniques like MineSpider, have been proposed to extract malicious URLs more efficiently, reducing false negatives and enhancing real-time detection capabilities. In addition to these techniques, the Random Forest algorithm emerges as a potent tool in the cybersecurity domain for its robustness against overfitting and its ability to handle high-dimensional data efficiently. Random Forest, by aggregating predictions from multiple decision trees, minimizes the risk of error from individual classifiers, thereby enhancing the overall detection performance. The advent of adversarial attacks presents a new frontier in the challenge of securing ML models against exploitation. These attacks, designed to deceive models through carefully crafted inputs, necessitate the integration of adversarial training techniques. Such methods, including Zeroth Order Optimization (ZOO) attacks, equip models to recognize and counter-

act attempts at evasion, significantly bolstering the resilience of ML systems against the sophisticated tactics employed by cybercriminals. S. C. Jeeva and E. B. Rajsingh, "Intelligent phishing url detection using association rule mining" J. Martinez-Romo and L. Araujo, "Web spam identification through language model analysis" M. Vasek and T. Moore, "Empirical analysis of factors affecting malware url detection" F. Vanhoenshoven, G. Napoles, R. Falcon, K. Vanhoof and M. Koppen, "Detecting malicious urls using machine learning techniques"

4. Data Collection

Our study leverages an extensive dataset of 972,019 URLs, diversified across benign, phishing, malware, and spam categories, sourced from reputable online platforms including Kaggle, the DMOZ open directory project, the Phish-Tank anti-phishing site, the URLhaus malware database, and the WEBSPPAM-UK datasets. This comprehensive collection ensures a robust foundation for developing a high-performing detection model. Dataset Composition: The dataset comprises 773,818 benign URLs, 140,608 phishing URLs, 54,470 malware URLs, and 3,123 spam URLs, highlighting the varied nature of cyber threats present in the digital domain. The data were curated from multiple sources to ensure a rich mix of URL types, enhancing the model's ability to generalize across different forms of malicious activities. To address imbalanced data and computational limitations, we defined specific sizes for each class, preserving all 'spam' samples and randomly sampling from the others to create a balanced training set. The final dataset included 58,132 benign, 14,374 phishing, 31,851 malware, and 3,123 spam URLs, suitable for efficient model training within computational constraints. Train-Test Split: A strategic train-test split of 20% was performed to facilitate model validation and enhance test accuracy.

5. Classification Techniques

Utilizing advanced machine learning algorithms via the Scikit-Learn library, our classification framework employs ensemble methods like XGBoost, AdaBoost, LightGBM, and CatBoost. These algorithms were selected for their efficiency in both binary and multiclass classification scenarios, offering a nuanced approach to malicious URL detection that capitalizes on the strengths of each technique.

XGBoost: A high-performance gradient

boosting framework that excels in computational speed and accuracy, renowned for its effectiveness in competitions. It's capable of handling missing values and offers strong regularization. AdaBoost: Enhances decision trees for binary tasks, developing a series of classifiers from previous outcomes. It's known for adaptive learning in multiclass scenarios, particularly using the SAMME.R algorithm, which leverages real-valued probabilities from weak learners to improve accuracy. This version provided better results in our study, demonstrating its effectiveness in multiclass URL classification alongside other advanced models. LightGBM: A Microsoft-developed, tree-based boosting framework that's fast and computationally efficient, designed to minimize losses more effectively but can overfit on small datasets. CatBoost: A Yandex-developed gradient boosting method on decision trees, known for fast prediction and excellent handling of categorical features, reducing overfitting through its tree structure. Random Forest: An ensemble of decision trees that enhances classification accuracy and prevents overfitting. It's effective for both binary and multiclass tasks, capable of handling large datasets with complex feature sets. Voting Ensemble: Combines predictions from multiple models to improve overall accuracy. It leverages the strengths of individual classifiers, like Random Forest, XGBoost, AdaBoost, LightGBM, and CatBoost, to make more accurate predictions by voting on outcomes for multiclass URL classification. The above mentioned boosting algorithms all have their own pros and cons. Most of the algorithms have been used in malicious URL binary classification tasks. However, we explore their performance in a multi-class URL classification setting. Boosting techniques prove to be an effective machine learning algorithm strategy.

Modeling and Evaluation: Feature Analysis and Performance Metrics

6. Feature Selection and Extraction

In line with methodologies outlined in foundational cybersecurity literature, our study meticulously selected and extracted features critical for URL classification. We identified key features from our dataset, ensuring comprehensive coverage of attributes relevant to malicious URL detection. Our feature extraction process centered on analyzing the characteristics of original URL strings we focused on lexical features, which encompass textual properties extracted

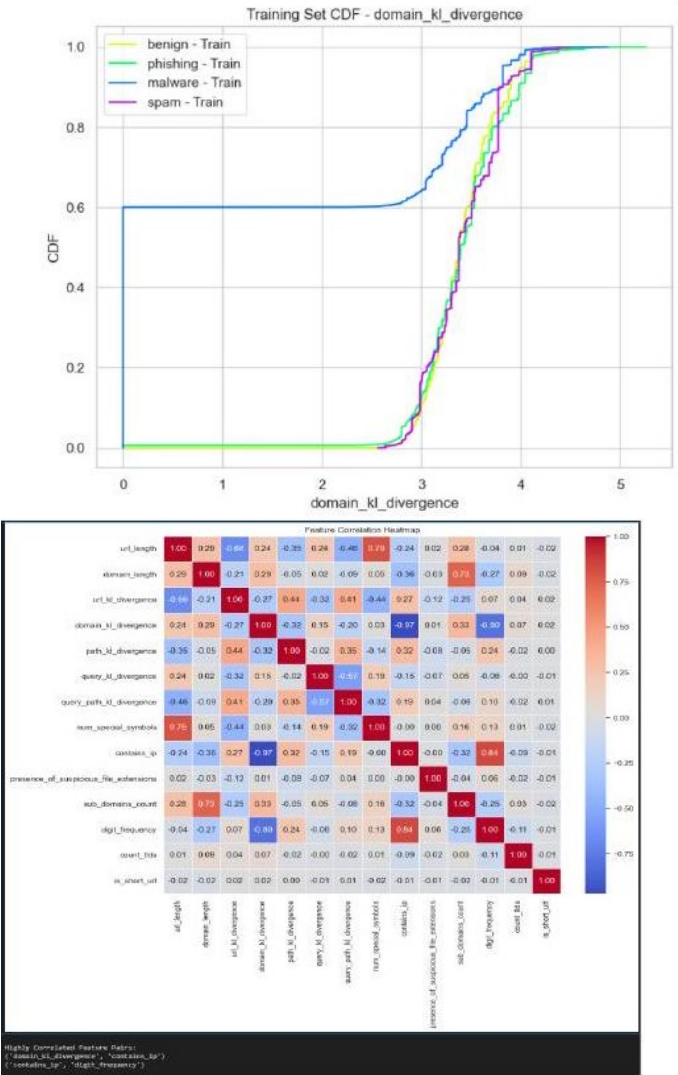
from various URL components, including the protocol, domain name, host name, top-level domain name, path, and parameters. Word-Based Features: We identified suspicious, security-sensitive, suggestive, and obfuscated words commonly found in URL strings. Special Character Features: Analysis of special characters in URLs, including common punctuation and mathematical symbols. KL Divergence Features: KL divergence, or relative entropy, was employed to measure the similarity between the character distributions in malicious and benign URLs, providing insights into their structural differences. File Extension Features: Presence of specific file extensions in URL strings, such as .exe, .js, and .pdf, was examined to identify potentially harmful URLs associated with executable files.

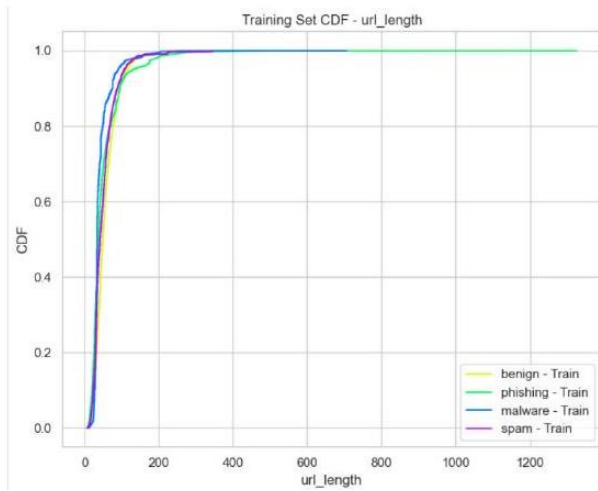
- URL length
- Domain name length
- URL KL divergence
- Special characters count
- Presence of an IP address
- Domain KL divergence
- Path KL divergence
- Query KL divergence
- Query + Path KL divergence
- Sensitive brand-name words
- Segmented bag of words
- Number of signs
- Presence of suspicious file extensions
- Sub domains count
- Frequency of digits
- Short URL features
- Top level domains count

7. Data Exploration and Preprocessing

Our exploration of the dataset unearthed crucial insights, notably the prevalence of URLs containing IP addresses—a potential indicator of malicious intent. An exhaustive examination of missing values guided our preprocessing strategy, ensuring data integrity. Utilizing cumulative distribution function (CDF) graphs, we analyzed feature distributions, gaining a nuanced understanding of our dataset’s characteristics. Correlation between features was visualized through heatmaps, elucidating relationships

and influences on URL classification. This exploration phase was pivotal in identifying patterns and anomalies within our data, informing our feature engineering and model development strategy. Preprocessing was a critical step in preparing our dataset for the modeling phase. We implemented several key preprocessing measures: normalization of data to ensure uniformity, handling of missing values to maintain data quality, and encoding of categorical variables for model compatibility. Special emphasis was placed on numerical features and the application of Bag of Words (BoW) and TF-IDF techniques for textual data, alongside label encoding and standard scaling. This comprehensive preprocessing workflow not only streamlined our dataset for efficient modeling but also enhanced the predictive capability of our classification models.



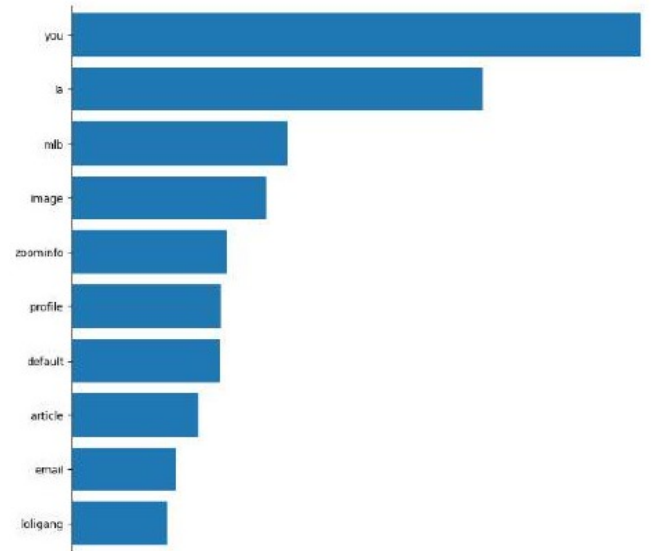


Enhanced Malicious URL Detection through Diverse Ensemble Learning and Feature Exploration

8. Evaluation Metrics

Accuracy: This metric measures the proportion of correctly classified instances out of the total instances. It provides an overall assessment of model correctness. **Error Rate:** Complementary to accuracy, the error rate indicates the proportion of incorrectly classified instances. **Precision:** Precision measures the accuracy of positive predictions, indicating the proportion of correctly predicted positive instances out of all instances predicted as positive. It is calculated as the ratio of true positives to the sum of true positives and false positives. **Recall:** Recall, also known as sensitivity, measures the ability of the model to identify all relevant instances. It is calculated as the ratio of true positives to the sum of true positives and false negatives. **F1 Score:** The F1 score is the harmonic mean of precision and recall, providing a balanced measure of a model's performance. **Feature Importance:** This metric indicates the significance of each feature in influencing the model's predictions. It helps in understanding which features contribute most to the model's performance. **Confusion Matrix (Multi-class):** The confusion matrix provides a detailed breakdown of model performance across different classes. For a multi-class classification problem with classes '0' (benign), '1' (malware), '2' (phishing), and '3' (spam), the confusion matrix displays the counts of true positives, false positives, true negatives, and false negatives for each class.

		True Class		
		A	B	C
Predicted Class	A	5	8	4
	B	1	3	5
	C	9	2	7



$$Accuracy = \frac{T_p + T_n}{T_p + T_n + F_p + F_n}$$

$$Precision = \frac{T_p}{T_p + F_p}$$

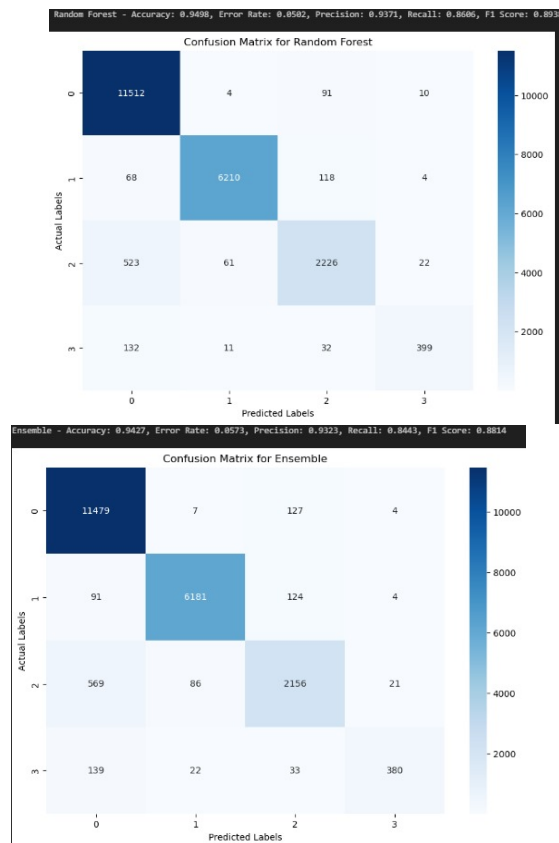
$$Recall = \frac{T_p}{T_p + T_n}$$

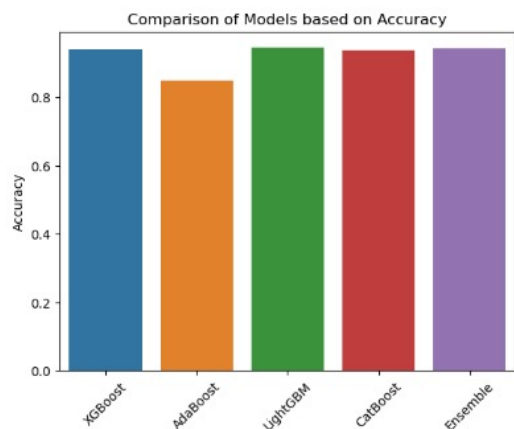
$$F_1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$

9. Modeling Workflow and Evaluations

Our research extends the current understanding of malicious URL detection by implementing and evaluating five distinct models, each with a unique approach to feature extraction and classification. The workflow encompasses several key steps: data preprocessing, feature extraction using both CountVectorizer/TfidfVectorizer and numerical features, model implementation, and model evaluation. Notably, we explored the impact of excluding spam classifications and employing SMOTE with AdaBoost to address class imbalance. Each model—XGBoost, AdaBoost, LightGBM, CatBoost, and a comprehensive Ensemble model—was carefully selected and optimized to enhance detection capabilities. Ensemble Model (IEEE paper): This model combines XGBoost, AdaBoost, LightGBM, and CatBoost through a VotingClassifier, aiming for a balanced approach to malicious URL detection. Using BoW and the whole features without removing any. Ensemble - Accuracy: 0.9440, AdaBoost - Accuracy: 0.8481 (IEEE 95.50%) TF-IDF & Data Exploration Ensemble: Switching to TF-IDF for this ensemble model emphasizes the importance of each term within the dataset, giving higher weight to terms that are rare across documents but prevalent within specific ones. This can be particularly useful for distinguishing malicious URLs. Ensemble - Accuracy: 0.9427, AdaBoost - Accuracy: 0.8502 (improving in AdaBoost) Ensemble Model (Omitting Spam Classification): By excluding spam from the classification, this model focuses on more distinct malicious categories, potentially improving specificity. This adjustment in the dataset leads to better performance metrics, as it simplifies the classification problem for the ensemble, and makes the dataset more balanced. Ensemble - Accuracy: 0.9526 Random Forest Model: Data Exploration and TF-IDF vectorized data aims to leverage the model’s capability for handling high-dimensional data effectively. Random Forest’s inherent randomness and decision tree ensemble approach make it adept at capturing complex patterns without overfitting, particularly valuable given the nuanced differences between various types of malicious URLs. Random Forest - Accuracy: 0.9498 AdaBoost with SMOTE: This model attempts to address class imbalance by applying SMOTE before using AdaBoost for classification. The rationale is to improve AdaBoost’s performance by ensuring it learns from a balanced dataset, given its sensitivity to imbalanced data. However, the results indicate that while SMOTE

can help address imbalance, it may not always lead to improved performance, highlighting the trade-off between addressing class imbalance and maintaining model accuracy. AdaBoost - Accuracy: 0.7270 The choice between BoW and TF-IDF hinges on the specific characteristics of the data and the model’s strengths. BoW is simpler and effective for capturing the presence of tokens, while TF-IDF provides nuanced insights into the importance of each term, potentially offering better performance for models that can leverage this detail. The diversity in modeling approaches and feature extraction techniques underlines the complexity of malicious URL detection. Each model offers unique strengths, with ensemble methods combining these to achieve robust performance. However, the trade-offs between simplicity and complexity, as well as precision and recall, remain central challenges in optimizing for different cybersecurity contexts.

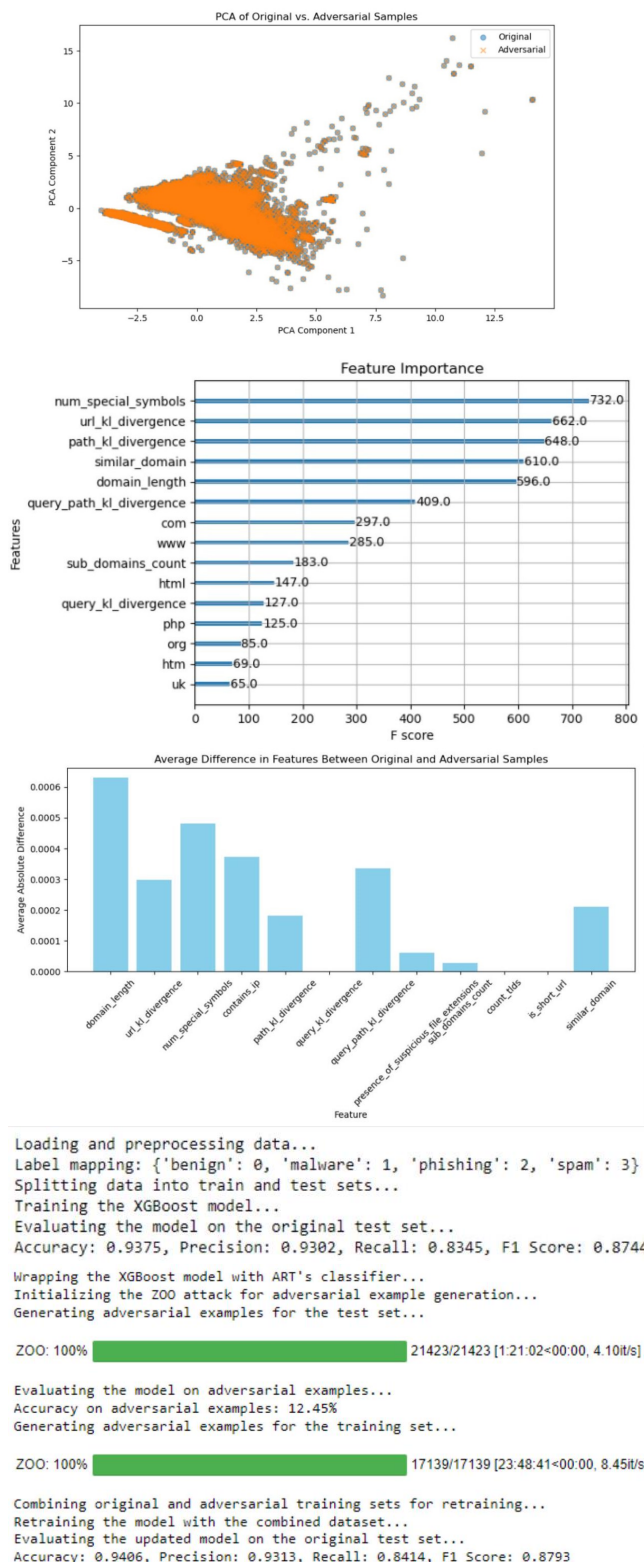




10. Enhancing Machine Learning Security through Adversarial Training: Insights from the Adversarial Robustness Toolbox and Zoo Attack s

With the pervasive integration of machine learning (ML) models in cybersecurity, the robustness of these models against adversarial threats has become paramount. This study explores the application of the Adversarial Robustness Toolbox (ART) and the Zeroth Order Optimization (ZOO) attack to enhance the security of ML models. By generating and defending against adversarial examples, our approach seeks to mitigate the vulnerabilities of ML systems in cybersecurity applications, particularly in the context of URL-based threat detection. The deployment of machine learning models in critical security domains necessitates a rigorous assessment of their vulnerability to adversarial attacks. Adversarial examples, crafted inputs designed to cause model misclassification, pose a significant threat to the integrity of ML systems. This research leverages the Adversarial Robustness Toolbox (ART), a comprehensive library for adversarial machine learning, to evaluate and fortify models against such threats. Specifically, we focus on the Zeroth Order Optimization (ZOO) attack, a sophisticated evasion technique, to probe and enhance the resilience of our ML models. The Adversarial Robustness Toolbox (ART) offers an extensive suite of tools for the evaluation, defense, and verification of machine learning models across various frameworks. Its framework-agnostic nature facilitates the deployment of advanced adversarial techniques in a diverse array of ML applications. By providing a unified interface for implementing both attack and defense mechanisms, ART plays a crucial role in bridging the gap between machine learning research and

practical cybersecurity defenses. The Zeroth Order Optimization (ZOO) Attack The ZOO attack embodies a novel approach to generating adversarial examples without direct access to the target model's gradients. By approximating gradients through finite differences, the ZOO attack enables the crafting of evasion attacks against black-box models. This characteristic is particularly relevant in cybersecurity, where attackers often operate with limited knowledge of the defensive models. The ZOO attack methodically perturbs inputs to find minimal modifications that lead to misclassification, thereby exposing potential vulnerabilities in the model's decision boundaries. Methodology Our methodology integrates the generation of adversarial examples via the ZOO attack within the ART framework to both test and improve the robustness of an XGBoost-based URL classification model. This process involves: Preprocessing and vectorization of URL data to transform textual information into a numerical format suitable for ML models. Training the baseline XGBoost model on the original dataset and evaluating its performance. Utilizing the ZOO attack to generate adversarial examples that simulate potential evasion attempts by malicious actors. Retraining the model on a combined dataset of original and adversarial examples to enhance its adversarial resilience. The initial evaluation of our XGBoost model revealed better accuracy in classifying URLs. However, the introduction of adversarial examples generated by the ZOO attack significantly degraded performance, underscoring the model's vulnerability to evasion attacks. Subsequent retraining with adversarial examples led to an improved model robustness, as evidenced by enhanced performance metrics on both original and adversarial datasets. These findings highlight the efficacy of adversarial training in fortifying machine learning models against sophisticated evasion tactics. This research underscores the critical importance of adversarial robustness in the deployment of machine learning models for cybersecurity applications. Through the strategic use of ART and the ZOO attack, we demonstrate a practical approach to identifying and mitigating the vulnerabilities of ML models to adversarial threats. Future work will focus on exploring a broader range of adversarial techniques and defense mechanisms within the ART framework to further advance the security and reliability of machine learning systems in adversarial environments.



Conclusion

This research significantly advances malicious URL detection by integrating ensemble machine learning models with advanced feature engineering and adversarial training techniques. Our methodology, highlighted by the deployment of a user-friendly Flask web application, demon-

strates not only the technical feasibility but also the practical applicability of real-time URL classification. Achieving an accuracy rate of approximately 0.95, our model markedly improves upon existing methods, providing a robust tool for enhancing cybersecurity measures against phishing, spam, and malware threats.

Future Work

Future directions for this research include exploring a deeper integration of adversarial training methods to further enhance the model's robustness against sophisticated cyber threats. Additionally, expanding the dataset to include a more diverse and balanced selection of recently discovered malicious URLs will improve the model's predictive accuracy and generalizability. Investigating the incorporation of neural network-based models could also offer new insights and potentially higher accuracy in malicious URL detection. In line with enhancing our model's feature set, we plan to incorporate an innovative approach by leveraging the VirusTotal API to enrich our detection capabilities. By integrating VirusTotal's classification score of URLs, our model will benefit from a comprehensive risk assessment derived from various antivirus engines and website scanners. This integration aims to augment our model's predictive power by utilizing external, real-time data on URL reputations, thereby improving its effectiveness in identifying emerging threats. Finally, continuous updates and improvements to the Flask web application, containerized within Docker, will facilitate ease of deployment and scalability, ensuring our solution remains at the forefront of technological advancements in malicious URL detection.

References

IEEE PAPER - T. Manyumwa, P. F. Chapita, H. Wu and S. Ji, "Towards Fighting Cybercrime: Malicious URL Attack Type Detection using Multiclass Classification," 2020 IEEE International Conference on Big Data (Big Data), Atlanta ART - adversarial-robustness-toolbox - <https://github.com/Trusted-AI/adversarial-robustness-toolbox> S. C. Jeeva and E. B. Rajsingh, "Intelligent phishing url detection using association rule mining" J. Martinez-Romo and L. Araujo, "Web spam identification through language model analysis" M. Vasek and T. Moore, "Empirical analysis of factors affecting malware url detection" F. Vanhoenshoven, G. Napoles, R. Falcon, K. Vanhoof and M. Koppen,

”Detecting malicious urls using machine learning techniques”