



T.C.
KARABÜK ÜNİVERSİTESİ
Lisansüstü Eğitim Enstitüsü
Tez Önerisi Savunma Tutanak Formu

Doküman No	UNİKA -FRM-0073
Yayın Tarihi	26.01.2022
Revizyon Tarihi	-
Revizyon No	0

TITLE OF THESIS

An Efficient and Safe Mechanism for Cyber Security Phishing Attacks Based on Machine Learning Algorithms

KEY WORDS:

Phishing Emails, Cyber Security, Machine Learning Classification, Attack Detection

AMAÇ VE HEDEFLER

Tezin amacı ve hedefleri ayrı bölümler halinde kısa ve net cümlelerle ortaya konulmalıdır. Amaç ve hedeflerin belirgin, ölçülebilir, gerçekçi ve tez süresinde ulaşılabilir nitelikte olmasına dikkat edilmelidir.

This research's primary purpose is to design and propose a brand-new mechanism for detecting phishing attempts utilizing a range of features and classification approaches, which will represent a significant enhancement in the field of machine learning algorithms. The following particular strategies will assist us reach our destination:

- Examine both manual and automatic approaches for finding features based on email structure in order to determine which properties of the dataset are most beneficial for recognizing phishing emails. You may then decide which features are most effective for identifying phishing emails.
- Develop a mechanism for discovering phishing emails using several classifications algorithms and evaluate its effectiveness.
- The purpose is to determine which of many classification algorithms better recognizes phishing attempts.

Doküman No	UNİKA -FRM-0073
Yayın Tarihi	26.01.2022
Revizyon Tarihi	-
Revizyon No	0

KONU, KAPSAM ve LİTERATÜR ÖZETİ

Tez önerisinde ele alınan konunun kapsamı ve sınırları, tezin araştırma sorusu veya problemi açık bir şekilde ortaya konulmalı ve ilgili bilim/teknoloji alan(lar)ındaki literatür taraması ve değerlendirilmesi yapılarak tez konusunun literatürdeki önemi, arka planı, bugün gelinen durum, yaşanan sorunlar, eksiklikler, doldurulması gereken boşluklar vb. hususlar açık ve net bir şekilde ortaya konulmalıdır.

Literature Review:

Contributions of this work include a three-stage phishing attack series for content-based phishing attacks that are accurate. The three input characteristics were Uniform Resource Locators, Web Traffic, and Web Content Based on Phishing Attack and Non-Attack Features. To apply the recommended phishing attack approach, a dataset of real phishing attempts is required. Detection of zero-day phishing attacks and phishing attacks in general was demonstrated to be more accurate than previous approaches using actual phishing possibilities. Using three different classifiers, the accuracy of recognizing phishing was assessed, and the findings revealed that NN had a 95.18% accuracy, SVM had an accuracy of 85.45%, and RF had an accuracy of 78.89% [1]. This paper investigates and implements a rule-based strategy for phishing detection utilizing three machine learning models that were trained using a dataset including fourteen (14) features. The machine learning approaches include k-nearest neighbor (KNN), random forest (RF), and support vector machine (SVM). It was shown that of the three applied algorithms, the Random Forest model produced the best results. PhishNet is a Google Chrome extension that includes Random Forest Model principles. Throughout this examination, web technologies like as HTML, CSS, and Javascript are used to construct PhishNet. Therefore, PhishNet enables very efficient web-based phishing detection [2].

In this study, the authors describe a phishing detection model known as PDGAN that runs based just on a URL. Researchers trained a convolutional neural network (CNN) to assess the phishing's features and a long short-term memory network (LSTM) to produce probable phishing URLs in order to identify whether a particular phishing attack is malicious or not. Approximately two million phishing URLs and legal URLs are extracted from PhishTank and DomCop data. The experimental results reveal that the PDGAN beats state-of-the-art models in terms of accuracy and achieves a detection accuracy of 97.58 percent and a precision of 98.02 percent without the need of external services [3]. This study proposes a machine learning attack for identifying phishing assaults. Over 4,000 phishing emails captured by the email system at the University of North Dakota were collected and analyzed by the authors of this research. After selecting the ten most relevant features to these risks, a vast dataset was generated and utilized to develop a model. Using this dataset, machine learning techniques were trained, verified, and evaluated. Probability of detection, miss-detection, false alarm, and accuracy have been employed as performance measurements. The results of the studies indicate that the quality of detection may be enhanced by applying a simulation of a neural network [4].

Intelligent techniques such as Machine Learning (ML) and Deep Learning (DL) are becoming more prevalent in the field of cybersecurity due to their capacity to learn from current data in order to extract pertinent information and predict future events. In this study, researchers examine the efficacy of using such ingenious ways to detect phishing websites. The authors examined two distinct data sets and selected the most strongly linked features, which included a mixture of content-based, URL-lexical, and domain-based features. Then, they implemented many ML models and compared the outcomes. The findings revealed the importance of feature selection for improving the performance of models. In addition, the results tried to discover which features are most useful in leading the model towards the detection of phishing websites. The Random Forest (RF) algorithm has the highest classification accuracy overall [5]. In this study, they describe a model of deep learning that outperforms the most recent developments. Their research, which employs content-based features in substitution of time-consuming text analysis

Doküman No	UNİKA -FRM-0073
Yayın Tarihi	26.01.2022
Revizyon Tarihi	-
Revizyon No	0

methods, is done on three benchmark datasets. Users may use their classifier to filter spam based on one of three major categories. Various performance indicators are used to verify and test models. The durations of the offline training and online detection stages are also specified. The proposed classifier is designed with validation accuracy in mind, with the objectives of achieving faster and competitive performance, which will stimulate its implementation in practical applications [6].

This study provides a unique way of feature selection that integrates the scores of several current approaches to eliminate discrepancies in feature selection results, hence bolstering the credibility of preprocessing. Application to the problem of website phishing classification has shown the benefits and drawbacks of the proposed method for selecting which features are more essential. Features conducted on a security dataset demonstrate that the proposed preprocessing technique successfully generates novel feature datasets that, when mined, provide highly competitive classifiers with respect to detection rate, compared to results obtained using existing feature selection features [7]. In this article, researchers introduce a unique classification strategy for phishing detection based on neural networks. This detection approach is able to achieve high accuracy and robust generalizability by using the concept of risk minimization in design. According to the Monte Carlo technique, the training approach for the unique detection model is straightforward and trustworthy. By comparing this novel phishing detection model to others, such as Naive Bayes (NB), Logistic Regression (LR), K-Nearest Neighbor (KNN), Decision Tree (DT), Linear Support Vector Machine (LSVM), Radial-Basis Support Vector Machine (RSVM), and Linear Support Vector Machine (LSVM), researchers discovered that it outperforms them in terms of Accuracy, True-positive rate (TPR), False-positive rate (FPR), Precision (LDA). Additionally, experimental findings demonstrate that the proposed detection model has a low false positive rate (FPR) of 1.7% and a high Accuracy (AUC) of 97.71%. This gives promising evidence that the proposed detection approach may be used to phishing detection [8].

The objective of this paper is to offer a complete analysis of a few of the most recent solutions for phishing protection. Researchers also present a high-level review of a variety of phishing prevention tactics, such as detection, offensive defense, correction, and prevention, since they believe it's crucial to demonstrate how the different phishing detection methods fit into the greater scheme [9]. Using a mix of Ensemble Learning methods and hybrid features, the authors of this paper propose a system named HELPFED for identifying phishing emails. By combining the content and linguistic properties of email communications, hybrid features provide an accurate representation of the messages. They propose two HELPFED techniques: the Stacking Ensemble Learning approach and the Soft Voting Ensemble Learning technique. Each solution employs two independent Machine Learning algorithms to concurrently process the hybrid data, therefore lowering the complexity of the features and enhancing the performance of the model [10].

The authors provide a novel detection for detecting phishing attacks by combining blacklist-based, online content-based, and heuristic-based strategies with machine learning (ML) algorithms that use a wealth of features. On the basis of Adaptive neuro-fuzzy inference system (AN- FIS), Nave Bayes (NB), PART, J48, and JRip with features, an extensive evaluation was undertaken using evaluation approaches (metrics) to evaluate the performance of the proposed strategy. Every classifier has an accuracy between 99.01% and 99.33%. The fact that PART achieved 99.33% accuracy in less than 0.006 seconds (secs) is a record-breaking accomplishment. Using real-world data, researchers demonstrate that the proposed accuracy can identify phishing websites in real-time and adapt to new phishing techniques [11]. The authors provide an LSTM-based detection for detecting phishing in huge email collections. The new strategy consists of two essential phases: sample increase and sufficient sample testing. In order to guarantee that the training data set is sufficiently big to allow deep learning, they combined KNN and K-means during the phase of sample expansion. Before the testing phase, these samples are first subjected to generalization, word segmentation, and word vector creation. After cleaning and preparing the data, an LSTM model is trained on

Doküman No	UNİKA -FRM-0073
Yayın Tarihi	26.01.2022
Revizyon Tarihi	-
Revizyon No	0

it. Finally, they classify the phishing emails based on the trained model. They undertake trials to determine the accuracy of the proposed approach, and the results indicate that the detection of phishing attempts is 95 percent effective [12]. In this study, the authors provide a phishing email classifier model that use deep learning techniques to recognize fraudulent emails by evaluating the email's body content using a graph convolutional network (GCN) and natural language processing (NLP). GCN has been shown to be successful for text classification in the literature, and this study confirms its efficacy for improving the accuracy of email phishing detection. A supervised learning technique was used to assess the categorization system. Experiments demonstrated that the classifier outperformed existing methods for detecting phishing emails using the body text by generating findings in a fraction of the time with a high accuracy rate (98.2%) and a low false-positive rate (0.015) [13].

In-depth study into what distinguishes legitimate and phishing websites led to the creation of algorithms which could extract 15 features from these pages. Using these heuristic results as input, a machine-learning system was subsequently trained to identify phishing websites. In this approach, two preliminary screening modules were used prior to applying algorithms to the webpages. The user's private whitelist is utilized by the first module, the preapproved site identification, to evaluate whether or not a certain website is safe to view. Based on the existence or absence of login forms, the second module, Login Form Finder, establishes the authenticity of a website. These modules help reduce the amount of needless processing inside the system and the number of false positives without compromising the number of false negatives. When categorizing websites using all of these modules, they achieve an accuracy of 99.8 percent and a false-positive rate of 0.4%. The results of the trial indicate that the method is successful in avoiding identity theft over the Internet [14]. In this study, they suggest a new category of descriptive properties that can be retrieved from the whole email, including the header, the body, and any attachments, in order to enhance the detection of malicious emails using machine learning methods. To meet the needs of real-time detection systems, the indicated features are taken straight from the email itself. All features are thus autonomous, since they do not need an active Internet connection or the use of any other services or tools. Using a dataset of 33,142 emails, of which 38.73% are malicious and 61.27 % are benign, researchers did a comprehensive evaluation of our new novel features in comparison to sets of characteristics provided by previous academic work. These results suggest that, when paired with machine learning approaches, our unique traits may effectively detect malicious emails. In addition, when paired with characteristics revealed in prior research, these novel factors significantly enhance the detection of malicious emails. The Random Forest classifier was the most successful in detecting the detection, with an AUC of 0.929, a TPR of 0.947, and an FPR of 0.03. Moreover, we offer the IDR (integrated detection rate), a new measure for adjusting the threshold of a machine learning classifier to acquire the greatest possible TP and FP rates, the two most important metrics for any operational cyber-security system [15].

The bulk of email classification approaches described in the literature use supervised learning algorithms, which need a substantial amount of labeled data for classification training. Due to the time-consuming and expensive nature of data labeling, the effectiveness of email classification would decrease significantly if just a small amount of data were available. In this article, researchers develop a semi-supervised system for classifying emails based on the idea of multi-view conflict in an effort to relieve this problem. The primary argument is that the literature often disregards the possibility that the multi-view approach may give more information for classification. You may use both labeled and unlabeled data using semi-supervised learning. We evaluate the performance of our technique using two distinct data sets and a real network environment. The experimental results demonstrate that our technique is more effective than several similar algorithms and that multi-view data may be used to generate more precise email classification than single-view data [16]. In this study, the authors propose a novel method for identifying phishing emails using a combination of Collective Learning approaches and hybrid characteristics; they refer to this method as HELPFED. By combining the content and linguistic characteristics of emails, the hybrid characteristics provide an accurate representation of the communications. As possible HELPFED implementations, they propose a Stacking Ensemble

Doküman No	UNİKA -FRM-0073
Yayın Tarihi	26.01.2022
Revizyon Tarihi	-
Revizyon No	0

Learning technique and a Soft Voting Ensemble Learning approach. Both solutions use two independent Machine Learning algorithms to simultaneously process the hybrid features, therefore minimizing the complexity of the features and optimizing the model's performance. To prevent distorted or misleading results, a full evaluation study using cutting-edge methods is done. Experimental investigations indicate that detection performance is enhanced when hybrid features are used with Ensemble Learning as opposed to utilizing just content-based or text-based features [17].

This paper presents the first fully automated method for identifying malicious emails by utilizing deep ensemble learning to analyze the whole email, from text to header to attachments. This eliminates the need for human intervention in the engineering of features. In this study, the authors demonstrate that popular email analysis methods that examine just a subset of the email for analysis may be surpassed by an ensemble framework of deep learning classifiers that have been trained on various email components (thus separately using the full email). A comprehensive study demonstrates that the proposed system outperforms the state-of-the-art methodologies for identifying malicious emails, including human expert feature-based machine learning models, by 5% TPR. The AUC for the new architecture is 0.993% [18]. In this study, researchers provide a model for phishing email categorization that combines deep learning techniques, a graph convolutional network (GCN), and natural language processing to examine the body of an email for indicators of malicious intent. This study demonstrates that GCN may also be used to increase the accuracy of email phishing detection. GCN has been found to be successful in text classification. A supervised learning technique was used to assess the classifier. Compared to previous detection approaches, experimental testing revealed that the classifier could rapidly and reliably identify phishing emails based on their body text with an accuracy rate of 98.2% and a false-positive rate of 0.015 [19]. The authors present a complete study of the most recent machine learning techniques for identifying and filtering spam emails. As part of their research, the authors examine the evolution of spam filtering, the most recent developments, and the future of the subject. In the background portion of the paper, the authors examine how major internet service providers (ISPs) such as Gmail, Yahoo, and Outlook use machine learning techniques to filter out spam emails. It was explained how spam is filtered out of emails in general, as well as the many efforts made by researchers to combat spam using machine learning methods. This paper compares the advantages and disadvantages of existing machine learning methods and highlights outstanding concerns about spam filtering. They recommended deep learning and deep adversarial learning as future techniques that may effectively combat the threat presented by spam emails [20].

This paper presents a novel architecture that combines neural networks with reinforcement learning to detect phishing attempts in real-time. Using the idea of reinforcement learning to dynamically enhance the system over time, the proposed model is able to adapt to the development of a new phishing email detection system that reflects changes in newly investigated behaviors. By adding more emails to the offline dataset in online mode, the proposed model circumvents the problem of a limited dataset. A fresh approach has been offered for investigating any unique phishing activities in the new dataset. Researchers demonstrate that the proposed strategy can handle zero-day phishing attacks with high performance levels, getting high accuracy, TPR, and TNR values of 98.63 percent, 99.07 percent, and 98.19 percent, respectively, following extensive testing using well-known data sets. And both the FPR and the FNR are rather low, at 1.81% and 0.93 %, respectively. When compared using the same dataset as its rivals, the proposed model outperforms previous methods [21]. In this proposal, scientists created an innovative method for improving the accuracy of the Naive Bayes Spam Filter in spotting textual changes and correctly categorizing emails as spam or ham. Their Python-based system increases Spamassassin's Naive Bayes accuracy by integrating semantic, keyword, and machine learning techniques. They also discovered that the length of an email is connected with its spam score, indicating that the disputed Bayesian Poisoning approach is used by spammers [22].

Researchers propose a new taxonomy of features that considers how characteristics are interpreted and used. Next,

Doküman No	UNİKA -FRM-0073
Yayın Tarihi	26.01.2022
Revizyon Tarihi	-
Revizyon No	0

they propose a benchmarking framework, dubbed "PhishBench," that would enable us to systematically and comprehensively evaluate the available features for phishing detection using the same experimental settings, i.e., a unified system specification, datasets, classifiers, and evaluation metrics. PhishBench is the first of its type in the field of phishing benchmarking research since it combines a comprehensive and systematic analysis with a comparison of features. To guarantee the robustness and scalability of the approaches outlined in the phishing literature, we use PhishBench to conduct tests on unique and new datasets. The authors examine how altering the ratio of valid to incorrect records and increasing the size of imbalanced datasets impacts classification accuracy. Their results demonstrate how the asymmetry of phishing attacks affects the effectiveness of detection systems, underlining the necessity for researchers to take this into account while creating fresh techniques. In addition, they discovered that retraining alone is insufficient to defend against new kinds of attack. Defeating attackers that are able to deceive detection systems would need new capabilities and strategies [23]. This paper aims to review and synthesize the available research on the use of natural language processing (NLP) to the challenge of spotting phishing emails. One hundred research articles published between 2006 and 2022 were discovered and assessed using these specified criteria. Researchers investigate the most important topics in the field of natural language processing for detecting phishing emails, such as the most popular machine learning algorithms for phishing email detection, the most prevalent text features found in phishing emails, relevant datasets and resources, and evaluation metrics. According to the data, feature extraction and selection is the major focus of phishing detection research, followed by strategies for classifying phishing emails and boosting their detection. Support vector machines (SVMs) are one of the most often utilized tools for recognizing phishing emails. The most prevalent natural language processing techniques are word embeddings and the TF-IDF approach. Moreover, the Nazario phishing dataset is the most often used benchmarking dataset for phishing email detection systems. Moreover, Python is often used as the main tool for spotting phishing emails. The authors expect that the findings of this paper will be valuable to scholars, especially those focusing on the role of NLP in addressing cybersecurity challenges. The connecting of works to their associated publicly available resources and tools is an additional distinguishing characteristic of this research. Based on the examined literature, it is evident that research into Arabic-written phishing emails that use natural language processing techniques is in its infancy. As a consequence, there are a great deal of inquiries about how to recognize Arabic-written phishing emails [24].

Problem Statement:

Phishing is a method used by identity thieves to get personal user information. To accomplish this objective, fraudulent emails appearing as those from respectable firms are sent. Emails that seem to have originated from a personal source are often employed in phishing attempts, since they provide the impression that the receiver is more motivated to open and click the attachment. Email phishing is the most rapidly expanding part of cybercrime. Their usage allows the theft of the identities and financial information of victims. When individuals react to phishing emails with personal or financial information through pop-up windows, websites, or emails, they put themselves and their companies at danger. There has been a significant amount of study on identifying phishing emails, but no single approach has yet emerged as the obvious winner. This also applies to the underlying categorization procedure, which is viewed as a nondeterministic condition. Lastly, there is an ongoing need to enhance the sensitivity of present detection technologies. In a summary, this research focused on the following significant concerns:

- Improving the performance of the most optimum features and classifiers.
- Which characteristics are most useful for recognizing phishing attempts?
- How to choose a proper classification strategy for phishing detection.

Doküman No	UNİKA -FRM-0073
Yayın Tarihi	26.01.2022
Revizyon Tarihi	-
Revizyon No	0

ÖZGÜN DEĞER

Tez önerisinin, özgün değeri (bilimsel kalitesi, farklılığı ve yeniliği, hangi eksikliği nasıl gidereceği veya hangi soruna nasıl bir çözüm geliştireceği ve/veya ilgili bilim/teknoloji alan(lar)ına metodolojik/kavramsal/kuramsal olarak ne gibi özgün katkılarda bulunacağı vb.) ayrıntılı olarak açıklanmalıdır.

The main goal of this research is to determine if data mining methods can be used to detect phishing attempts. Among the contributions of this thesis to the discussion are the following.

- Select the efficient feature sets for recognizing phishing attempts manually or automatically.
- Evaluate the relative effectiveness of both automatically and manually chosen feature sets.
- Conduct tests to determine the usefulness of your classification system in identifying excellent fishing spots.
- Providing and proposing a detection approach that utilizes many classification types.

YÖNTEM

Tezde uygulanacak yöntem ve araştırma teknikleri (veri toplama araçları ve analiz yöntemleri dahil) ilgili literatüre atıf yapılarak (gerekirse ön çalışma yapılarak) belirgin ve tutarlı bir şekilde ayrıntılı olarak açıklanmalı ve bu yöntem ve tekniklerin tezde öngörülen amaç ve hedeflere ulaşmaya elverişli olduğu ortaya konulmalıdır.

In order to stay up with the ever-changing strategies used by attackers, the feature selection algorithms currently utilized in the literature need additional improvement. In order to increase the accuracy with which phishing emails may be identified and to keep up with the evolution of phisher methods, we recommend the development of a new mechanism to extract new information from raw emails.

We will investigate the machine learning techniques Naïve Bayes, Bagging Decision Tree, Adaptive Boosting used to identify phishing attempts.

1. The Proposed Methodology:

Phishing Attempts may be broken down into three basic groups, which are as follows:

- URL Analysis based approach
- Header Analysis based approach
- Body Content Analysis based approach

2. Architecture of the proposed phishing attack mechanism:

Using these three methodologies and the dataset's feature selection process, we may utilize the newly proposed phishing attack detection mechanism. Below is an illustration of its architecture:

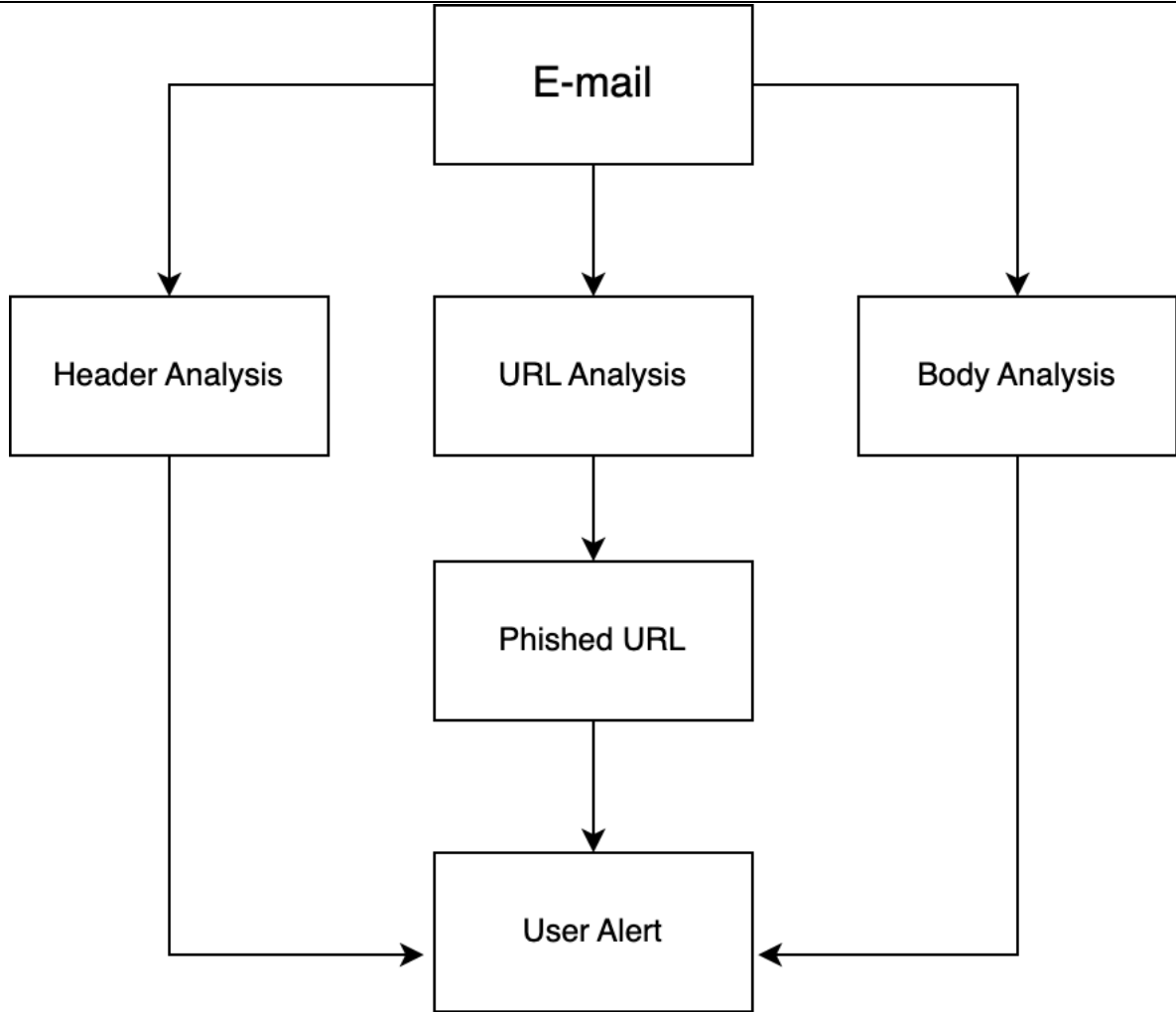


Figure 1 Architecture of the Proposed Phishing Attack Mechanism

3. Machine Learning Classifiers that are used for Phishing Detection

By merging many models into a "ensemble," and using it, the expected performance on a modeling task may be enhanced in comparison to a single model. Machine Learning Classifiers that are used for phishing detection are described here:

- Naïve Bayes

Based on Bayes' Theorem [27], it is a classification technique that employs the idea of predictor independence. A Naive Bayes classifier operates on the premise that the presence of one feature in a class is independent of the presence of other features, for the sake of simplicity. The implementation of the Naive Bayes model needs less data and is advantageous when dealing with enormous datasets. Because of its simplicity and reputation for outperforming more sophisticated methods, Naive Bayes is a commonly used classification technique.

- Bagging Decision Tree

It was first presented by Leo Breiman [26], and it is one of the most effective classifiers for unbalanced data classification. Bagging is comprised of three major components: Using the bootstrapping technique, we first drew random samples from the original training data to build the ensemble. Next, we utilized a decision tree classifier to separately train each of these data.

- Adaptive Boosting

It also called AdaBoost algorithm presented by Yoav Freund and Robert Schapire [25]. All accessible data are

weighted evenly throughout the model-building phase of this approach. Later on, it gives incorrectly classified data points more weight. In the future model, all of the points with higher weights will be given more attention. It will continue training models until the error rate decreases.

4. Features Selection in Email

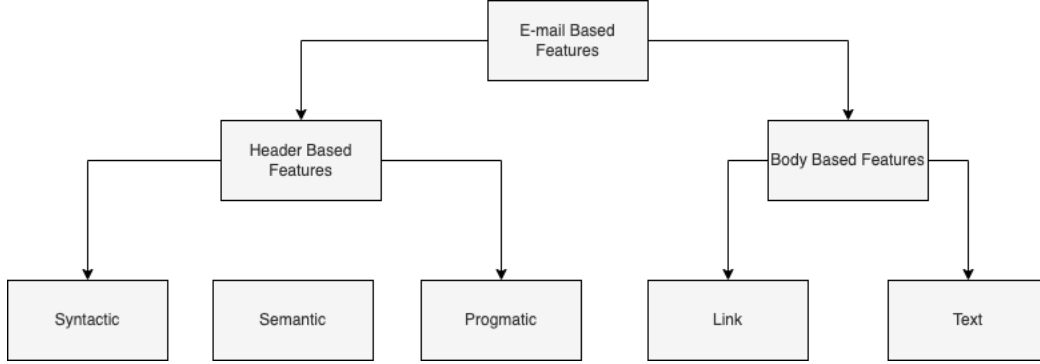


Figure 2 Feature Selection in Email

Data Set:

The employed data set contains a total of 4800 emails, of which 2400 are phishing emails and 2400 are authentic. Valid emails were acquired from the spam Assassin website for the data mining competition, while phishing emails were collected from the monkey website (Monkey, 2022). Both the "monkey website" (known for sending phishing emails) and the "spam Assassin" website were responsible for the emails (Apacheorg, 2022).

Spam Assassin is a website that provides users with a database of valid email addresses. There are two categories of valid emails: those that are straightforward to detect and those that are more likely to be deemed spam.

Using feature extraction, each email is transformed into an array of 47 carefully selected features, which are then utilized to define the data set. Samples of Dataset feature selection shown in figures below.

Y	Z	AA	AB	AC	AD
Subject fwd Word	Subject Num Chars	Subject Num words	Subject Reply word	Subject Richness	Subject Verify word
0	21	4	1	0.19047619	0
0	21	4	1	0.19047619	0
0	37	6	0	0.162162162	0
0	21	3	0	0.142857143	0
0	51	7	0	0.137254902	0

Figure 3 Sample of the data set features of email Subject

A	B	C	D	E	F	G
Email Number	body Dear Word	Body Form	Body HTML	Body Multipart	Body NumberChart	Body Num Function Words
1	0	0	0	0	4522	0
2	0	0	0	0	890	1
	0	0	0	0	3931	12
	0	0	1	0	2995	19
4801	1	0	1	0	1382	15

H	I	J	K	L
Body Num Uniq Words	Body Num words	Body Richness	Body Suspension Word	Body verify your account phrase
374	931	0.205882353	0	0
124	198	0.22247191	0	0
499	864	0.219791402	0	0
195	642	0.214357262	1	0
164	327	0.236613603	0	0

Figure 4 Sample of the data set features of email Body

Tools:

Google Colab - What is Google Colab?

Colab created by Google's corporate research team. It is perfect for machine learning, data analysis, and education since it allows anybody to build and execute arbitrary Python code in the web browser.

Colab's pre-installed machine learning libraries, such as TensorFlow, PyTorch, and Keras, will be used to develop the recommended mechanism. It will be useful for working with data libraries such as NumPy, Matplotlib, and Pandas since it is a hosted Jupyter notebook service that requires no installation and provides users with free access to computer resources such as GPUs.

ÇALIŞMA TAKVİMİ

Tezin başlıca iş paketleri ve tahmini süreleri verilen İş-Zaman Çizelgesinde verilmelidir. (Aşağıda bir örnek verilmiştir.)

İş Adı/Tanımı	AYLAR																																				
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	
Survey and preliminary study																																					
Definitions and identification of the features selection and classification techniques																																					
Performance evaluation using theoretical modeling of detection algorithms																																					
Simulation/Hardware investigation of the proposed detection model																																					
Research and Develop of Phishing Attack Method																																					
Final testing and comparison studies																																					
Thesis Writing and Publications																																					

KAYNAKLAR

1. Mohamed, G.; Visumathi, J.; Mahdal, M.; Anand, J.; Elangovan, M. An Effective and Secure Mechanism for Phishing Attacks Using a Machine Learning Approach. *Processes* 2022, 10, 1356. <https://doi.org/10.3390/pr10071356>.
2. T.O. Ojewumi, G.O. Ogunleye, B.O. Oguntunde, O. Folorunsho, S.G. Fashoto, N. Ogbu, "Performance evaluation of machine learning tools for detection of phishing attacks on web pages, Scientific African", Volume 16, 2022, ISSN 2468-2276, <https://doi.org/10.1016/j.sciaf.2022.e01165>.
3. S. Al-Ahmadi, A. Alotaibi and O. Alsaleh, "PDGAN: Phishing Detection With Generative Adversarial Networks," in IEEE Access, vol. 10, pp. 42459-42468, 2022, doi: 10.1109/ACCESS.2022.3168235.
4. Salahdine, Fatima, Zakaria El Mrabet, and Naima Kaabouch. "Phishing Attacks Detection A Machine Learning-Based Approach." 2021 IEEE 12th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON). IEEE, 2021.
5. M. Aljabri and S. Mirza, "Phishing Attacks Detection using Machine Learning and Deep Learning Models," 2022 7th International Conference on Data Science and Machine Learning Applications

Doküman No	UNİKA -FRM-0073
Yayın Tarihi	26.01.2022
Revizyon Tarihi	-
Revizyon No	0

- (CDMA), 2022, pp. 175-180, doi: 10.1109/CDMA54072.2022.00034.
6. Magdy, S.; Abouelseoud, Y.; Mikhail, M. Efficient spam and phishing emails filtering based on deep learning. *Comput. Netw.* 2022, 206, 108826.
 7. Rajab, K.D. New Hybrid Features Selection Method: A Case Study on Websites Phishing. *Secur. Commun. Netw.* 2017, 2017, 1–10.
 8. Feng, F.; Zhou, Q.; Shen, Z.; Yang, X.; Han, L.; Wang, J. The application of a novel neural network in the detection of phishing websites. *J. Ambient Intell. Humaniz. Comput.* 2018, 1–15.
 9. Khonji, M.; Iraqi, Y.; Jones, A. Phishing detection: A literature survey. *IEEE Commun. Surv. Tutor.* 2013, 15, 2091–2121.
 10. Bountakas, Panagiotis, and Christos Xenakis. "HELPHED: Hybrid Ensemble Learning PHishing Email Detection." *Journal of Network and Computer Applications* (2022): 103545.
 11. Barraclough, Phoebe A., Gerhard Fehrer, and John Woodward. "Intelligent cyber-phishing detection for online." *Computers & Security* 104 (2021): 102123.
 12. Li, Qi, et al. "LSTM based phishing detection for big email data." *IEEE transactions on big data* (2020).
 13. Alhogail, Areej, and Afrah Alsabih. "Applying machine learning and natural language processing to detect phishing email." *Computers & Security* 110 (2021): 102414.
 14. Gowtham, R., and Ilango Krishnamurthi. "A comprehensive and efficacious architecture for detecting phishing webpages." *Computers & Security* 40 (2014): 23-37.
 15. Cohen, Aviad, Nir Nissim, and Yuval Elovici. "Novel set of general descriptive features for enhanced detection of malicious emails using machine learning methods." *Expert Systems with Applications* 110 (2018): 143-169.
 16. Li, Wenjuan, et al. "Design of multi-view based email classification for IoT systems via semi-supervised learning." *Journal of Network and Computer Applications* 128 (2019): 56-63.
 17. Bountakas, Panagiotis, and Christos Xenakis. "HELPHED: Hybrid Ensemble Learning PHishing Email Detection." *Journal of Network and Computer Applications* (2022): 103545.
 18. Muralidharan, Trivikram, and Nir Nissim. "Improving malicious email detection through novel designated deep-learning architectures utilizing entire email." *Neural Networks* 157 (2023): 257-279.
 19. Alhogail, Areej, and Afrah Alsabih. "Applying machine learning and natural language processing to detect phishing email." *Computers & Security* 110 (2021): 102414.
 20. Dada, E. G., et al. "Machine learning for email spam filtering: review, approaches and open research problems. *Heliyon*, 5 (6), e01802." (2019).
 21. Smadi, Sami, Nauman Aslam, and Li Zhang. "Detection of online phishing email using dynamic evolving neural network based on reinforcement learning." *Decision Support Systems* 107 (2018): 88-102.
 22. N. Mageshkumar, A. Vijayaraj, N. Arunpriya, A. Sangeetha, Efficient spam filtering through intelligent text modification detection using machine learning, *Materials Today: Proceedings*, Volume 64, Part 1, 2022, Pages 848-858, ISSN 2214-7853.
 23. El Aassal, Ayman, et al. "An in-depth benchmarking and evaluation of phishing detection research for security needs." *IEEE Access* 8 (2020): 22170-22192.
 24. Salloum, Said, et al. "A systematic literature review on phishing email detection using natural language processing techniques." *IEEE Access* (2022).
 25. Breiman, Leo. "Bagging predictors." *Machine learning* 24.2 (1996): 123-140.
 26. Freund, Yoav, and Robert E. Schapire. "A decision-theoretic generalization of on-line learning and an application to boosting." *Journal of computer and system sciences* 55.1 (1997): 119-139.
 27. Joyce, James (2003), "Bayes' Theorem", in Zalta, Edward N. (ed.), *The Stanford Encyclopedia of Philosophy* (Spring 2019 ed.), Metaphysics Research Lab, Stanford University, retrieved 2020-01-17.