

Deep Learning Approach in Anti-Phishing Technology

Anju Gopi¹, Sajin R Nair²

^{1,2,} Assistant Professor, Malla Reddy College of engineering, Telengana, India

Abstract— A significant rise in electronic trade, or consumer-to-consumer online transactions, has occurred recently as a result of advances in Internet and cloud technologies. Malware, phishing, ransomware, manin-the-middle attacks and other techniques are just a few of the ways cybercriminals utilize to begin a cyber attack. Phishing is a form of social engineering or a cyber security assault in which nefarious actors send messages acting as reliable individuals. **Phishing attacks** use complex strategies, techniques, and technologies to obtain sensitive content injection. data. including online social networks etc. Real-time identification of phishing sites is challenging with the available anti-phishing tools. Therefore, a protected and open system is required that allows the user to directly verify the site. The long shortterm memory (LSTM) algorithm and the convolutional neural network (CNN) were combined to develop a hybrid deep learning model (LCR). The goal of this research was to develop a web browser named Phishing Detection Browser (PD Browser) that can identify phishing attacks using deep learning.

Keywords— Phishing, Anti-Phishing, Deep learning, Hybrid Deep Learning Model

I. INTRODUCTION

Phishing is a term that is frequently used in traditional media. social media, and academic publications nowadays. The term "phishing" is defined by many researchers in different ways and several different definitions can be found in the literature. In a Usenet newsgroup called AOHell on January 2, 1996, a group of hackers first used the term "phishing" to describe their theft of users' login information from America Online (AOL). Since then, the sophistication and scale of phishing attacks have grown, causing significant financial and reputational harm to online users. Phishing is characterized as a fraudulent practise used to get sensitive user data, including passwords, credit card numbers, and login credentials. Typically, it is accomplished through the use of email or other electronic contact while masquerading as a trustworthy company entity. Phishing websites frequently share a common set of objectives; they are created to steal or collect sensitive information from a target. This frequently takes the form of credit card information theft or credential harvesting. Usually, phishing websites and emails are used in conjunction with one another to achieve these objectives [1].

Anti-phishing describes measures taken to thwart phishing attempts. Anti-phishing software consists of computer programs that make an effort to recognize phishing content in emails, websites, and other kinds of data access and block it, typically with a warning to the user. In addition to enhancing

awareness and promoting the use of anti-phishing toolbars that are designed to prevent users from accessing phishing web pages where their sensitive information would be requested and then transmitted to criminals, security professionals are working to lessen the impact of phishing. In order to counteract phishing attempts, a variety of strategies have been used up to this point, but their accuracy is poor. Solutions for detecting phishing frequently rely on dynamic black lists or supervised machine learning models that have been trained on datasets that contain actual data. Some phishing detection models (for training and prediction) rely solely on URLs, whereas others also employ HTML contents to extract feature data. Three main drawbacks exist for them: First, due to biased phishing datasets used for training. Second, a constant need for big, labelled phishing datasets to keep the model current and third, the absence of an explanation for the projected outcomes, biased models are produced. Moreover, keep in mind that incredibly diverse HTML scripts might be used to produce webpages with identical visuals. Technical difficulties arise as a result, reducing the precision of detection when attempting to infer the visual semantics of webpages. Attackers can very readily use evasion strategies to trick such solutions. As a result, the goal of this research was to provide a solution that would make it possible to detect phishing attempts more accurately and quickly as well as raise awareness among active Internet users about how they may defend themselves.

The deep learning algorithm has been applied in a few of these systems to deal with this issue. In recent years, a type of machine learning called deep learning has developed as a potentially effective method for phishing detection. This algorithm falls under the category of an unsupervised machine learning algorithm since it makes its own discoveries from existing data before applying them to new data so this kind of algorithm has a lot of potential for detecting things. Additionally, keep in mind that sites with identical aesthetics could be created using a wide variety of HTML programs. Due to technical challenges, attempts to infer the visual semantics of webpages suffer from decreased detection precision. Attackers can easily deceive such solutions by using evasive techniques. In order to create a powerful approach for phishing website detection, this study combines the long shortterm memory (LSTM), recurrent neural network (RNN), and convolutional neural network (CNN) algorithms.

LSTMs are often used to learn from, analyses, and classify sequential data because of their capacity to comprehend longterm connections between data time steps. Language modelling, speech recognition, sentiment analysis, and video analysis are examples of common LSTM applications.

MINTERNATIONAL JOURNAL OF SCIENTIFIC RESEARCH IN ENGINEERING AND MANAGEMENT (IJSREM)

VOLUME: 07 ISSUE: 09 | SEPTEMBER - 2023

SJIF RATING: 8.176

ISSN: 2582-3930

Recurrent Neural Networks are a Deep Learning method for modelling sequential data (RNN). Prior to the development of attention models, RNNs were the go-to recommendation for handling sequential data. A deep feed forward model can need particular parameters for each component of the sequence.

RNN extensions that expand the memory are known as long short-term memory (LSTM) networks. Building pieces for an RNN's layers are called LSTM. By giving data "weights," LSTMs enable RNNs to accept new information, ignore it, or give it sufficient relevance to affect the result. CNN is a popular and efficient pattern detection and image processing approach. It offers several attributes, including adaptability, a straightforward structure, and fewer training requirements. It is currently a prominent issue in image and voice recognition[1][2].

II. TYPES OF PHISHING ATTACKS

Through phishing emails, many online criminals disseminate harmful attachments and links to lure unwary users into downloading malware. There are other kinds of phishing, while email-based fraud is the most common type [3].

1. Email Phishing

The most common phishing method is email. Scammers create bogus domains that mimic legitimate businesses and bombard their targets with a large number of requests. Character substitutions, such as placing "r" and "n" next to each other to create "rn" in place of "m" are frequently used in fake domain names. Additionally, they might use a legitimate company name in the local section of an email address with the sender's name visible in the inbox.

2. Spear Phishing

Similar to other phishing attempts, spear phishing involves deceiving targets by sending messages from an official-looking source. A spear phishing assault, on the other hand, targets a particular person or group of people rather than sending out generic communications to a large number of users in the hopes that one of them will fall for the ruse. Because they have greater access to the larger firm, HR personnel and IT managers are often targets. Whaling is the term used when the goal is exceptionally lofty. While whaling targets high-value people like the CEO, standard spear phishing targets IT or management figures). Attackers frequently use their ability to pass for other top executives or employees of different businesses to persuade their target to reveal confidential and priceless information.

3. Vishing and Smishing

Smishing (SMS phishing) and Vishing use mobile devices instead of email (voice phishing). Smishing is the practise of attackers sending text messages with misleading information akin to phishing emails. Vishing entails phone conversations during which the victim is spoken to directly by the con artist. In one well-known vishing scam, the con artist poses as a fraud investigator working on behalf of a bank or credit card company. The perpetrator notifies the victims of an account breach and requests their credit card information to confirm their identity. Alternately, the attacker can demand that the victim send money to a particular account.

III. RELATED WORKS

There have been many phishing detection models established throughout this.

The web service is one of the crucial software services for Internet communications. Web phishing is one of the security hazards to web services on the Internet. In this study focused primarily on developing a deep learning infrastructure to detect phishing websites. The study first divided the components of web phishing into two groups: special features and interactive features. A Deep Belief Network-based detection model makes use of these traits (DBN). The DBNbased detection algorithm produced encouraging results in tests utilising actual IP streams [4].

In order to more accurately capture the core issue and enhance classifier performance in recognizing fraudulent URLs, a combination of linear and nonlinear domain conversion algorithms were developed. As a two-step distance measure learning method, they applied the singular value decomposition technique to linear transformation in order to create a perpendicular space and linear programming in order to solve an optimal distance measure. The Nystrom method for the kernel approach was suggested for nonlinear transformation [5].

Attacks were divided into the following categories based on mobile phishing protection strategies: Phishing via social engineering (e.g., SMS, VoIP, website, e-mail), phishing via mobile application (e.g., similarity attack, notification attack, floating attack, forwarding attack), phishing via malware (e.g., ransomware, botnet, key loggers), phishing via online social network (e.g., malicious link, fake profile), phishing via content injection (e.g., cross-site scripting attack), phishing In addition, they explored text mining, URL-based, mobile Quick Response Code (QR-code), machine learning, optical character recognition, and blacklist-based anti-phishing techniques [6].

The current method for detecting phishing websites locates the site using a blacklists/whitelists strategy. The available detection techniques cannot determine with accuracy whether the website is a phishing site or not a phishing site. The system makes use of the device that utilising deep learning and learning algorithms to train the phishing website on the computer system. These two methods are employed to categories the URL using the training dataset and determine whether the website is legitimate or a phishing site. The usage of the deep learning algorithm increases the precision of the forecast. Study conducted between the Deep learning and machine learning algorithms, research has shown that deep learning algorithms provide high accuracy that compared to the random forest approach. The deep learning algorithm provides greater accuracy. However, the detection technique is

VOLUME: 07 ISSUE: 09 | SEPTEMBER - 2023 SJIF RATING: 8.176 ISSN: 2582-3930

not 100% accurate, and it is unable to identify phishing websites with high accuracy [7].

Based on the experiment, it can be concluded that CANTINA is effective at identifying phishing sites and that about 95% of them are appropriately identified. Based on the content of phishing websites, a framework of file matching algorithms is constructed using a unique data set of 17,684 phishing attempts targeted at 159 different brands. According to the findings of trials performed on various algorithms, some phishing detection techniques can yield a detection rate of more than 90% [8].

An effective method [9] using a single-layer neural network to identify phishing websites. The suggested method, in particular, determines the worth of heuristics objectively. Then, a single-layer neural network produces the heuristic weights. The proposed method is assessed using a data set of 10,000 real websites and 11,660 fraudulent websites. The findings demonstrate that the approach can identify more than 98% of phishing websites.

IV. EXISITING SYSTEM

The Existing system detection mechanism for phishing attempts is described in this section. When the victim opens the URL supplied to them via email, a phishing assault takes place and they are taken to a phishing website that looks exactly like the real one. In this detection technique, we merely pay attention to the URL to distinguish between authentic and fraudulent websites. The random forest machine learning algorithm and the deep learning algorithm are used to process the URL.



Fig.1. Existing System Flow Chart

The classification of URLs is based on a number of factors, including IP address, URLs with "@" or "//" symbols, extended URL addresses, and URLs with prefixes or suffixes. In order to identify the URLs of phishing websites from these criteria, the random forest algorithm and deep learning algorithm are used. Using data from the UCI machine repository, the random forest algorithm and deep learning algorithm are utilized to train the detection system. The system using trained and provided input datasets if the system works as intended, a new dataset is used for testing. Then the detection system anticipates the phishing website can halt the training if the detecting system is not otherwise only after a positive evaluation outcome will the trainees be given additional training. Training has ended, and the detecting system is prepared to identify phishing websites [7].

The detecting system gains knowledge using machine learning and deep learning algorithms, determine information about the URL or phishing website. The detection system is prepared to identify the phishing website by detecting the URL when the training procedure is over. The user's input is collected by the detecting system, and evaluates the trained detection against the provided input system and determines whether a given input URL is a trustworthy website or a scam site.

V. PROPOSED SYSTEM

Deep learning algorithms can adapt to the characteristics of the data they are trained on by "thinking" like a human brain using deep neural networks. Because of this, it can more easily adjust automatically to the numerous risks that exist [11].

Within this research, a substantial amount of background knowledge, experience, and associated facts about phishing were gained. When creating phishing model classifiers, the utilisation of high-quality datasets in phishing detection classification is crucial. Various pieces of literature were examined for this study, and they all revealed a wide range of phishing prevention methods. To get statistical findings, a quantitative approach was also applied. These efforts led us to create the Phishing Detection Browser (PD Browser).

INTERNATIONAL JOURNAL OF SCIENTIFIC RESEARCH IN ENGINEERING AND MANAGEMENT (IJSREM)

OLUME: 07 ISSUE: 09 | SEPTEMBER - 2023

SJIF RATING: 8.176

ISSN: 2582-3930





In our method, a phishing website detection system based on LSTM and the CNN was implemented using a hybrid deep learning algorithm. In order to better reliably identify phishing websites, the CNN and LSTM were integrated in this research to detect a wide range of website components. The text and frame content of online pages were detected using the LSTM algorithms, while the website's image features were examined using the CNN algorithm. The Fig. 2 illustrates the flowchart of proposed system. It briefly outlines the conceptual framework of action underlying the suggested system. Firstly, involves collecting the URL that will be used to determine whether it has been phished or not. Our own web browser, the PD Browser, connects to the server. Secondly an evaluation of the URL's data variables is employed the LCR-Hybrid Deep Learning model. Afterward, a result is produced demonstrating the legitimacy of looked up URL If the result indicates that the URL is phished; the user can decide how to avoid it. If the result indicates that it is not phished, browsing may continue. If it is decided not to limit access also generate a Pop up Message "Phishing Website", users may still view the web data on that resource locator (the system data are open to hackers). If the user chooses to block rather than mark something as spam, control is transferred to the system properties (the system data are safe). It should be noted that access to the banned site is prohibited not only from our browser, but also from other browsers, and this is the main benefit of the suggested solution.

The goal of this work is to detect the status of a given URL using minimal distinctive features with deep learning classifiers. Comprises feature extraction, feature selection and classification methodologies. A set of webpage URLs are fed as an input into the feature extractor, which extracts required features from three sources (URL obfuscation, hyperlink and third-party- based). The outcome of the algorithm helps in selecting the best performing features by a clear investigation in considering the dependences. The best performing 10 features are further trained through different deep learning methodologies to output the status of the URL as legitimate or phishing. The objective of this work is to use deep learning classifiers and minimal distinguishing features to determine the status of a given URL, includes methods for feature extraction, feature selection, and classification. The feature extractor receives a list of website URLs as input and pulls the necessary features from three sources (URL obfuscation, hyperlink, and third-party based) as needed. By doing a examination and taking dependences thorough into consideration, the algorithm's output assists in choosing the best-performing features. Characteristics that obfuscate URLs These qualities are those that may be deduced from the URL itself. The incorporation of website content on third-party services is not a part of these functionalities. Before specifying various URL-based features, it is necessary to comprehend how an average URL is constructed. To find existing resources on the Internet, a URL is a particular Uniform Resource Identifier (URI). When a web client asks the server to provide resources like HTML, CSS, photos, videos, or other hypermedia, it is used. A URL typically has four or five elements. To determine whether a URL is real or phishing, the top 10 attributes are subsequently trained using various deep learning approaches.

VI. CONCLUSION

Anti-phishing is the technique of preventing or resolving phishing attacks or scams in which attackers try to get sensitive data or personal information by pretending to be a reliable source. In this research, we develop a browser named PD Browser to recognize phishing websites. The CNN and LSTM were used as a combination classifier in a Novel method known as the LCR to investigate the prospect of distinguishing distinctive authentic URLs from phishing URLs. The main contribution of this research is the incorporation of hybrid features that were taken from text, images, and frames and used to create a strong deep learning solution which examined the best ways to combine picture, text, and frame characteristics with a deep learning algorithm to develop an unified scheme for phishing detection.

Ι

INTERNATIONAL JOURNAL OF SCIENTIFIC RESEARCH IN ENGINEERING AND MANAGEMENT (IJSREM)

OLUME: 07 ISSUE: 09 | SEPTEMBER - 2023

SJIF RATING: 8.176

ISSN: 2582-3930

REFERENCES

- [1] Tanimu, Jibrilla & Shiaeles, Stavros. (2022). Phishing Detection Using Machine Learning Algorithm. 10.1109/CSR54599.2022.9850316.
- [2] Lin, Yun & Liu, Ruofan & Divakaran, Dinil Mon & Ng, Jun & Chan, Qing & Lu, Yiwen & Si, Yuxuan & Zhang, Fan & Dong, Jin. (2021). Phishpedia: A Hybrid Deep Learning Based Approach to Visually Identify Phishing Webpages (USENIX Security).
- [3] Bhavsar, Vaishnavi & Kadlak, Aditya & Sharma, Shabnam. (2018). Study on Phishing Attacks. International Journal of Computer Applications. 182. 27-29. 10.5120/ijca2018918286.
- [4] Ping Yi, Yuxiang Guan, Futai Zou, Yao Yao, Wei Wang, Ting Zhu, "Web Phishing Detection Using a Deep Learning Framework", Wireless Communications and Mobile Computing, vol. 2018, Article ID 4678746, 9 pages, 2018. https://doi.org/10.1155/2018/4678746.
- [5] Li T, Kou G, Peng Y (2020) Improving malicious urls detection via feature engineering: Linear and nonlinear space transformation methods. Inf Syst 91:101494

- [6] Goel D, Jain AK (2018) Mobile phishing attacks and defence mechanisms: state of art and open research challenges. Computer Security 73:519–544
- [7] Navin R 1, Dr Yuvaraj N (2019)Identification of Phishing Website using Deep Learning Algorithm,International Research Journal of Engineering and Technology (IRJET),Volume: 06 Issue: 01, pp 1230-1235.
- [8] Anjali Gupta, Juili Joshi, Khyati Thakker, Chitra bhole, "CONTENT BASED APPROACH FOR DETECTION OF PHISHING SITES" Apr-2015.
- [9] S. Jagadeesan, AnchitChaturvedi, Shashank Kumar (2018) "URL Phishing Analysis using Random Forest".
- [10] Adebowale, Moruf & Lwin, Khin & Hossain, Alamgir.
 (2020), Intelligent Phishing Detection Scheme Algorithms Using Deep Learning. Journal of Enterprise Information Management. ahead-of-print. 10.1108/JEIM-01-2020 0036.