

Innovate a Model of Phishing Website And detection With Features Tools

Roselinourd.J¹, Abdhul Rahman.M.A², Keshor. M³, Jafferset.S⁴

¹Assistant Professor

^{2, 3, 4}Dept of Computer Science and Engineering and Technology

^{1, 2, 3, 4}RAAK College of Engineering and Technology, Puducherry, Pin-605010, India

Abstract- *The phishing email is one of the significant threats in the world today and has caused tremendous financial losses. Although the methods of confrontation are continually being updated, the results of those methods are not very satisfactory at present. Moreover, phishing emails are growing at an alarming rate in recent years. Therefore, more effective phishing detection technology is needed to limit the threat of phishing emails. In this article, we first analysed the structure of the email Then, based on an improved Recurrent Convolution Neural Network (RCNN) model with multilevel vectors and attention mechanisms, we proposed a new named phishing email detection model, to be used to model emails at the subject level, email body level, character level, and word level at the same time. To evaluate its effectiveness, we use an unbalanced dataset that presents the actual ratio of phishing emails to legitimate emails. Experimental results show. Meanwhile, the ensure that the filter can identify phishing emails with high probability and filter out legitimate emails as little as possible. This promising result is superior to the existing detection methods and verifies the effectiveness of in detecting phish.*

Keywords- phishing website detection R-CNN algorithm, website analyses, splitting clone and original website.

I. INTRODUCTION

Nowadays Phishing becomes a main area of concern for security researchers because it is not difficult to create the fake website which looks so close to legitimate website Experts can identify fake websites but not all the users can identify the fake website and such users become the victim of phishing attack. Main aim of the attacker is to steal banks account credentials. In United States businesses, there is a loss of US\$2billion per year because their clients become victim to phishing. In the 3rd Microsoft Computer Safer Index report published in February 2014, it was estimated that the annual global impact of fraud could be as high as \$5 billion. Phishing attacks succeed due to a lack of user awareness. Because phishing attacks exploit user weaknesses, mitigating them is difficult but improving phishing detection techniques is important. The general method of detecting phishing websites

by updating the blacklisted Internet Protocol (IP) URLs in the anti-virus database, is also known as the "blacklisting" method. and many other simple techniques including: fast-fluxion which a proxies are automatically generated to host the web-page; algorithmic generation of new URLs; etc. Major drawback of this method is that, it cannot detect zero-hour phishing attack. Heuristic based detection which includes characteristics that are found to exist in phishing attacks in reality and can detect zero-hour phishing attack, but the characteristics are not guaranteed to always exist in such attacks and false positive rate in detection is very high. To overcome the drawbacks of blacklist and heuristics based method, many security researchers now focused on machine learning techniques. Machine learning technology consists of a many algorithms which requires past data to make decision or prediction on future data. Using this technique, algorithm will analysed various blacklisted and legitimate URLs and their features to accurately detect the phishing websites including zero- hour phishing websites.

Planning

This initial step is used to collect confidential data of users in the form of e-mail lists, templates of scam pages as well as retrieving information from consumers of phishing identifications. Through various techniques and Trojan malwares the computers can easily be compromised (also known as Roots). Through various platforms the scammers get access to proof of notion exploits which enable the scammers to gain admittance to vulnerable computes.

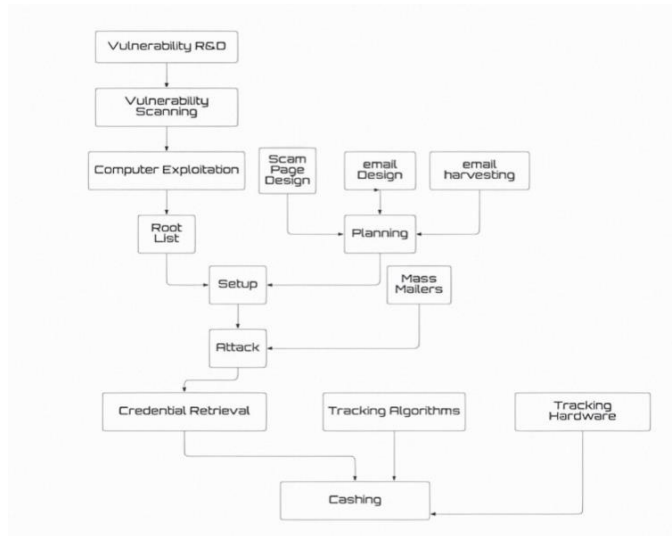
Setup

The further steps involves ensuring the proper scam pages infrastructure on the compromised hosts used in the phishing attack.

Attack

There are millions of programs which have been written to handle mass mails, which enable a scammer to send out e-mails en masse using readily available right tools. The

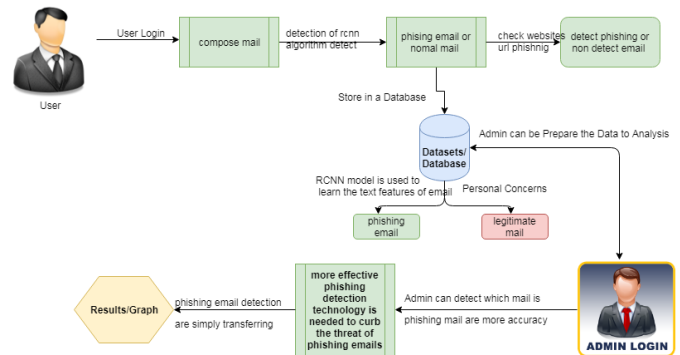
scammers learn through many online tutorials which provide easy explanations on how to send fake e-mails via different programs.



II. EXISTING WORK

Various techniques for detecting phishing emails are mentioned in the literature. In the entire technology development process, there are mainly three types of technical methods including blacklist mechanisms, classification algorithms based on machine learning and based on deep learning. From previous work, the existing detection methods based on the blacklist mechanism mainly rely on people’s identification and reporting of phishing links requiring a large amount of manpower and time. However, applying artificial intelligence to the detection method based on a machine learning classification algorithm requires feature engineering to manually find representative features that are not conducive to the migration of application scenarios. Moreover, the current detection method based on deep learning is limited to word embedding in the content representation of the email. These methods directly transferred natural language processing (NLP) and deep learning technology, ignoring the specificity of phishing email detection so that the results were not ideal. Given the methods mentioned above and the corresponding problems, we set to study phishing email detection systematically based on deep learning. Specifically, this paper makes the following contributions.

Architecture



Disadvantages

- With respect to the particularity of the email text, we analysed the email structure, and mine the text features from four more detailed parts: the email header, the email body, the word-level, and the char-level.
- The RCNN model is improved by using the Then, the email is modelled from multiple levels using an improved RCNN model. Noise is introduced as little as possible, and the context information of the email can be better captured.

III. PROPOSED SYSTEM

With the emergence of email, the convenience of communication has led to the problem of massive spam, especially phishing attacks through email. Various anti phishing technologies have been proposed to solve the problem of phishing attacks. studied the effectiveness of phishing blacklists. Blacklists mainly include sender blacklists and link blacklists. This detection method extracts the sender’s address and link address in the message and checks whether it is in the blacklist to distinguish whether the email is a phishing email. The update of a blacklist is usually reported by users, and whether it is a phishing website or not is manually identified. At present, the two well-known phishing websites are Phish Tank and Open Phish. To some extent, the perfection of the blacklist determines the effectiveness of this method based on the blacklist mechanism for phishing email detection. The current situation is that new threats may not only cause severe damage to customers’ computers but also aim to steal their money and identity. Among these threats, phishing is a noteworthy one and is a criminal activity that uses social engineering and technology to steal a victim’s identity data and account information. According to a report from the Anti-Phishing Working compared with the fourth quarter of According to the striking data, it is clear that phishing has shown an apparent upward trend in recent years. Similarly, the harm caused by phishing can be imagined as well

PYTHON

Python is a general-purpose interpreted, interactive, object-oriented, and high-level programming language. An interpreted language, Python has a design philosophy that emphasizes code readability (notably using whitespace indentation to delimit blocks rather than curly brackets or keywords), and a syntax that allows programmers to express concepts in fewer code than might be used in languages such as C++ or Java. It provides constructs that enable clear programming on both small and large scales. Python interpreters are available for many operating systems. CPython, the implementation of Python, is source software and has a community-based development model, as do nearly all of its variant implementations. CPython is managed by the non-profit Python Software Foundation. Python features a type system and automatic memory management. It supports multiple programming paradigms, including object-oriented, imperative, functional and procedural, and has a large and comprehensive standard library.

DJANGO

Django is a high-level Python Web framework that encourages rapid development and clean, pragmatic design. Built by experienced developers, it takes care of much of the hassle of Web development, so you can focus on writing your app without needing to reinvent the wheel. It's free and open source.

Django is a high-level Python Web framework that encourages rapid development and clean, pragmatic design. Built by experienced developers, it takes care of much of the hassle of Web development, so you can focus on writing your app without needing to reinvent the wheel. It's free and open source. Django's primary goal is to ease the creation of complex, database-driven websites. Django emphasizes reusability and "plug ability" of components, rapid development, and the principle of don't repeat yourself. Python is used throughout, even for settings files and data models.

Advantages

- Phishing email refers to an attacker using a fake email to trick the recipient into returning information such as an account password to a designated recipient.
- Additionally, it may be used to trick recipients into entering special web pages, which are usually disguised as real web pages, such as a bank's web page, to convince users to enter sensitive information such as a credit card

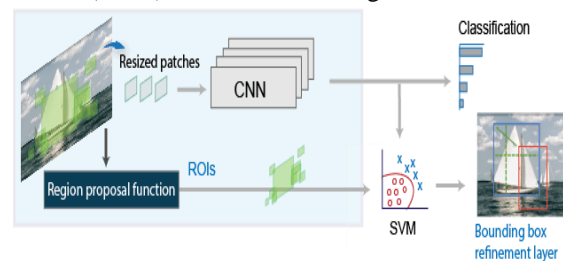
or bank card number and password. Although the attack of phishing email seems simple, its harm is immense.

III. ALGORITHM

Let's quickly summarize the different algorithms in the R-CNN family (R-CNN, Fast R-CNN, and Faster R-CNN) that we saw in the first article. This will help lay the ground for our implementation part later when we will predict the bounding boxes present in previously unseen images (new data). R-CNN extracts a bunch of regions from the given image using selective search, and then checks if any of these boxes contains an object. We first extract these regions, and for each region, CNN is used to extract specific features. Finally, these features are then used to detect objects. Unfortunately, R-CNN becomes rather slow due to these multiple steps involved in the process. Fast R-CNN, on the other hand, passes the entire image to Convnet which generates regions of interest (instead of passing the extracted regions from the image). Also, instead of using three different models (as we saw in R-CNN), it uses a single model which extracts features from the regions, classifies them into different classes, and returns the bounding boxes. All these steps are done simultaneously, thus making it execute faster as compared to R-CNN. Fast R-CNN is, however, not fast enough when applied on a large dataset as it also uses selective search for extracting the regions.

Detection Using R-CNN Algorithms

The R-CNN detector first generates region proposals using an algorithm such as Edge Boxes. The proposal regions are cropped out of the image and resized. Then, the CNN classifies the cropped and resized regions. Finally, the region proposal bounding boxes are refined by a support vector machine (SVM) that is trained using CNN features.



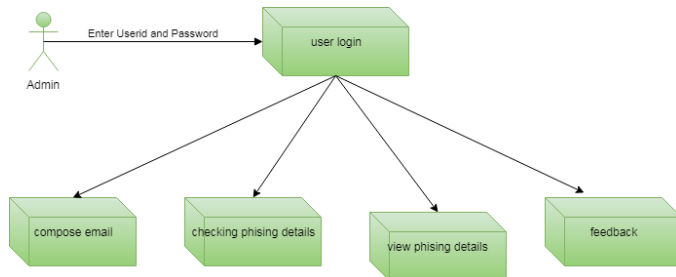
Fast R-CNN

As in the R-CNN detector the Fast R-CNN detector also uses an algorithm like Edge Boxes to generate region proposals. Unlike the R-CNN detector, which crops and resizes region proposals, the Fast R-CNN detector processes the entire image. Whereas an R-CNN detector must classify each region, Fast R-CNN pools CNN features corresponding

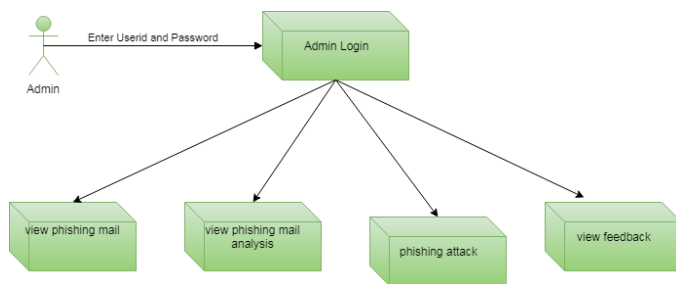
to each region proposal. Fast R-CNN is more efficient than R-CNN, because in the Fast R-CNN detector, the computations for overlapping regions are shared.

COMPONENT DIAGRAM

A) User



B) Admin



IV. REQUIREMENT ANALYSIS

The project involved analysing the design of few applications so as to make the application more users friendly. To do so, it was really important to keep the navigations from one screen to the other well-ordered and at the same time reducing the amount of typing the user needs to do. In order to make the application more accessible, the browser version had to be chosen so that it is compatible with most of the Browsers

V. MODULES

DATASET

The dataset has been divided into a training set and testing set. Both the training set and the testing set contain emails without header and emails with header. In this paper, we only focus on email data with the header. Due to the irrationality of the segmentation of the training set and the testing set in the original dataset, after merging the two datasets, the training-validation set and the testing set are reinvented. The dataset is divided by stratified random sampling; that is, random samples are taken from legitimate email and phishing email at the same proportion. This ensures

that the two datasets used in training and testing phases are well.

USER QUERIES

Users can have queries about the process. This part of project is dedicated to make and get response for queries that are needed to answerable. The major part of the modules is making project as interactive one, queries have been very normally arise to users regarding different details about the process.

GRAPH ANALYSIS

Graph analysis is the part where admin can know the statistics about process of details. The data are taken from the project flow and it shows until updated value. The data are given clear solution to admin that part of improvement and user satisfaction and other factors.

ANALYSIS

Analysis of email structure. a circle represents a character, and a rectangle represents word. A rectangle is filled with an indefinite number of circles, indicating that the word consists of an indefinite number of characters.

VI. SYSTEM STUDY

FEASIBILITY STUDY

The feasibility of the project is analysed in this phase and business proposal is put forth with a very general plan for the project and some cost estimates. During system analysis the feasibility study of the proposed system is to be carried out. This is to ensure that the proposed system is not a burden to the company. For feasibility analysis, some understanding of the major requirements for the system is essential.

Three key considerations involved in the feasibility analysis are

ECONOMICAL FEASIBILITY

This study is carried out to check the economic impact that the system will have on the organization. The amount of fund that the company can pour into the research and development of the system is limited. The expenditures must be justified. Thus the developed system as well within the budget and this was achieved because most of the technologies used are freely available. Only the customized products had to be purchased.

TECHNICAL FEASIBILITY

This study is carried out to check the technical feasibility, that is, the technical requirements of the system. Any system developed must not have a high demand on the available technical resources. This will lead to high demands on the available technical resources. This will lead to high demands being placed on the client. The developed system must have a modest requirement, as only minimal or null changes are required for implementing this system.

SOCIAL FEASIBILITY

The aspect of study is to check the level of acceptance of the system by the user. This includes the process of training the user to use the system efficiently. The user must not feel threatened by the system, instead must accept it as a necessity. The level of acceptance by the users solely depends on the methods that are employed to educate the user about the system and to make him familiar with it. His level of confidence must be raised so that he is also able to make some constructive criticism, which is welcomed, as he is the final user of the system.

VII. SYSTEM TEST

The purpose of testing is to discover errors. Testing is the process of trying to discover every conceivable fault or weakness in a work product. It provides a way to check the functionality of components, subassemblies, assemblies and/or a finished product. It is the process of exercising software with the intent of ensuring that the Software system meets its requirements and user expectations and does not fail in an unacceptable manner. There are various types of test. Each test type addresses a specific testing requirement.

VIII. CONCLUSION

We use a new deep learning model named to detect phishing emails. The model employs an improved RCNN to model the email header and the email body at both the character level and the word level. Therefore, the noise is introduced into the model minimally. In the model, we use the attention mechanism in the header and the body, making the model pay more attention to the more valuable information between them. We use the unbalanced dataset closer to the real-world situation to conduct experiments and evaluate the model. The model obtains a promising result. Several experiments are performed to demonstrate the benefits of the proposed model. For future work, we will focus on how to improve our model for detecting phishing emails with no email header and only an email body.

REFERENCES

- [1] Korkmaz, M., Sahingoz, O. K., & Diri, B., "Detection of Phishing Websites by Using Machine Learning- Based URL Analysis", IEEE IIT - Kharagpur, Istanbul, Turkey, 1-Jul-2020.
- [2] Arun Kulkarni & Leonard L. Brown, "Phishing Websites Detection using Machine Learning", International Journal of Advanced Computer Science and Applications, Tyler, TX, 2019.
- [3] A, A. A., & K, P. "Towards the Detection of Phishing Attacks", IEEE 4th International Conference on Trends in Electronics and Informatics, Coimbatore, India, July 27-2020.
- [4] Ma, Justin, et al. "Beyond blacklists: learning to detect malicious web sites from suspicious URLs." Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, 2009.
- [5] Khonji, Mahmoud, Youssef Iraqi, and Andrew Jones. "Phishing detection: a literature survey." Communications Surveys & Tutorials, IEEE 15.4 (2013): 2091-2121.
- [6] H. N. A. Pham and E. Triantaphyllou, "The impact of overfitting and overgeneralization on the classification accuracy in data mining," in *Soft Computing for Knowledge Discovery and Data Mining*, O. Maimon and L. Rokach, Eds. 2008.
- [7] R. M. Mohammad, F. Thabtah, L. McCluskey, "An Assessment of Features Related to Phishing Websites Using An Automated Technique," International Conference for Internet Technology and Secured Transactions, pp. 492-497, IEEE, 2012.
- [8] I. Qabajeh, F. Thabtah, "An Experimental Study for Assessing Email Classification Attributes Using Feature Selection Methods," 3rd International Conference on Advanced Computer Science Applications and Technologies (ACSAT), pp. 125-132, IEEE, 2014.
- [9] H. Liu, J. Li, L. Wong, "A comparative study on feature selection and classification methods using gene expression profiles and proteomic patterns," *Genome informatics*, vol. 13, pp. 51-60, 2002.
- [10] Visualizing Models, Data, and Training With TensorBoard—PyTorch Tutorials 1.9.1+Cu102 Documentation. PyTorch.org. Accessed: Oct. 15, 2021.
- [11] N. Gupta, S. Mujumdar, H. Patel, S. Masuda, N. Panwar, S. Bandyopadhyay, S. Mehta, S. Guttula, S. Afzal, R. Sharma Mittal, and V. Munigala, "Data quality for machine learning tasks," in *Proc. 27th ACM SIGKDD Conf. Knowl. Discovery Data Mining*, Aug. 2021,
- [12] S. Seo, C. Kim, H. Kim, K. Mo, and P. Kang, "Comparative study of deep learning-based sentiment

- classification,” IEEE Access, vol. 8, pp. 6861–6875, 2020, doi: 10.1109/ACCESS.2019.2963426.
- [13] J. Chung, C. Gulcehre, and K. Cho. (Dec. 2014). Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. Accessed Aug./29/2021. [Online] Available: <https://ashutoshtripathicom.files>.
- [14] Huang, Huajun, Junshan Tan, and Lingxi Liu. “Countermeasure techniques for deceptive phishing attack.” New Trends in Information and Service Science, 2009. NISS’09. International Conference on. IEEE, 2009.
- [15] [Anti Phishing Working Group. (2015. March.) APWG Phishing Activity Trend Report 2nd Quarter 2010. [Online]. Available: http://docs.apwg.org/reports/apwg_trends_report_q2_2014.pdf.
- [16] A. Shewalkar, D. Nyavanandi, and S. A. Ludwig, “Performance evaluation of deep neural networks applied to speech recognition: RNN, LSTM and GRU,” J. Artif. Intell. Soft Comput. Res., vol. 9, no. 4, pp. 235–245, Oct. 2019.
- [17] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using RNN encoder-decoder for statistical machine translation,” 2014.
- [18] S. Smys, J. I. Zong Chen, and S. Shakya, “Survey on neural network architectures with deep learning,” J. Soft Comput. Paradigm, vol. 2, no. 3, pp. 186–194, Jul. 2020.
- [19] Satish S, Suresh Babu K (2013) “Phishing Websites Detection Based On Web Source Code And Url In The Webpage” Aburrou, Maher & Hossain, Mohammed & Dahal, Keshav & Thabtah, Fadi. (2010).
- [20] V. Babel, K. Singh, S. K. Jangir, B. Singh, and S. Kumar. (2019). Journal of Analysis and Computation (JAC) Evaluation Methods for Machine Learning. Accessed: Oct. 14, 2021.
- [21] El Aassal, A., Baki, S., Das, A., & Verma, R. M., “An In-Depth Benchmarking and Evaluation of Phishing Detection Research for Security Needs”, IEEE Access, Houston, U.S., 5-Feb-2020.
- [22] S. Kumar. (2019). Malicious and Benign URLs. kaggle.com. Accessed: Oct. 20, 2021. [Online]. Available: <https://www.kaggle.com/siddharthkumar25/malicious-and-benign-urls>.
- [23] J. Zhang, Y. Ou, D. Li, and Y. Xin, “A prior-based transfer learning method for the phishing detection,” J. Netw., vol. 7, no. 8, p. 1201, Aug. 2012, doi:10.4304/jnw.7.8.1201-1207.