# A Novel Machine Learning & NLP based approach for Analysing Phishing Attack

Dr. Savita Sangam[1], Rajshree Shejwal [2], Gouri Shelar [3]
[1, 2, 3]Dept. of Information Technology Engineering, SSJCOE, Maharashtra, India

Abstract: Phishing refers to a sophisticated social networking technique used to target sensitive victims. Now the days of people who use the internet to do many things like shopping, advertising, shipping etc. Phishing is a major attack on the website, people facing in their day to day life. Phishers uses custom-looking websites similar to those of original websites. As technology continues to grow, phishing scams are beginning to develop rapidly and this needs to be prevented by using phishing techniques. Machine learning is a powerful tool used to combat phishing attacks. This paper looks at the materials used to derive detection and detection methods using machine learning and NLP
Keywords: Phishing, Detection, Machine Learning, NLP

## I. INTRODUCTION

In our society, the security of confidential information is a matter that concerns everyone. Phishing is a type of attack on social engineers that focuses on obtaining sensitive information in disguise as an organization worthy of trust. Electronic communications, such as email or text messaging are common platforms for sending phishing attacks. Attackers are often disguised as popular social websites, banks, and executives from the IT department or online shopping websites. Most methods check the URLs that contain the message. Our approach performs a non-linear analysis of the text transmitted by the attacker to verify the validity of each paragraph. Sentence is considered dangerous if it requires sensitive information or dictates the performance of an action that would disclose personal information. The roles of each word in a sentence, the way we use it determines whether a sentence is a question or a command. Much research has been ongoing to prevent phishing attacks by various communities around the world. Phishing attacks can be prevented by detecting websites and notifying users to detect phishing websites. Machine learning algorithms have been one of the most effective techniques for detecting phishing websites. In this study, various ways to access phishing websites were discussed.

## II. LITERATURE SURVEY

A. TianruiPeng & et. All Data Hacking Attacks sensitive data using natural language processing and machine learning. We are introducing a way to detect cybercrime attacks. Our approach relies on text analysis, rather than the metadata that may be associated with emails. As a result, our method is effective for obtaining phishing emails containing clean text. Our results on cybercrime emails show the most advanced recollections that indicate that semantic information is a strong indicator of communication engineering. Electronic communications, such as email or text messaging are common platforms for sending phishing attacks. Phishing has been shown as an effective attack throughout the years, fooling many people. Attackers are often disguised as popular social websites, banks, and executives from the IT department or online shopping websites. These emails may urge users to click links to perform malware downloads, or enter personal information on a malicious website that has the same legal appearance. Many ways to get automated email automation are based on email metadata, data associated with emails not related to the semantic meaning of the text. Most methods check the URLs that contain the message.

B. Phish Storm: Identifying Phishing Analytics by Streaming Analytics Samuel Marshal & et all Despite the rise of prevention methods, phishing is currently a significant threat because the primary accounting used is based on active URL testing. This method does not work well due to the short time of designing Web sites, making recent methods dependent on the use of real-time URLs or techniques for obtaining relevant criminal information. In this paper we introduce Phish Storm, a phishing scheme that can process real-time any URL to identify phishing sites. Phish Storm can interact with any email server or HTTP proxy. We argue that phishing URLs generally have little to do between the part of the URL to be registered (the base domain) and the remainder of the URL (the top-level domain, method, query). We show in this paper that experimental evidence supports this observation and can be used to identify phishing sites. For this purpose, we define a new concept, related to URL combinations and test it using features extracted from URL names based on query data from Google and Yahoo search engines. These features are used for machine learning classification to retrieve phishing URLs from sensitive data.
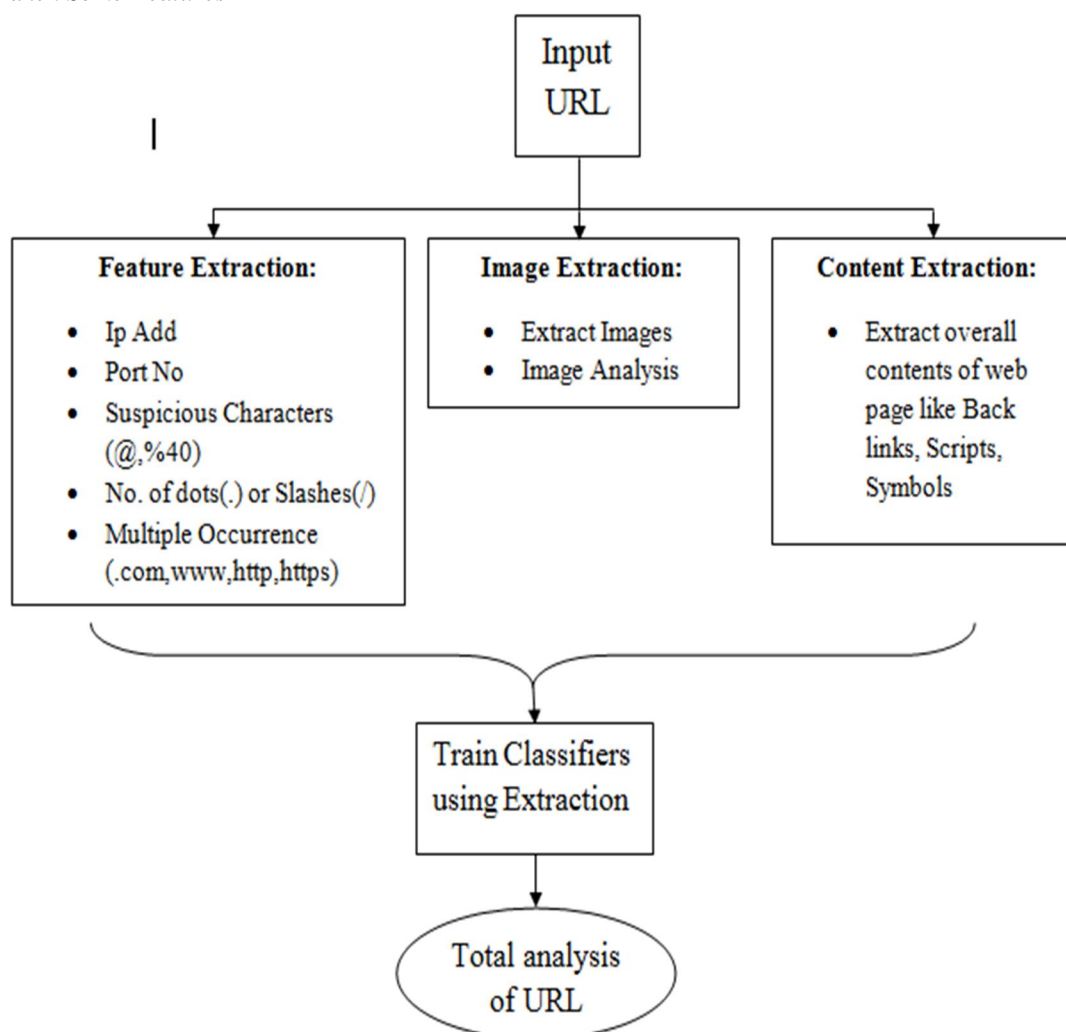
C. Phishing URL discovery: web mining learning method. B.E. Sananse & et all 2016 the identification of malicious queries and commands depends on the presence of a list of topics in the verb-direct object in pairs where their presence of query command is intended to be malicious. To generate a list of adults in an article we use machine learning, build a Naive Bayes data-driven class that is distributed across multiple countries, and is used to classify text. We used the Multinomial () function from the Scikit-learn Python library that uses this algorithm. This algorithm generates each prediction label (the verb of a specific object), and produces a confidence note in the prediction. Range of confidence scores is 0 to 1, with a validation verification number. We used 1000 phishing emails from Nazi phishing emails. We tested our results on all 5014 emails in the Nazi phishing email. As well as about 5000 phishing emails on Enron Corpus. Before applying for machine learning, we used Stanford typed water parser to extract all (verb-direct object) from all expressions by identifying "nsubjpass" and "dobj" based on the findings of each sentence.
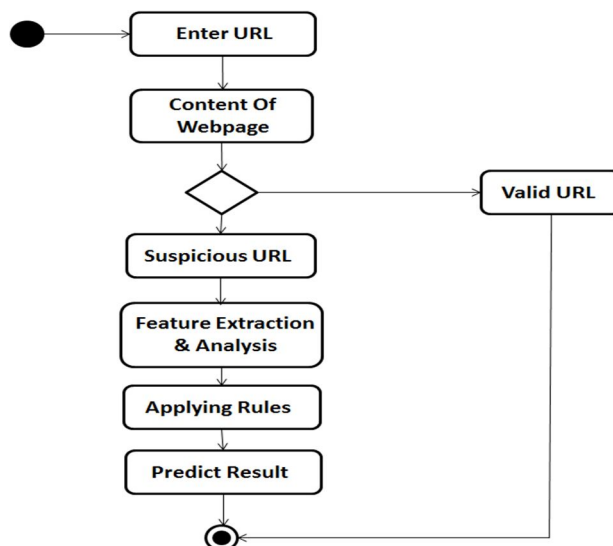
### III. IMPLEMENTATION

Phishing attacks are one of the most common and least defended security threats today. We present an approach which uses natural language processing techniques to analyze text and detect inappropriate statements which are indicative of phishing attacks. Our approach is novel compared to previous work because it focuses on the natural language text contained in the attack, performing semantic analysis of the text to detect malicious intent. To demonstrate the effectiveness of our approach, we have evaluated it using a large benchmark set of phishing URL.

A machine learning based anti-phishing system which is based on Uniform Resource Locator (URL) features.

*A. We Have Taken Some Features*

International Journal for Research in Applied Science & Engineering Technology (IJRASET)
ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.429
Volume 8 Issue V May 2020- Available at www.ijraset.com

Below is a only shows the flow of data between user and the system. The user enter URL .the feature extraction are applying to URL and check the URL length, number of suspicious character, multiple occurrence of symbols and protocol.



*B. Detecting Accuracy of URL*

*1) IP Address:* Example: " http://66.135.200.145 "

*2) Protocol:* Example: http, https, ftp, SSH, Telnet Natural Language Processing (NLP)

*a) Step1:* The input to natural language processing will be a simple stream of Unicode characters (typically UTF-8). Basic processing will be required to convert this character stream into a sequence of lexical items (words, phrases, and syntactic markers) which can then be used to better understand the content.

*b) Step2:* Decide on Micro Understanding It is used for: clustering, categorization, similarity, topic analysis, word clouds and summarization.

*C. Numbers of Dots and Slashes*

Example: *https://www.**idbi**.com*

https://www.ldbi.com/.in

*1)* Multiple Occurrence(.com ,https ,http):

*2)* Example: *http://www.google.com/url?q=http://www.badsite.com*

## IV. RESULT

*1) User Browser:* In this diagram user enter the URL and perform their activity but this time more chances of fake URLs and redirected links are attached in valid URL, so our project help us to user to detect the valid URLs
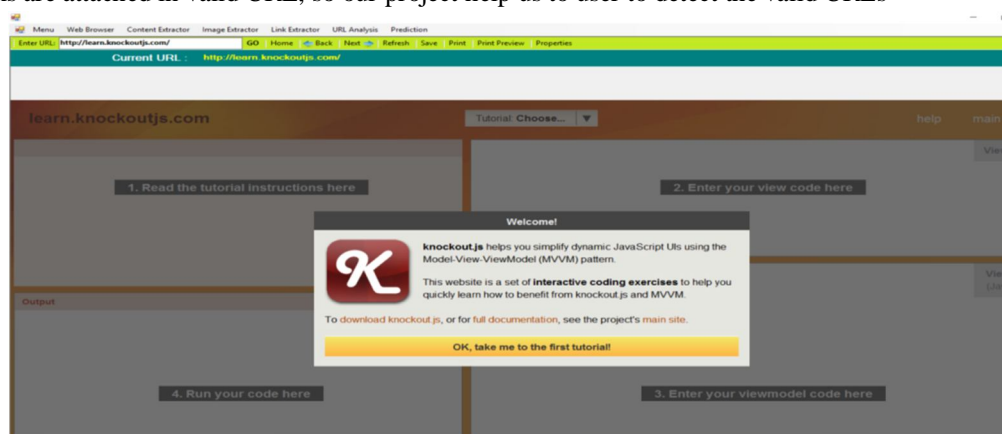


Fig 4.1 User Browser

2) *Content Extractor:* In this diagram current URL are used and all content are collected in currant URL then by using DOM Href , Links, Symbols are separated and showing the count.
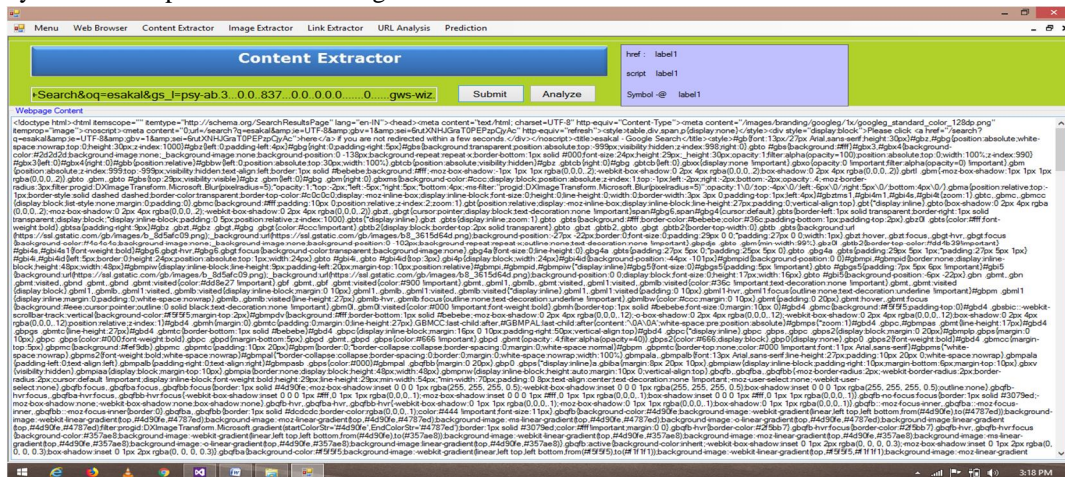


Fig 4.2 Content Extractor

3) *Image Extractor:* When the user checking content extractor this time all link, symbols, slashes and images are separated by using the DOM tree and all images all display in this Image Extractor.
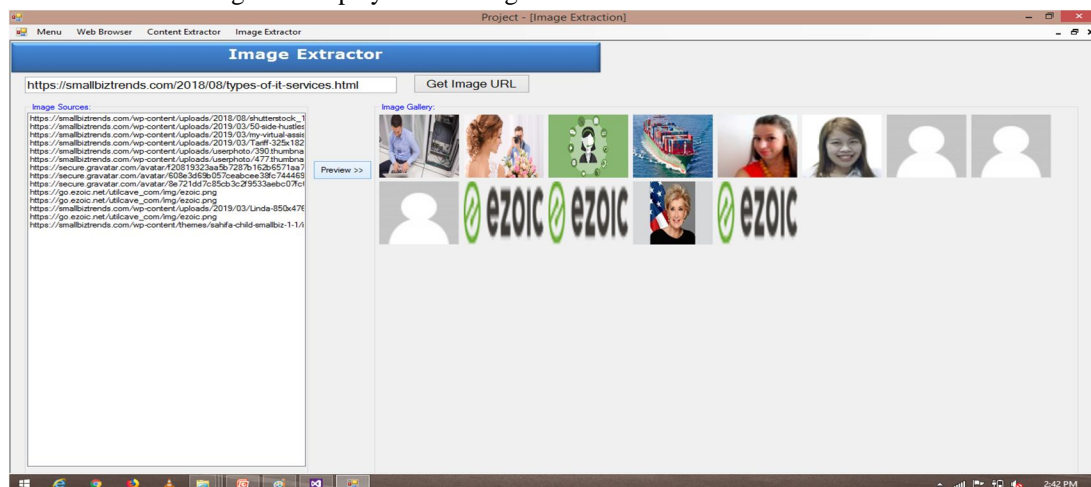


Fig 4.3 Image Extractor

4) *Link Extractor:* In this diagram all links are extracted and showing the total no of links contain this URL
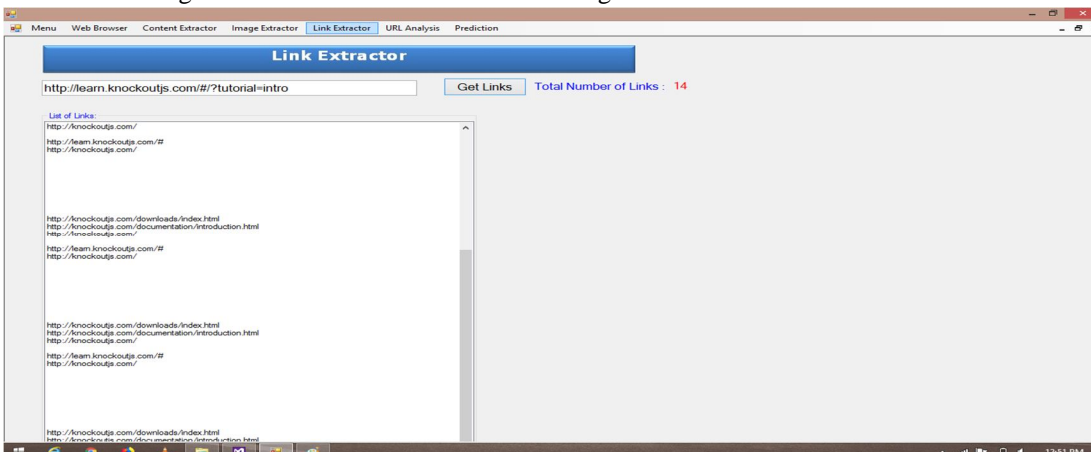


Fig 4.4 Link Extractor

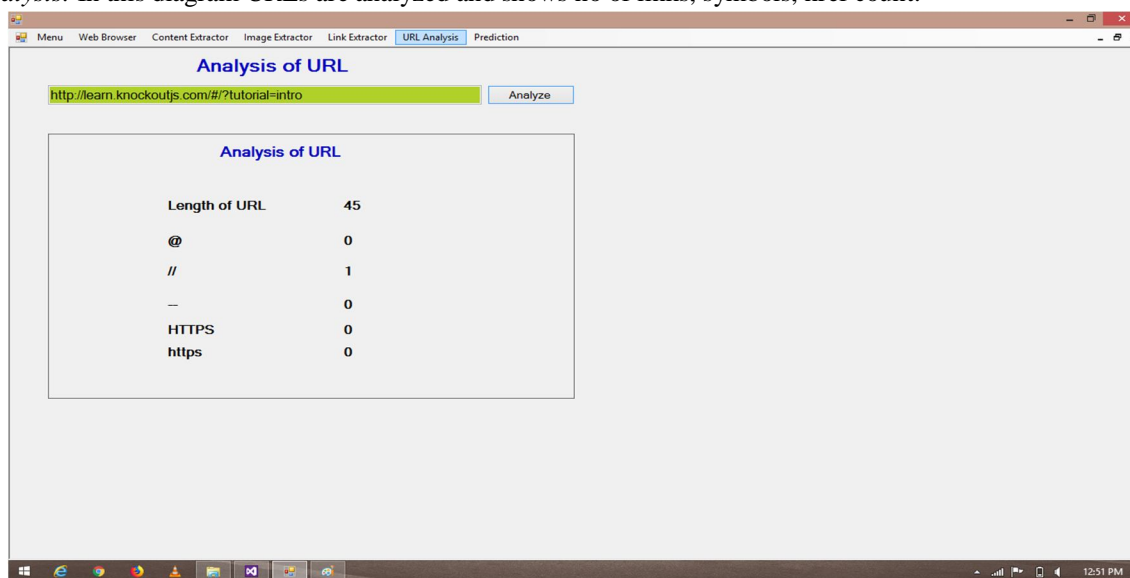5) *URL Analysis:* In this diagram URLs are analyzed and shows no of links, symbols, href count.



Fig 4.5 URL Analysis.

6) *URL Validate:* In this diagram shows the final output of our project. here showing complete analysis of URL and then validate the URL and check the accuracy of URL
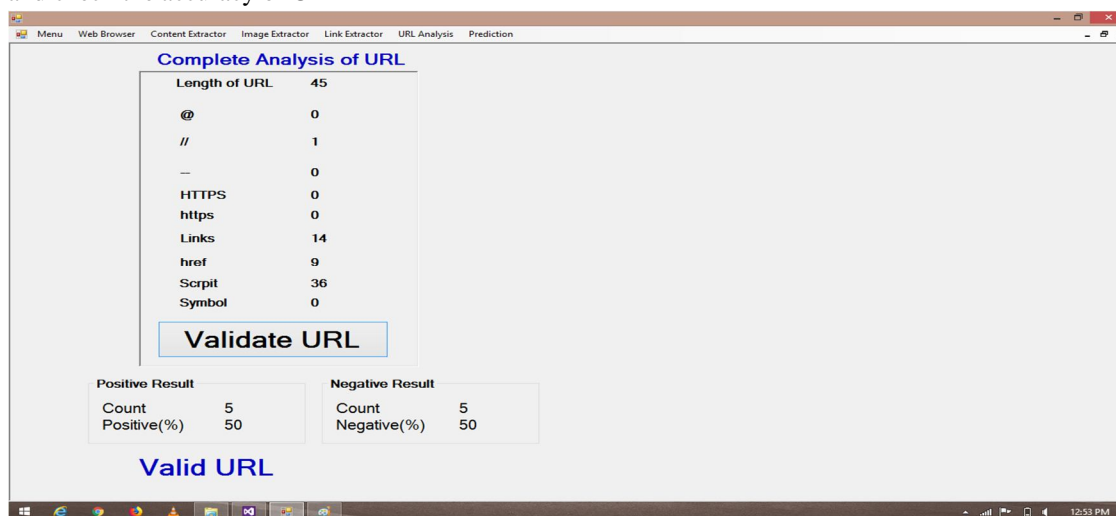


Fig 4.6 URL Validate.

## V. ANALYSIS

| Sr.No | Papers | Techniques | Accuracy |
|---|---|---|---|
| 1. | Detecting Phishing Attacks using NLP and Machine Learning | SEAHound Algorithm | Legitimate URL(90% Accuracy) |
| 2. | Fresh Phish :A Framework for auto detection of phishing website | Tensorflow | Legitimate URL(70% Accuracy) |
| 3. | Phishing URL detection machine learning web mining approach | SVM Algorithm | Legitimate URL(65% Accuracy) |
| 4. | Phishstrom: Detection Phishing with Streaming analysis | Search Engine Feature Extraction | Legitimate URL(60% Accuracy) |
| 5. | A novel machine learning and NLP based approach for detecting phishing attacks | NLP & Feature Extraction | Legitimate URL(75% Accuracy) |

## VI. CONCLUSION

After studying different research papers published nationally and internationally, we have seen many types of techniques. Web browser which successfully browse any kind of websites & work on user behaviour. Our results on phishing URLs demonstrate significantly improved recall which demonstrates that semantic information is a strong indicator of social engineering .Thus the user can check all the contents of web page of that entered URL to decide sharing of personal Information.

## REFERENCES

[1]   R.kiruthiga, D.Akila, Phishing website detection using machine learning International Journal Of Recent Technology and Engineering, September 2019.

[2]    Ram Basnet, Srinivas Mukkamala, Detection of Phishing Attacks: Machine Learning Approach, New Mexico Tech, New Mexico, USA, 2008.

[3]   Michael C.Kotson, MIT Lincoln Laboratory Lexington, MA, Characterizing Phishing Threats with Natural Language Processing, 31 Aug 2015.

[4]   Dr.G.Ravi Kumar, Dr.S.Gunasekaran, Detection of Phishing Attacks using Machine Learning Techniques, Department of CSE Coimbatore Institute of Engineering,Dec2008.

TOGETHER WE REACH THE GOAL

# INTERNATIONAL JOURNAL
# FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 ⊙ (24*7 Support on Whatsapp)