

# SPAM DETECTION OF PHISHING WEBSITES USING ML

Dr. J. Selvakumar<sup>+</sup>, Mr. R.Prithiviraj<sup>+</sup>, Mr. Joshua Jafferson<sup>+</sup>, Mr.S.Bashyam<sup>+</sup>

<sup>+</sup>,*ECE Department, SRM Institute of Science and Technology, Chennai, Tamil Nadu.*

[selvakuj@srmist.edu.in](mailto:selvakuj@srmist.edu.in), [prithivir@srmist.edu.in](mailto:prithivir@srmist.edu.in), [joshuaj@srmist.edu.in](mailto:joshuaj@srmist.edu.in), [bashyams@srmist.edu.in](mailto:bashyams@srmist.edu.in)

## ABSTRACT

In today's internet era various websites through which a number of individuals purchase items. There are certain online forums which request their users to provide confidential data such as card number, cvv, pin number etc. for various malicious practices. These websites are referred as Phishing Websites. Therefore, to distinguish between the authentic website and the malicious website we suggested an intelligent, adaptable, and efficient model that utilizes Machine learning techniques. We carry through the project using the algorithm of classification and different methods to gather the phishing websites dataset to verify its validity. These spoofing websites are differentiated on certain significant attribute such as encryption standards, Domain Identity, URL and security. The project will utilize machine learning concept thus informing the user if the website is legal or not. This software is highly secured and can be utilized by many E-commerce ventures so as to provide hassle free transaction. Machine Learning design utilized in the project gives good results when compared with other standard classification algorithms. Detection of Phishing web site is ML intelligent and effective model that's supported victimization classification or association data processing algorithms. The algorithms we are using here is logistic regression. We are also using decision tree classifier so that we can make a point-to-point comparison between them which will help us to know parameters like accuracy and time taken.

**Keywords---** Phishing, Phishing Websites, Detection, Machine Learning.

## 1. INTRODUCTION

In the field of software security, phishing is the reprehensibly malicious method of making an attempt to amass sensitive info like usernames, passwords and banking account details by masquerading as a trustworthy entity in an electronic communication. A phishing web site can be a widely projected online strike that tries to swindle user of distinctive information as well as their confidential data like credit card details, account info etc. which can be used falsely to fraud the user. Spoofing causes immense impact on company's revenues, client relationships, selling efforts and overall corporate image. Spoofing is typically dispensed by electronic mail or electronic communication, and it typically sends users to fraudulent website which seem virtually similar to the authentic one. Phishing is a perfect instance to explain social exploitation to fool the users and various network security advancements. In the present generation the internet is used as a platform for various online activities such as bank transactions, bill payments and mobile recharges etc. Due to this the user details are becoming easily available which leads to cybercrime. Phishing is one of the most common cybercrime activities happening presently. In the last six years some anti phishing groups have detected a growth of 36% in the rise of phishing websites. This numbers are increased to 97.3% in last two years. As the rates are increasing day by day it is necessary for the cyber security community to form a system for the detection of these phishing websites [1]. Internet has become a valuable platform for people to interact. Because of which sensitive, private and sometimes confidential data has become easily accessible on the net and it's been an important issue[2].Spoofing is the process of creating a fake website which resemble the legitimate website of the firm with the view of getting hands on the user's private data like credit card details, passwords, username etc[3].Many of the present day techniques for the detection of the phishing sites are weak in opposition to domain name

system which are generally attacks based on poisoning. . A potential process is proposed for the spotting of these critical strikes: characteristics related to the website's performance are used for the classification [4]. Computer systems security quite often depends on action and the decision of user. We explore into users' sensitivity to online strikes by focusing on the basic components that governs end user that is - the human brain [5].

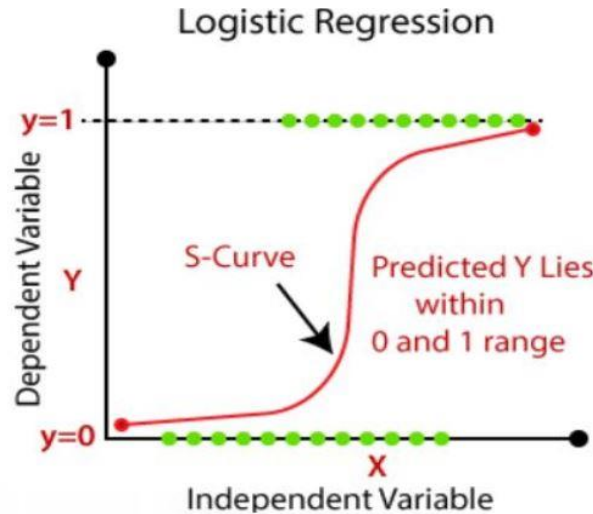
## 2. ALGORITHM

### A. Logistic Algorithm

The calculation that we are utilizing is Logistic Regression. Let us examine the fundamental ideas of Logistic Regression as shown in fig.1. in subtleties and realize what sort of issues it would be able to support us to tackle. The calculated capacity, additionally describes as the capacity of sigmoid that was created by huge analysis so as to depict features of popularity developing in the environment, increasing rapidly. It is a Stormed incline that is used to take any legitimated esteemed number hence guiding it in a transformation somewhere that is in the possibility of '0' and '1', however never precisely at these limits.  $1 / (1 + e^{-value})$ . The utilized condition is helpful in the Logistic relapse shows the portraying, basically the relapse of straight type. Here, input esteems (x) are consolidated straight using coefficients esteems or loads in order to anticipate an esteem to be yielded. The important contrast that is inferred from relapse of direct type is that instead of numerical qualities the yield esteem being demonstrated is a double quality (0 or 1).

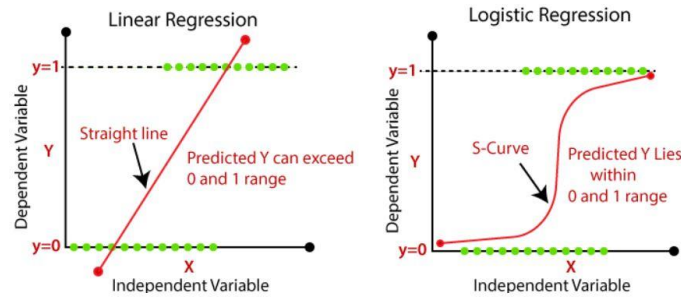
$$Y = e^{(a_0+a_1*x)} / (1 + e^{(a_0+a_1*x)}) \tag{1}$$

Here 'y' is the anticipated yield, 'a<sub>1</sub>' used for the single information esteem (x) coefficient mapping. Each and every block in our prior information usually has a mapped to the coefficient of 'b' that needs to be inferred from our preparation information and also the a<sub>0</sub> is the predisposition or block term.



**Fig. 1. Logistic Regression Curve**

Logistic regression is different from linear regression as shown in fig.2.. Linear regression is used when the variable is continuous in nature where as logistic regression is used when the variable is binary in nature(0 or 1).

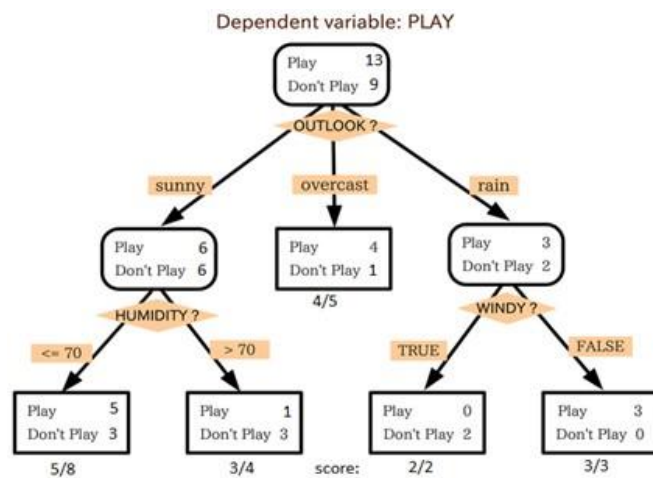


**Fig.2. Linear and Logistic regression graphs**

**B. Decision Tree**

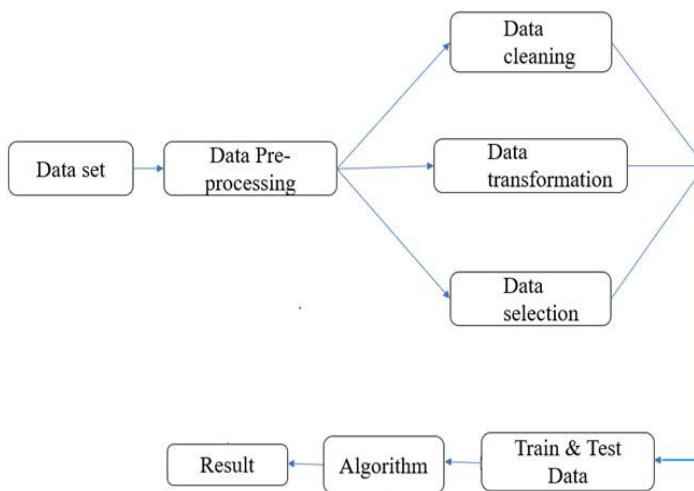
This is one of the regular calculations and it is utilized habitually. Characterization issues are for the most part done by administered learning calculation. Amazingly, it works for both consistent and clear-cut ward factors. In this calculation, the populace is part into at least two homogeneous sets. This is done dependent on generally huge free factors to make as particular gatherings as could be expected under the circumstances.

In the fig.3, you can see that populace is characterized into four unique bunches dependent on various ascribes to distinguish 'on the off-chance that they will play or not'. For the populace to part into various heterogeneous gatherings, different procedures like Chi-square, Gini, Information Gain, entropy are utilized. The most ideal route for getting Decision tree working is to play Jezzball – a great game. Basically, you have a room which contains moving dividers and the dividers should be made in with the end goal that most extreme zone gets tidied up without the balls.



**Fig.3. populace decision tree**

### 3. METHODOLOGY



**Fig.4. Basic flow diagram for spam detection**

Data Collection is an important task in building a machine learning model.

It is the process of gathering information(data) based on some targeted variables to examine and produce some effective outcome. However, some of the data may contain inaccurate values, incomplete values or incorrect values.

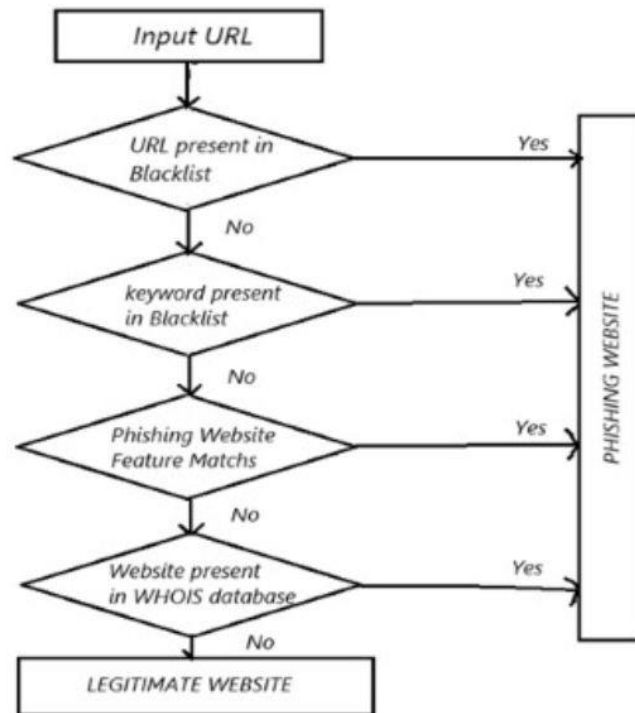
Hence, it is mandatory to remove all the disturbance from the dataset, in order to produce accurate output.

Data pre-processing done as follows: Data cleaning; Data transformation and Data selection.

Data cleaning: Filling the missed values, noisy data to be smoothed, inconsistencies to be resolved and remove or identify outliers. Transformation of data may include smoothing, generalization, transformation which improves the quality of the data

Data selection: Includes some methods or functions which allow us to select the useful data for our system.

#### IV. ARCHITECTURE



**Fig. 5. Architecture Block of Detection algorithm**

Composition is a numerous level process that centres around data such as procedural subtleties, structure programming design, system etc. and makes bond among section. PC programming configuration change constantly as novel techniques such as improved examination and boundary knowledge approach. The uprising of Design proposition is at primary stage.

Thus, design plot procedure doesn't have profundity, adaptability and calculable description which is generally combined along increasingly ordinary designing orders. Anyhow the methods of program structures techniques get over, plan characteristics models are existing and model annotations are applied.

#### 4. RESULTS

Finally, we get the result based on our proposed ML (Logistic regression) algorithms used. From the obtained results it is proved that our proposed algorithm accuracy is close to '1' and 11% (approx.) higher than conventional algorithms.

The accuracy of code using Decision Tree Classifier is

Accuracy 0.8582640647854162

The accuracy of code using Logistics Regression is

Accuracy 0.9616377106298979

The final output of our proposed algorithm is in fig. 5.

```
x_predict = ["https://en-gb.facebook.com/login/",  
            "http://www.discretepackagemovers.com",  
            "https://twitter.com/i/flow/signup",  
            "http://www.skylinexpresscourier.com",  
            "https://www.geeksforgeeks.org/"  
            ]  
  
print(New_predict)  
  
['good' 'bad' 'good' 'bad' 'good']
```

**Fig.5. Output Results of proposed algorithm**

## 5. CONCLUSION

For the detection of various phishing websites logistic regression methodology which is used to obtain accuracy. The algorithm of classification analysis the performance based on the literature review is proof to give 97% accuracy. Hence logistic regression is selected for the analysis which performs better when compared to the Decision Tree algorithm. Logistic Regression has an accuracy of around 96 to 97 percentage and also it is time saving when compared to Decision Tree algorithm. Decision Tree algorithm takes up a time more than five minutes to result the output whereas in Logistic regression it only takes very few seconds. Thus, it can be summarized that better algorithm is chosen and the experimentally successful technique in detecting phishing website.

## 6. REFERENCE

- [1] Zuochao Dou, Issa Khalil, Abdallah Khreishah, Ala Al-Fuqaha, Mohsen Guizan, "Systematization of Knowledge (SoK): A Systematic Review of Software Based Web Phishing Detection", IEEE Communications Surveys & Tutorials, Volume: 19 , Issue: 4, pp: 2797 – 2819, Sep. 2017.
- [2] JIAN MAO, WENQIAN TIAN, PEI LI, TAO WEI, ZHENKAI LIANG, "Phishing-Alarm: Robust and Efficient Phishing Detection via Page Component Similarity", IEEE Access vol. 5, pp: 17020 – 17030, 2017.
- [3] Rami M. Mohammad, Fadi Thabtah, Mee McCluskey , "Intelligent rule based phishing websites classification" IET Information Security, 2014, vol. 8, Issue. 3, pp. 153–160, 2013.
- [4] Kumar, V.D.A., Kumar, V.D.A., Malathi, S, **Vengatesan.K**, Ramakrishnan.M "Facial Recognition System for Suspect Identification Using a Surveillance Camera" **Pattern Recognition and Image Analysis** ,July 2018, Volume 28, Issue 3, pp 410–420
- [5] Vengatesan K., Kumar A., Chavan V.T., Wani S.M., Singhal A., Sayyad S. (2020) Simple Task Implementation of Swarm Robotics in Underwater. In: Hemanth D., Kumar V., Malathi S., Castillo O., Patrut B. (eds) Emerging Trends in Computing and Expert Technology. COMET 2019. Lecture Notes on Data Engineering and Communications Technologies, vol 35
- [6] H.Kim, J.H. Ruh,"Detecting DNS-poisoning-based phishing attacks from their network performance characteristics", Eletronics Letters, vol.47, issue 11, pp:656-658, 2011.
- [7] Ajaya Neupane, Nitesh Saxena, Jose O Maximo, and Rajesh Kana, "Neural Markers of Cybersecurity: An fMRI Study of Phishing and Malware Warnings", IEEE Transactions on Information Forensics and Security, vol.11, issue 9, sep.2016.

- [8] Y. Cao, W. Han, and Y. Le, “Anti-phishing based on automated individual white-list,” in Proceedings of 4<sup>th</sup> ACM Workshop on Digital Identity Management, 2008, pp. 51–60.
- [9] E. Saravana Kumar, **K.Vengatesan**, R. P. Singh, C.Rajan,” Biclustering of Gene Expression data using Biclustering Iterative Signature Algorithm and Biclustering Coherent Column, International Journal of Biomedical Engineering and Technology, vol.26, issue3-4,pp. 341-352, 2018.
- [10] **K.Vengatesan**,R.P.Singh, Mahajan S. B , Sanjeevikumar P,Paper entitled “Statistical Analysis of Gene Expression data using Biclustering Coherent Column” International Journal of Pure and Applied Mathematics , Volume 114 No. 9 2017, 447-454  
B.Narmadha, M.Ramkumar, **K.Vengatesan**, M.Srinivasan, "Household Safety based on IOT", International Journal of Engineering Development and Research (IJEDR), ISSN:2321-9939, Volume.5, Issue 4, pp.1485-1492, December 2017.