# Global Journal of Advance Engineering Technologies and Sciences

## FEATURE EXTRACTION FOR MALICIOUS URL DETECTION IN DATA MINING

**Mr. Jadhav Bharat S.[1], Dr. Gumaste S.V.[2]**

[1] M.E. student Department of Computer Engineering, Sharadchandra Pawar College of Engineering, Dumbarwadi, Outr, Pune, Maharashtra, India

[2] Associate Professor Department of Computer Engineering, Sharadchandra Pawar College of Engineering, Dumbarwadi , Outr, Pune, Maharashtra, India

[1] bharatjadhav754@gmail.com , [2] svgumaste@gmail.com,

## ABSTRACT

In this paper we have discussed the section of feature extraction with cost sensitivity. In This Section of Feature Extraction We have proposed different features for detecting Whether the URL is malicious or not . The different features are:- Length of the URL, Number of full-stops in the URL, TTL of the URL, and Get Info which gives information about the Registrar of the URL. Date of the URL, the malicious URL discovery framework uses genuine dataset .and with the help of the above feature extraction we will be able to detect whether the URL is Malicious or not.

*Keywords*— TTL, Cost Sensitivity. Cost sensitive Classification , online anomaly detection, online learning

## INTRODUCTION

This paper is related to Malicious URL detection with the help of cost sensitivity, section of feature extraction.

Malicious URL is detected on the basis of The features which we have set. They are Length of the URL, then number of dots present in the URL, then on the basis of Time to live (TTL). And on the basis of the date on which it was created and about the information about the registrar of that malicious URL.

## COST SENSITIVE CLASSIFICATION:-

Class considers classification of particular incorrect sorting. A price grid the regulation of requesting case encodes from one class as a substitute. Many real-world classification problems, such as fraud detection and medical diagnosis, are naturally cost-sensitive . to solve such tasks researchers have proposed several cost sensitive metrics . The Famous examples include the weighted sum of sensitivity and specificity and the weighted misclassification cost that takes cost into consideration when measuring classification performance.

## ANOMALY DETECTION

Anomaly detection is additional and can say that outlier detection on the other hand interest acknowledgment. The Goal of abnormality recognition is to find surprising information designs which failed to notice with normal patterns. Anomaly detection has been considered broadly from most recent many of years. In previous task, innovation recognition in semi supervised setting is naturally solved by reducing to a binary classification issue. An identifier which has desired false positive rate can be achieved by reducing it into Neyman-Pearson classification. Interestingly of inductive technique, semi-supervised novelty detection (SSND) concedes finders that are ideal despite of the circulation on novelties. In curiosity identification, there is a considerable blow on the imaginary properties of the choice principle of unlabeled information.

Although Anomaly detection is well studied from many years out still it remains a difficult task still today. it is due to several reasons .first is it is often a highly class-imbalanced learning problem as the number of anomalies is considerably smaller than that of normal patterns, which brings a critical challenge to many schemes using regular classification techniques

Secondly it is very costly to collect the labeled data mainly the positive training data ("anomalies"), which reduces the application of some classical supervised classification approaches. Moreover, in a real-world application, data usually arrives in a sequential fashion and the size of data patterns can be very large, leading to a big challenge for developing efficient and scalable algorithms for anomaly detection.

## MALICIOUS URL DETECTION:

In the Malicious URL detection is used to detect Malicious URLs consequently or semi-normally, which has been generally inspected in web and data mining groups for quite a long time when all is said in done, which is segment the current work

into two classifications: (i) non-machine learning procedures, for instance, blacklisting or principle based; and (ii) machine learning methodologies. The non-machine learning strategies for the most part experience the ill effects of poor generalization to new malicious URLs and masked spiteful patterns In the captivating after, this is focus on investigating necessary related work utilizing machine learning techniques. In writing, an collection of machine learning plans have been proposed for malicious URL detection, which can be assembled into two characterization: (i) regular batch machine learning systems , and (ii) online learning techniques Most of the existing malicious URL recognition techniques utilize customary regular batch classification methods to learn a classification model (classifier) from a preparing information set of named examples and after that applies the model to classify a test/unremarkable case. With everything, the categorization problem can be formed as either binary classification (normal vs. abnormal) or multi-class classification (accepting typical examples originate from numerous classes). In literature, a variety of classification systems have been connected, such as Support Vector Machines (SVM) , Logistic Regression , most extreme entropy Naive Bayes, and so forth. moreover, these calculations commonly require collecting and storing all the preparation occasions ahead of time and manufacture the models in a batch learning fashion, which are both time and memory wasteful and experiences extremely costly retraining taken a toll at whatever point any new preparing data arrives. Not at all like the clump machine learning calculations,, online Learning  has been recently proposed as a scalable way to deal with handling hugh scale online malicious URL recognition tasks .In general, online learning systems are more suitable for huge scale, genuine online web applications in light of the fact that their high ability and flexibility. In any case, most of the past online learning algorithms were intended to upgrade the order exactness, normally expecting fundamental preparing distribution. This un-mistakenly unseemly for online malicious URLs recognition undertakings since this present reality URL information delivery is regularly profoundly class-imbalanced, i.e., the quantity of malicious URLs is typically fundamentally littler than the quantity of amiable URLs on the WWW. Along these

lines, it is imperative to consider this issue at the point when outlining a machine learning and data mining algorithm for understanding a functional URL identification task. Finally, all the existing learning methodologies ordinarily need to mark a authentically considerable measure of preparing cases with a specific end goal to manufacture a sufficiently great grouping model, which is impractical as the naming expense is regularly wasteful in a authentic word application. This consequently drives us to study a brought together learning plan, which not just has the capacity minimize the naming expense, additionally augment the perceptive execution with the given measure of marked preparing examples.

## SYSTEM ARCHITECTURE

Fig.1 :-We have shown that How we have used the real time features like Length of the URL, Time to live of the URL , number of Full-stops present in the URL, and get info which gives information of the Registrar of the URL, for detecting whether the  URL is malicious or not.
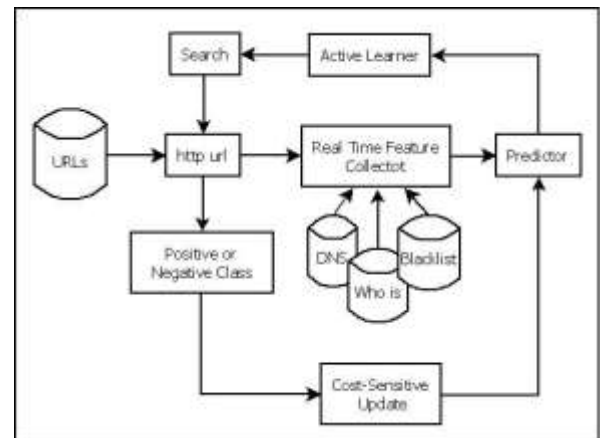


*Fig 1: System Architecture*

## IMPLEMENTATION

In this paper we have proposed a system For detecting the malicious URL. To practically implement the proposed System and detect malicious URL we have taken total 700 website URL   in that 500 websites URL are Real website URL i.e. without any malicious data in it. And remaining 200 website URL are malicious ones. but they are  randomly placed in our  dataset .so to detect which website URL is malicious URL , we have proposed  different features of extraction they are length , Number of dots ,TTL, get info.

*Length***:** On the basis of the length of the website URL we can detect that the URL is real or malicious one

*Number of full-stop:-* We can also detect the URL is malicious or not on the basis of dots present in the whole URL

*TTL*-: It stands for Time to live, Time-to-live (TTL) is a value in an Internet Protocol (IP) packet that tells a network router whether or not the packet has been in the network too long and should be discarded. For a number of reasons, packets may not get delivered to their destination in a reasonable length of time.

It is one of the important parameter for every website URL, We can detect The URL is malicious or not even with the help of TTL

Practical implementation of TTL is given below. We have use 'Ping' to implement TTL, Ping is use to connect one machine to another, Ping is a computer network administration software utility used to test the reach ability of a host on an Internet Protocol (IP) network and to measure the round-trip time for messages sent from the originating host to a destination computer and back it uses the ICMP protocol which has been created to check IP connectivity and get information about other machines in an IP network.ICMP is summed up in an IP packet, but is considered part of the IP or Internet.

Ping functioning:-Ping sends very small packets to an IP host who will n of which we can detect the URL is Malicious or not .

*Date* is also one of the feature for detecting the URL is Malicious or not . Date on which the Website URL was launched can help in detecting the website is malicious or not

*Who is Connect*: - It gives the information about the server i.e. when it is registered date of the registration name of the registrar, i.e. all kinds of information of website

## CONCLUSION

IN this paper we have successfully implemented the different features of extractions like length of the URL, Number of Full-stops in the URL, TTL Date and Get info for detecting malicious URL. And with the help of the above features we are successfully able to detect Whether the URL is malicious or not.

## REFERENCES

[1] Don Lancaster, TTL Cookbook, Howard W. Sams and Co., Indianapolis, 1975, ISBN 0-672-21035-5Jump up^ RCA COSMOS

[2] The Engineering Staff, The TTL Data Book for Design Engineers, 1st Ed., Texas Instruments, Dallas Texas, 1973, no ISBN, pages 59 , 87

[3] Paul Horowitz and Winfield Hill, The Art of Electronics 2nd Ed. Cambridge University Press, Cambridge, 1989 ISBN 0-521-37095-7 table 9.1 page 570

[4] MA, J., SAUL, L. K., SAVAGE, S., AND VOELKER, G. M. 2009a. Beyond blacklists: Learning to detect malicious web sites from suspicious URLs. In Proceedings of the SIGKDD Conference.

[5] Salus, Peter (1994). A Quarter Century of UNIX. Addison-Wesley. ISBN 0-201-54777-5.

[6] "RFC 1122 - Requirements for Internet Hosts -- Communication Layers". p. 42. Retrieved 2012-03-19. Every host MUST implement an ICMP Echo server function that receives Echo Requests and sends corresponding Echo Replies.

[7] "ICMP: Internet Control Message Protocol". repo.hackerzvoice.net. January 13, 2000. Retrieved December 4, 2014.

[8] Mike Muuss. "The Story of the PING Program". Adelphi, MD, USA: U.S. Army Research Laboratory. Archived from the original on 8 September 2010. Retrieved 8 September 2010. I named it after the sound that a sonar makes, inspired by the whole principle of echo-location.

[9] P. Domingos, P. Metacost, "Metacost: a general method for making classifiers cost sensitive, in: Advances in Neural Networks", International Journal of Pattern Recognition and Artificial Intelligence, San Diego,CA, 1999, pp. 155-164.

[10] N. Abe, B. Zadrozny, J. Langford, "An iterative method for multiclasscost-sensitive learning", in: Proceedings of the tenth ACN SIGKDD International Conference on Knowledge Discovery and Data Mining, Seattle, WA, August 2004, pp. 3

[11] F. Rosenblatt. "The perceptron: A probabilistic model for information storage and organization in the brain".Psychological Review, 65:386-407, 1958

[12] K. Crammer and Y. Singer. Ultraconservative online algorithms for multiclass problems. JMLR, 3:951-991,2003.

[13] N. Cesa-Bianchi, A. Conconi, and C. Gentile. "On the generalization ability of on-line learning algorithms".IEEE Trans. on Inf. Theory, 50(9):2050-2057, 2004

[14] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz,and Y. Singer. "Online passive-aggressive algorithms".JMLR, 7:551-585, 2006.

[15] P. Zhao, S. C. H. Hoi, and R. Jin."Double updating online learning". Journal of Machine Learning Research, 12:1587-1615, 2011

[16] Michael Attig, Sarang Dharmapurikar, and John Lockwood. Implementation Results of Bloom Filters for String Matching. In Proceedings of: IEEE Symposium on Field-Programmable ustom Computing Machines (FCCM), Napa, CA, April 20-23, 2004.

[17] Z. Baker and V. Prasanna. Automatic Synthesis of Effcient Intrusion Detection Systems on FPGAs.In Proceedings of FPL'04, 2004.

[18] Z. Baker and V. Prasanna. Time and Area Effcient Pattern Matching on FPGAs. In Proc. of FPGA '04. 2004.[4] Dollas, A. et al., Architecture and Applications of PLATO, Reconfigurable Active Network Platform. 1999, Department of ECE, Technical University of Crete: Greece.

Pimaplwandi as lecturer in Computer Technology Department,.Now he is currently working in Tikona Digital Networks as Network Support Engg. Also he is pursuing Master Of Engineering in Sharadchandra Pawar College of Engineering, Dumbarwadi,Otur, University Of Pune .

**Dr. S.V.Gumaste**, currently working as Professor and Head, Department of Computer Engineering, SPCOE-Dumberwadi, Otur. Graduated from BLDE Association's College of Engineering, Bijapur, Karnataka University, Dharwar in 1992 and completed Post- graduation in CSE from SGBAU, Amravati in 2007. Completed Ph.D (CSE) in Engineerng & Faculty at SGBAU, Amravati. Has around 22 years of Teaching Experience.

**Authors Profile :**

**Mr. Bharat S. Jadhav** received the BE degree in Information Technolgy from Pravara Rural Engineering College in 2012. During 2013-2014, he stayed at Late Hon.D.R.Kakade Polytechnic