

Website Reputation System

B. Amutha, Prabhav Gupta, Himanshu Kumar

Because of the fast development of the web, sites have turned into the interloper's principle target. As the quantity of web pages expands, the vindictive pages are likewise expanding and the assault is progressively turned out to be modern developing different ways to trick a client into visiting malicious websites extracting credential information. This paper presents a detailed account of ensemble based machine learning approach for URL classification. Models already existing either use outdated techniques or limited set of features in their attack detection model and thus leads to lower detection rate. But ensemble classifiers along with a selection of robust feature list for single and multi attack type detection outperform all the previous deployed techniques. Focus of the study is being able to come up with a system model that yields us better results with a higher accuracy rate.

Keywords: Data Mining, Classifiers, Multiclass, Layered approach, Multi-Classificat

I. INTRODUCTION

Present day online peer to peer networks continue to be the most popular way for the exchange of distinct pieces of information. When the distinct information gains some characteristics it is called as data. Over the recent years volume of data is continuously expanding as well and doesn't ask show any signs of slowing down. The growth of data in volume gives rise to it's heterogeneous nature which makes managing it so tedious and a daunting task. Rigorous and excessive demands for making sense from this data, such that someone can gain actionable insight has been steeply increasing as well. Hence this gives rise in demand for producing and extracting the noteworthy knowledge by any means possible. The bridge of communication is becoming more and more obscure for the consumer's understanding. This obscurity occurs in the form of data manipulation during the communication phase, clearly implying precipitous expansion of data has lead to rise of attack vectors that are associated to web.

Giving assurance for a safe communication by deploying systems equipped with procedural and comprehensive security measures becomes a must for the vendors. But securing the client pool can be devastating when people well versed in the niche of black hat hacking crack the security mechanisms with no effort whatsoever. For this attacker

the most lucrative way to spread his mass attack is through web. The most potent way to do so is by the means of "URL" which is the global address of documents[1] and resources matching your search query.

These URLs contain malicious content in them, making the timely detection of attacks a necessity before there is a potential breach or a potential threat in the system. To tackle this problem the best possible arrangement is to create a machine learning model that can succeed in dealing with the pitfalls from the past strategies that were being used.

Vindictive site pages are those, which contain content that can be utilized by assailants to misuse end-clients user visiting these web pages on a daily basis. This incorporates website pages with phishing URLs, spam URLs, JavaScript malware contents, Ad-ware, and some more. These days, it is winding up extremely hard to distinguish such vulnerabilities due to the persistent improvement of new methods for conveying out such assaults. Furthermore, not every one of the clients know about the diverse sort of adventures which aggressors can exploit of and also it is hard for a user with less technical knowledge to observe such harmful contents on a glance. In this manner, when there is a defenselessness in a site page, which the client doesn't know about, this instrument will enable him to remain safe regardless of his absence of information of the site. Also, if the URL itself has indications of being a phishing URL, at that point the client will be shielded from that site.

AI is a field of software engineering which is connected in pretty much every zone like science, liquid elements, agribusiness nowadays. It is conceivable to do so in light of the fact that AI systems gain from the information accessible and after that plan a well-suited model to aid choice making by handling the future information. What's more, the information need not be from a particular field, it tends to be anything. This flexible nature of AI is the explanation behind it to end up well known in the current world. In this work, we have utilized calculated relapse utilizing bi gram as the model in managed AI to test the exactness acquired to characterize if a given URL is pernicious or not. However, the disadvantage of AI is that the highlights which alone decides the outcome should be picked physically.

II. BACKGROUND

Proficient hackers can easily bypass existing strategies and evade the security system to a point where users privacy might be at risk. Personal risks over the internet are rampantly spreading and the attacks are carried out by

Revised Manuscript Received on September 22, 2019.

B. Amutha, is with faculty of Computer Science, SRM Institute of Science and Technology, Chennai, India.

Prabhav Gupta, is with the Department of Computer Science, SRM Institute of Science and Technology, Chennai, India.

Himanshu Kumar, is with the Department of Computer Science, SRM Institute of Science and Technology, Chennai, India.

redirecting "URLs". URLs hold contents ranging from confidential files to exploit kits. A page that seems to be totally safe might hold a FUD (Fully Undetectable exe) file attached to it which on execution is responsible for silently running drive by download attacks.[2] This never ending range of content held by URLs creates a inconsistency between a malicious identity and a legitimate identity. To handle the complexities many computational techniques have been devised in the recent years. Few of them being:-

Using a blacklist which basically is a repository of all the malicious sources reported and now labeled under the malicious class. While using blacklisting for URL detection there were very low false positive. But extensive usage of blacklist, turned out to be exhaustive since the collection needed to be updated continuously.

Signature based techniques and conceptually simple blacklisting execution in VM box were used which monitored anomalous state changes by blocking them with high accuracy.

Current approaches create a nuance for enterprises and vendors. Since a wide variety of attacks are platform agnostic making the timely detection of paramount importance. Over the recent times Machine learning has emerged as a prominent solution for solving Info Sec problems ,few of the use cases seen over the recent ones being deployed in malware analysis and network intrusion detectors.Hardly a few years back the process of extracting features from static and dynamic analysis of malware required the engagement from the user but that same task doesn't anymore require attention to every minute detail.Training the Machine Learning Model allows for smoother and efficient string extraction of keywords from a malicious executable. making the very laborious task of classification very easy and fast to be analyzed such that insightful decisions at the time of a security breach.Thus machine learning seems to be the perfect resort for tackling these constantly developing attacks over never ending attack surfaces on the web due to 2 main reasons:-

Machine Learning is based on feature extraction as discussed in the next section which just requires a small labeled data set and gives considerably good results at the time of testing.

Zero days couldn't be detected with approaches and methodologies that existed previously.Hence a model with correct feature representation can help detect URLs not known to the blacklist and overcome the shortcomings.

III. ENHANCED FEATURE EXTRACTION

Performance and effectiveness of the proposed system can only be at par when the feature extraction has been done right and in a way that minimizes the noise from the data set. With a multi-labeled classification method detection of URLs become more accurate since it covers all the possibilities for attack identification which includes parameters like platform,vector, and in some cases if the attack has already taken place , generating a detailed log file from that time will help in boosting the accuracy rate by a significant margin.

The remarkable and noteworthy property of the model being able to use the identical data set for classifying the attack type and URL classification.

A. Lexical Features

Distinguished patterns can be observed in the URLs of phishing attacks.[3] A URL might be delimited by special characters such as ., /, ?, =, -, '@' and the main trait of the URL string being Domain name closely resembles to a vendor having a wide audience. So a difference of insignificant character such as a '.' can end up redirecting the innocent user to page embedded with a form grabber in it , such that users every activity from logging in to his keystrokes are now visible to the person with a wrong intent.

B. Content based Features

Progressing rapidly web technologies which have made daily consumer life effortless by online banking systems for wire payments, bank payments or even ticket booking. But this leads to a whole new attack surface where the hacker can either jack your session by session jammers or injecting the malicious code, later hiding it in webpage content.[4] Hence these set of features become of prime importance to be extracted since they run on the client side code of the web.

No.	Feature
1	Domain token count
2	domain length
3	presence of special characters
4	Brand name presence

C. Host based features

This class of features is related to the DNS of features related to the domain of the URL.[5]Usually spammers come from a small pool of autonomous systems. These features are the most stable to identify scammers since it becomes tough for them to continuously get a white listed AS and hold malicious content on it.

Feature	Description
IP Address	AS number, MX records
WHOIS information	Registrar info
Geographical	Country, city, zip,
Connection speed	Response time

D. Additive features

All these features above can be used for communication over the web despite of protocol taken into account. Broader picture from the web infers that 90 percent traffic is encrypted i.e. travelling through a secure channel using HTTPS protocol where S stands for secure[6]. Not only this, such classification can help in extracting issuer name,age,expiry date,days to expire etc which can help in the better execution of the model



and even assist in practical real time encrypted traffic analysis helping internet services detect frauds in real time.

IV. MACHINE LEARNING BASED APPROACHES

Classifiers run from cutting edge AI strategies to straightforward perceptrons. All ordering calculations share a typical structure dependent on the structure of the perceptron, yet vary in feed forward spread and back propagation learning calculations. While the structure of these straightforward counterfeit neural systems is shared, the learning instrument speaks to a huge contrast and incredibly influences the adequacy of the calculation. Order calculations which just utilize the lexical highlights are known as lightweight characterization calculations and have a significant preferred standpoint over completely included characterization calculations as they don't require the utilization of outside data sources to make a forecast about the security of a URL. This is an imperative factor, particularly when due thought is given to execution on the customer's side of the procedure, as no extra inertness is presented through the execution of queries. Another positive factor in such manner is the utilization of a feed forward component of a neural system, which might be executed proficiently, in this way accelerating the procedure considerably.

A. Single attack type detection

Machine learning and other artificially intelligent learning based models has been utilized in many ways to deal with vindictive URLs to detect malicious web links and preventing data infection. Web spam is done when many of the artificially developed web pages are injected into the web to boost up the profit or sometimes created for redirecting the web traffic towards a single source. Content based spam analysis detected method was developed by Ntoulas et al. [7]. These single attack based methods uses the particular attack based features and model coded using these features have capabilities to detect proximal feature attack sources. Phishing websites mimic a trusted confider to gain access to client's private data and information. Versatile machine learning classifiers are developed to detect such pages and also Whittaker et al. [8] used extracted features to develop an model that detects these intruders using an machine learning based classifier and regularly updates Google's blacklisting sources. Content based features can also be used to develop a training data set and thus a signature based system was proposed by Hou et al. [9]. Proposed system can not only detect significant malicious content but also the obfuscation of DHTML malicious codes are detected more accurately.

B. Multi attack type detection

Identification of multi web attacks comprises of extracting various host and lexical based features and then the dataset requires cleaning to remove unlabeled URLs and thus strategic methods are utilised to categoize out the unstructured

variables[10]. All of these method used machine leaenring to tune thei models and are optimized using optimization techniques such as genetic stuctures and in the ultimatum finds out the lowest error giving maximum accuracy. After gathering the labeled URL dataset and combining the efforts

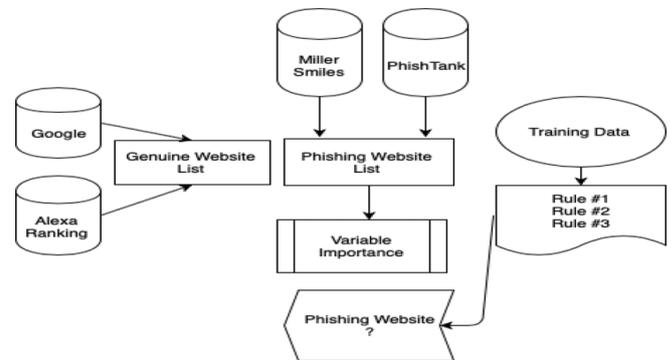


Fig. 1. Classification Workflow

from previous research works then integrated with experimental study done by us, helps extracting best possible features for URL classification. Some features provide extensive information about URL but highlighting some could turn out to be resource extensive task, whereas some have their own shortcomings ending up as noise in the dataset.

Implementing an ensemble based classifier have shown more reduction in errors and developing an idea that an ensemble learning uses combination of multiple trees in series and parallel learning modes can also lower the error rates of previous tree nodes reaching the end[11]. In this learning based method new instances are classified and categorized as malicious or benign.

Bayes Optimal Classifier

The Bayes Optimal Classifier is a course of action strategy. It is a gathering of the impressive number of theories in the hypothesis space. All things considered, no other gathering can outmaneuver it. Bayes Optimal Classifier is a version of this that acknowledge that the data is prohibitively self-ruling on the class and makes the estimation progressively achievable [12]. Each hypothesis is given a vote relating to the likelihood that the planning informational collection would be inspected from a system if that theory were legitimate.

Bootstrap aggregating

Bootstrap accumulating, regularly condensed as stowing, includes having each model in the outfit vote with equivalent weight. So as to advance model change, bagging trains each model in the outfit utilizing an arbitrarily drawn subset of the preparation set.

For instance, the irregular backwoods calculation joins arbitrary choice trees with stowing to accomplish extremely high characterization precision.

Bagging

Boosting includes gradually assembling a group via preparing each new model case to stress the preparation occasions that past models mis-arranged. Sometimes, boosting has been appeared to yield preferable exact-ness over packing, yet it likewise will in general be bound to over-fit the preparation information. By a long shot, the most widely recognized usage of boosting is Ada support, albeit some more current calculations are accounted for to accomplish better outcomes.

Bayesian model combination

Bayesian model blend (BMC) is an algorithmic acclimation to Bayesian model averaging (BMA). As opposed to inspecting each model in the outfit solely, it tests from the space of possible gatherings. This change beats the penchant of BMA to join toward giving most of the weight to a single model and giving better results. The results from BMC have been seemed, by all accounts, to be better generally speaking (with genuine significance) than BMA, and packing. The inclusion of Bayes' law to process loads requires figuring the probability of the data given by each model. Usually, none of the models in the gathering taking in are really the scattering from which the planning data were made, so all of them precisely get an almost zero term.

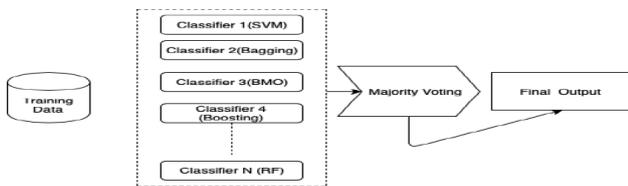


Fig. 2. Ensemble Learning Workflow

C. Noisy Data

The nature of the preparation information significantly impacts the nature of the forecasts. Assessing the preparation information quality is troublesome, yet it constantly includes precise naming of preparing occurrences and a preparation set that is illustrative of the cases 109 the classifier is approached to name. We are regarding this issue as a paired arrangement issue in which an area can either be malignant or obscure. The marks on the vindictive informational collection are profoundly precise: it is impossible that in excess of a couple of the areas that showed up on the boycotts were false positives as these boycotts are curated by capable security network individuals and have seen across the board selection.

There are a sizable number of malignant spaces in the obscure set that are never boycotted. Weeding these out from the really enlisted spaces is troublesome, best case scenario. at the point when a space is first enrolled there is a shortage of data about that area's notoriety. Likewise, we are as often as possible managing the spaces far down the long-tail of an area prevalence bend. That is, well known spaces are

moderately steady and show up with less recurrence than recently enrolled areas in the day by day DNS diff we compute. The sorts of areas that update their DNS information additionally contrasts radically with every preview, as one may expect given that the standard deviation in the quantity of changes is the greater part the normal.

The number what's more, sort of DNS changes every day fluctuates because of mass DNS exchanges, area sniffing (area name hypothesis) and occasional cleansing of lapsed spaces. To relieve the way that action in a DNS zone can change significantly between ages, we chose to utilize the arrangement of obscure spaces from the past age as preparing information for the present age. The subsequent preparing set is without a doubt uproarious yet, the boycotted set contains few or no bogus positives, enabling us to distinguish a huge number of boycotted spaces far ahead of time of the time they are boycotted. So as to investigate how frequently one needs to retrain the classifier so as to be viable, we took a gander at the consequence of preparing on obscure information from 1 age earlier presents an outline of the outcomes. Location of assaults debases as the preparation set.

Not with standing, the span of the competitor set stays unaffected. Accordingly, we choose to retrain the classifier each age. Before outlining our outcomes, we see that the estimation of our framework is very identified with the structure of the applicant set. On the off chance that the obscure areas added to the hopeful set are generally false positives, for example they were malevolent spaces that were not boycotted, Prospector is of sketchy esteem. Accordingly, we next look at the synthesis of the competitor set.

The last stage includes assessing classifier execution in information excluded in the preparation set by any stretch of the imagination. Practically speaking, classifier execution once in a while achieves the dimensions shown amid the preparation stage and this is the reason it is basic to have an unmistakable set of information to treat as a test set. Countless have been proposed, each of which depends on an alternate hypothetical premise. The accompanying segments are planned to clarify in some detail the techniques by which every classifier works.

V. RESULTS AND DISCUSSION

A. Error Analysis

False positives: A false positive occurs when a benign URL falls under malicious class. According to the model proposed based on feature extraction this is most likely to occur when the Domain or the Domain Registrar is not recognized or a URL containing the public ip can also be wrongly classified. Some more cases arise when the registrar ends the certificate issued, existing for long times site won't be able to hold content lowering the rank of website in the system.



False negatives: False negatives occur when a malicious URL goes undetected. Majority of the false negatives are those links which have very high resemblance with the domain name of a famous target company or the one's using '//' in their URL which makes the first link redundant. So, for establishing the grounds for a better detection system in general metrics associated with features should also be evaluated.

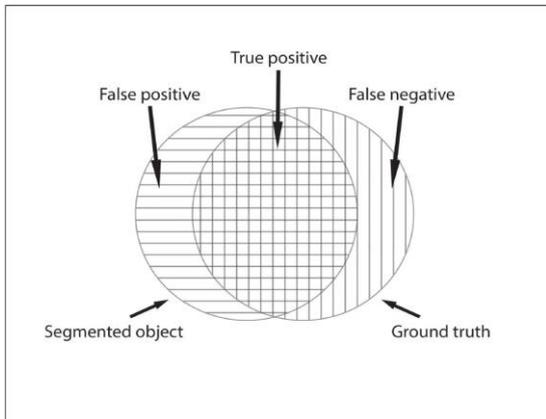


Fig. 3. Error Metrics Analysis B. Proposed

B. Proposed system drawbacks

URL obfuscation: Attackers and generation algorithms produce domain names having random string values and obfuscating those character codes which were the prime reason for their detection as a malicious content[14] on the web. Also referred to as hyperlink trick this link lands the user on a page or a spoof website which can be a gateway for malware to enter in the consumers system. Even though it is an anti-hacking procedure attackers tend to use it for their advantage

Javascript obfuscation: Malicious javascript on the web is the most dangerous piece of code and obfuscating it further leads to impossible detection. Particularly, this kind of obfuscation is used to embed a piece of malicious code that makes the static and dynamic malware analysis beyond the scope of preventing the attack in a real time system[15]. Even Machine learning solutions in this field of InfoSec world are yet to show considerable results.

```
<html><body><textarea id="KSHD29FC" style="display:none;">115,52,81,71,66,64,50,
42,114,83,93,102,43,40,51,119,76,34,111,123,74,75,97,96,90,32,44,45,37,104,57,58,94,116,
82,108,87,112,85,110,117,86,121,103,124,63,79,84,62,38,89,88,105,91,35,126,68,65,69,55,
118,100,125,53,107,92,113,47,54,106,122,99,95,39,72,77,120,61,78,70,73,98,36,46,59,67,
80,49,41,56,48,109,60,101,33</textarea>
<div style="display:none;" id="KSHD29FCa">+
%5D%60H%20sG8jg%3AWALPT1X210G@Q9%3CKExH%5DP%3D+%3Dn%26%2CM-Hixd%
3A3_W/eCGZM%5D%7D%209jm6%7C9G%2CjT^HGZ%7E%5DLw2%20/Gt%3CK9O%3B
/H%3C%3Do%21TQ%3DC%3Dy%7D_9l%3E-lc%7E%23G%24F%5DkSsu%60EP%216
...
%5EdG219R6zwted%5EPxEga%5DIV+%23%22-c%29xo%58t</div>
<script>-function skfjAd08djs(Cd9fs){return String["fromCharCode"][(Cd9fs)];}</script>
<script>-function xh7z8ckkMGe(FFzD54hch5){var Cc9yP8XdVx8M=0,rcA29Zt,
TbxK0cw=FFzD54hch5.length,MMsk77R3=0,Eki59FcVW="";
Cc9yP8XdVx8M=Cc9yP8XdVx8M++;if(!MMsk77R3)
{eval["rcA29Zt"+"t=125^FFzD54hch5.cha"+"rCo"+"deAt(Cc9yP8XdVx8M%62"+"55)"];
eval("Eki59FcVW"+"="+skfjAd08djs(Ceup1xb(125^"+"rcA29Zt)+"32)");
MMsk77R3=2;}else{MMsk77R3--;}return (Eki59FcVW);}var
li2nzK=window["unescape"]+"ape";(document.getElementById("KSHD29FCa").innerHTML);
eval(xh7z8ckkMGe(li2nzK));</script></body></html>
```

Fig. 4. Obfuscated Java script Codes Once

the these algorithms have been tried, the devices referenced in Applications will be actualized. The principal will be a program module for the Google Chrome browser. The second will be an intermediary module pursued by an extension add on. These executions will be refreshed through a focal server that will inquiry outside information sources intermittently for new preparing information with which to refresh its grouping models. At long last, an apparatus will be built up that will empower these classifiers to be sent as modules for a criminological occurrence investigation system.

C. Classification result

There exist two distinct conventions on which the order task is performed exclusively. Both the datasets were isolated into an ideal 70/30 split so as to choose the best fit and abstaining from underfitting or overfitting. K-crease cross approval technique is utilized to prepare every one of the troupes and network scan is utilized for advancing the hyper parameters for best model execution.

Table I Ensemble Training Accuracy (Https V/S Http)

Classifier	HTTPS	HTTP
Boosted Trees	97.4%	91.7%
Bagged Trees	97.5%	91.1%
Subspace Discriminate	93.3%	89.6%
Subspace KNN	80.0%	79.4%
RUS Boost Trees	96.3%	92.1%

While preparing all the troupe classifiers Bagged trees demonstrated the most noteworthy execution on HTTPS convention in contrast with RUS helped trees indicating more noteworthy exactness on the HTTP convention. The dissipate plot for preparing information indicates cross as the misclassified perceptions and specks as the core.

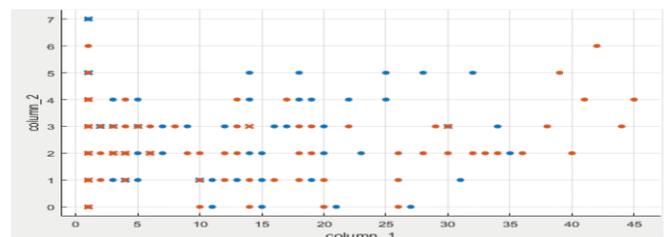


Fig. 5. Scatter Plot

VI. APPLICATIONS

One of the biggest favorable circumstances of the lightweight order of URLs is that it presents minimal overhead regarding dormancy and handling. While this reality has been referenced a few times in a great part of the work introduced on the subject what's more, the precedent utilization of customer side program modules has been referenced, there are different utilization's to which this work can be put. Right off the bat, a web



intermediary inside an association may utilize the classifier to channel malevolent URLs from being visited. This may be utilized on an independent premise or it tends to be utilized to increase the intermediary's current boycott of URLs. On the off chance that the intermediary gets a solicitation which the classifier predicts to be vindictive, it might deny the solicitation and add the URL to the boycott or another database framework for further investigation. Another utilization for these lightweight classifiers is that they might be utilized as an underlying first go of huge URL logs when attempting to break down a system occurrence. The URLs may then be hailed for further examination for example, data gathering from outer sources, for example, WHOIS information.

VII. CONCLUSION

Web attacks being one of the most prominent categories of threats are very important to be addressed and to develop the classifying architecture that can detect most of the malicious ones giving low false positives rate was in dire need. This paper gives the broad overview about the different range of attacks possible discussing how to extract most valuable features and using them to structure our machine learning model[16]. In addition, single and multi attack based detection techniques were used involving different algorithms that can work and overcome the earlier drawbacks[17] of the previous models. Most of the focus was given to the ensemble based learning approach[18] as verifying how a voting based ensemble classifier can detect the type of the attack and making websites more secured.

REFERENCES

- 1) Castillo, C. Donato, D. Gionis, A. Murdock, V. Andsilvestri, "Know your neighbors: web spam detection using the web topology", ACM SIGIR: Proceedings of the conference on Research and development in Information Retrieval (2007).
- 2) CISCOIRONPORT. Iron Port Web Reputation: Protect and defend against URL-based threat.
- 3) Fette, I. Sadeh, N. Antomasic, "Learning to detect phishing emails", WWW: Proceedings of the international conference on World Wide Web (2007)
- 4) GY Ongyi, Z. Andgarcia-Molina, "Link spam alliances", VLDB: Proceedings of the international conference on Very Large Data Bases
- 5) MA. J. Saul, L. K. Savage, S. Andvoelker, "Identifying suspicious URLs: an application of large-scale online learning", ICML: Proceedings of the International Conference on Machine Learning (2009)
- 6) Provos, N., Mavrommatis, P. Rajab, M. Andmon Rose, "All your frames point to us", Security: Proceedings of the USENIX Security Symposium (2008).
- 7) Ntoulas, A. Najork, M. Manasse, M. Andfetterly, "Detecting spam web pages through content analysis", WWW: Proceedings of international conference on World Wide Web (2006).
- 8) Whittaker R, C. Ryner and B. Nazif, "Large-scale automatic classification of phishing pages", NDSS: Proceedings of the Symposium on Network and Distributed System Security (2010).
- 9) Houy. T. Chang, Y. Chen, T. Lai and M. Chen, "Malicious web content detection by machine learning", Expert Systems with Applications (2010), 5560.
- 10) A. L. Buczak and E. Guven, A survey of data mining and machine learning methods for cyber security intrusion detection, IEEE Communications Surveys & Tutorials, vol. 18, no. 2, pp. 1153-1176, 2016
- 11) Hu, H. Gao, Z. Li, and Y. Chen, Poster: Cud: crowd sourcing for url spam detection, in Proceedings of the 18th ACM conference on Computer and communications security.

- 12) Y. Yan, Q. Wu, M. Tan, M. K. Ng, H. Min, I. W. Tsang. Online Heterogeneous Transfer by Hedge Ensemble of Offline and Online Decisions. IEEE Transactions on Neural Networks and Learning Systems, 2017.
- 13) R. M. Karp, S. Shenker, C. H. Papadimitriou. A simple algorithm for finding frequent elements in streams and bags. ACM Transactions
- 14) S. C. Hoi, J. Wang, P. Zhao, and R. Jin, Online feature selection for mining big data, in Proceedings of 1st intl. Workshop on big data, streams and heterogeneous source mining: Algorithms, systems, programming models and applications.
- 15) S. Popescu, D. B. Prelicpean, and D. T. Gavrilut, A study on techniques for proactively identifying malicious urls, in 2015 17th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC).
- 16) International Joint Conference on Artificial Intelligence, 2001, Vol. 17
- 17) M. N. Feroz and S. Mengel, Examination of data, rule generation and detection of phishing urls using online logistic regression, in Big Data (Big Data), 2014 IEEE International Conference
- 18) R. McDonald, K. Hall, and G. Mann, Distributed training strategies for the structured perceptron, in Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics.