

Machine Learning Based Web Risk Detection

N.VenkataVinod Kumar¹, Senthilnathan Palaniappan^{2*}, Naresh K³, Anitha K⁴

¹ Dept of CSE, Annamacharya Institute of Technology and Sciences, Tirupati, Andhra Pradesh, India

^{2,3} School of Computer Science and Engineering, Vellore Institute of Technology, Vellore, India

⁴ Dept of CSE, Saveetha School of Engineering, Saveetha Institute of Medical and Technical Sciences, Chennai, India

Article Info

Volume 83

Page Number: 5048 - 5058

Publication Issue:

March - April 2020

Abstract

Web applications have increasingly become one of the customary platforms for service releases and representing data and information over the worldwide web. And thus, many security susceptibilities have controlled to various types of attacks in web applications. This paper is thus intended to look into the use of supervised-machine learning techniques which include; genetic algorithms and support-vector machines, as they are used in detecting some of the key web application layer threats. Some of the most common application layer web threats include; remote file inclusion attacks, SQL injections, and cross-site scripting.

As the internet keeps growing each other day, it has become very important to detect the web threats as well as leveraging the powers of machine learning which is one of the various prospective methods in order to make the detection more effective. We will look into how we would use genetic algorithm and support-vector machines to detect the above-mentioned threats in the application layer due to the millions of data requests send every second. From this information, I will come up with a conclusion where I will be able to state the effectiveness, viability and weaknesses of each of the key techniques from which support-vector machines proved to be more effective compared to genetic algorithm in terms of performance and viability

Article History

Article Received: 24 July 2019

Revised: 12 September 2019

Accepted: 15 February 2020

Publication: 27 March 2020

Index Terms; *SQL injection, RFI, XSS, GA, SVM, Parser, Gathering Test Data, Parsing the Requests*

I. INTRODUCTION

In the recent past, the level of technology has been growing very fast, especially the number of devices launched every day, and hence increasing the usage of the internet which these devices heavily rely on. The heavy usage of the internet and production of new devices connected to it has provided new environment for attacks in the web. The major issue which has brought about this severe attack is because the new internet-connected devices being developed do not take into consideration the numerous and very common attacks which should always be the first priority.

Therefore, to be able to fight these attacks and threats, we should not only focus on prevention as detecting these threats is equally important. For instance, there is a very severe attack which was

named “zero-day” which has caused serious damages for quite a long time. From all the attempts made by various people like Arisawa to prevent it, it was not possible to stop it. Therefore, for this particular instance, the best approach to handle it is first detecting it and then finding any possible ways of stopping it after it has been identified.

Talking about detection, there are various ways of detecting vulnerable threats in the application layer like manual creation of detection signatures which are made from known attacks in order to identify the threats which are similar to that particular nature. This method was used so many years ago and hence it cannot work in the current environment because it has a number of drawbacks (Arisawa, 2014). First of all, there are advanced threats which do not have

similar nature to the ones which were available in the early days.

The other issue is that the whole procedure of coming up with these signatures is very slow and it needs a previous attack of the same nature which acts as the reference. In addition, these signatures only detect the attacks which are of the same nature like the reference. And finally, it has a very risky assumption that the former attack will not change its nature and hence the assumption that it will be redetected using the created signature. Therefore, this brought the urgent necessity to use improved and more conventional ways of detecting attacks and this is when the people who were involved with the research came up with machine learning techniques.

1.2 Purpose of the Study

The major purpose of this study is thus to analyze the effectiveness and performance of the genetic algorithm as well as the latest method, support-vector machine (SVM). In addition, we will also look at how effective both methods are at detecting threats, how frequent each can detect a threat as well as how each one of them can misidentify a particular attack as either being false positive or the wrong attack type. By doing this, when a person is doing the tests, they should ensure that they have sufficient data to use as well very accurate for both algorithms. For this paper, I will provide the procedure of using both algorithms so that one can do their own tests to ascertain that they are efficient.

1.3 Web Threats

Web threats are any mischievous attack that make use of the internet as their main method of distribution and hence making them very wide and varied. These web threats can be separated into two major categories namely, push and pull. The push based attacks are the ones which use luring techniques to get a user to fall as a victim to the attack, a good example of this type is the phishing emails. On the other hand, the pull based threats are

attacks that can affect any visitor to the service or website (Oppliger, 2011).

The main aim of all these web attacks is to get access to confidential information which can be used by the people doing the attacks to ask for ransom or use that particular information to blackmail them (Oppliger, 2011). These attacks are thus increasing everyday as the number of internets connected devices produced increases.

When doing their attacks, the attacker aim at the most vulnerable targets, the weakest links as they majorly focus on public-facing applications as it is very simple. The application layer is very vulnerable to these attacks and some of the reason why most attackers focus on it is due to:

- i) Runtime overheads
- ii) Incomplete implementations
- iii) False positives and false negatives
- iv) Inherent limitations
- v) Complex frameworks

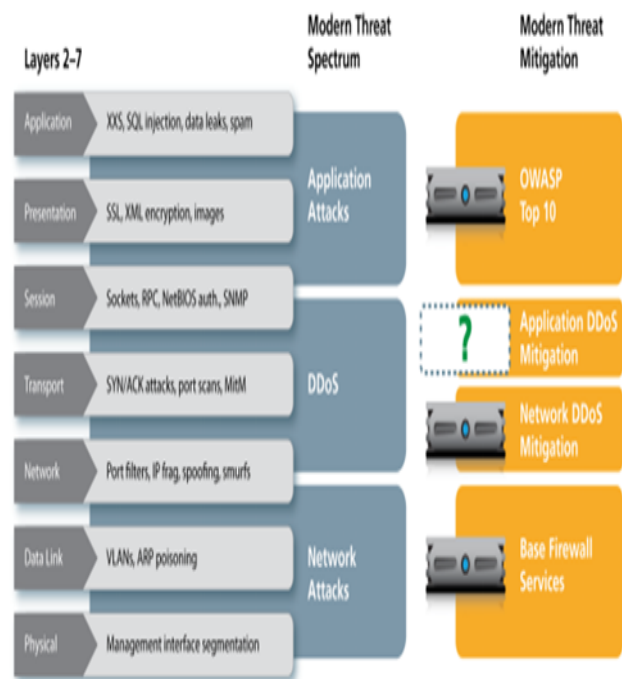


Fig1: How attacks occur in the application layer

II. CURRENT DETECTION AND PREVENTION METHODS

Web assaults is honestly harming to Associate in Nursing association from multiple points of read therefore identification and offsetting action is of high significance. These assaults and ensuing problems don't seem to be simply restricted to merely very little time associations additionally, with various high professional autoimmune disorder sites are becoming to be unfortunate casualties to expansive but comparable assaults too. Many investigations state over ninetieth of net applications are overwhelmed against SQL infusions alone [11]. Moreover, several net applications are overwhelmed as a result of their utilization of end-of-life programming. PHP is that the most prevailing programming idiom used for sites representing over eightieth of the piece of the pie, anyway of the sites that utilization PHP, ninety fifth of them are still on rendition five [12]. As of the beginning of 2017 PHP five isn't once more effectively upheld and it's prescribed to alter to rendition seven that is by and enormous effectively bolstered, refreshing can carry with it the bulk of the foremost recent security XSS that might create a web site effortlessly exploitable typically [13, 14]

2.1 Signature Based Detection

A traditional and basic methodology for recognizing security dangers is that the utilization of marks, anyway a substantial heap of those mark based mostly instruments are additional suited to the lower layers of the OSI show rather than the upper layers, as an example, the applying and introduction layers. These devices ar alluded to as Intrusion Detection Systems (IDS) and plenty of rely on customary articulations and alternative example coordinating procedures with marks created utilizing past assaults techniques, one case of such AN instrument is Snort [17]. during this approach, to that extent as there's a comfortable range of marks that cowl the expansive vary of conceivable assaults then the popularity approach are often terribly effective.

2.2 Modern Methods of Detection

With the top goal to defeat these difficulties forrecognizing internet dangers, some folks have projected a variedsuperimposed methodology can offer the simplest barrier. Such a framework wouldn't simply have varied layers of location together with marks and totally different ways however in addition input circles to refresh the insurance frameworks for increased future discovery. A multi-layered methodology would likewise have the capability to {deal with to handle} all dimensions of the system as opposition a solitary framework that may simply deal with a set of the layers, for instance, the system or application layers. this system wouldlikewiseempower components of the making ready to be focused and within the cloud whereas totallydifferentterritories are nearer to the endpoints. standard systems like mark based mostly identification would in any case be useful, but with the growth of additional ways, for instance, conduct examination are valuableasinternetassaultstogetherwith mammoth specification endeavors and not solely a solitary awful demand is closing additional typical place. Nonanal purpose is such a solution would take into thought worldwide joint effort to feature to ill fame records, whitelists, then forth to in addition tackle the difficulty of a developing risk as opposition having similar instruments sent in numerous regions and poor answers for inevitable updates [4].

A multi-layered methodology joins the simplest of the previous customary ways with new conceivably higher arrangements, and significantly additional curiously proposes a framework that may innately learn and enhance as a middle characteristic; basically, a similar because the procedure being stated during this examination.

III. PROPOSED METHODOLOGY

3.1 Machine Learning Approaches

In the recent past, machine learning and pattern-recognition performances have been adopted in

security applications like network intrusion detection, spam filtering and malware detection as they have the capability of potentially detecting fresh attacks and also generalizing. Genetic algorithm and support-vector machines are the only techniques which have proved to be more successful compared to the ancient way of creating signatures (Mitchel, 2017). However, support-vector machines have proved to be more efficient compared to the genetic algorithms.

Even though these learning techniques, especially support-vector machines assume stationary, in both the operational data and the data used to train the classifier as they are selected from the same distribution. However, this scenario can be avoided when intelligent and adversarial settings are applied as they can control the data hence hindering the stationary to feat the existing susceptibilities of the learning algorithms (Oppliger, 2011). This approach hence poses another problem which makes it doubtful whether the machine learning techniques can really be effective in security-sensitive tasks or they needed some more modification for them to be effective in that task. Some of the keys issues noted include: -

- i) Analyzing the key vulnerabilities of learning algorithms
- ii) Assessing their security by executing the corresponding attacks
- iii) Designing appropriate countermeasures.

These issues are thus being researched on in the newly emerging research area of adversarial machine learning, as it comes in as the connection between machine learning and computer security. Generally, machine learning is the process by which data is utilized to form a proposition that accomplishes better than an a priori theory designed without the data. Therefore, the large amount of data requires intelligent and very advanced ways to process it as well as big incentives driving the plea to do so. Due to this, through the utilization of machine learning techniques and patterns, more

advanced meaning can be extracted from the data which can thus be used to give analysis and thus drawing conclusions (Mitchel, 2017).

Due to the a lot of involvements from machine learning techniques in data mining, it has become one of the greatest ways of coming up with new data patterns and gaining meaningful information from it and hence predictions and this data is from the numerous web requests. This then give it an upper hand in in security applications such as web threat detection and hence giving the system an upper hand in improving threats detection and hence enhancing prevention (Michalski, Carbonell, Mitchell, 2013). The web applications which cannot detect threats in real time can only use the ancient ways of detecting attacks but they are so time consuming, and the attacks get complex with time and hence beating them.

Basically, the machine learning works in a much-esteemed way by pinpointing a series of features in the data set which is presented for investigation, these very specific features are thus the constraints used by these algorithms to learn and hence enabling them to make better decisions. In this case, the algorithms in machine language are not dictated on what to do but rather they can make their own decisions based on some certain criteria like performance.

3.2 Types of Machines Learning

In machine learning algorithms, there are two key distinctions, there can either be supervised learning or unsupervised learning. In supervised learning, the system in this case starts with a labelled set of data, and this data states or rather dictates what the end result will be and thus the training of this particular system can be very accurate provided that there are no interventions (Mitchel, 2017). The other type of machine learning, unsupervised learning, which is also called reinforcement learning, has some external informers which updates the system on how well its working as it progresses with the work.

Supervised learning has its own drawbacks which needs to be overcome for it to work perfectly. First of all, collecting the initial data set is a problem as it needs prior research of getting the informed sources which is a very trifling exercise. However, the features of informed sources can be obtained through the use of brute force method but there is an issue with using this method as it requires more computation time and the data can be noisy and hence making it to lose very key features in the process and this can lead to more issues (Michalski, Carbonell, Mitchell, 2013). So as to avoid the above-mentioned issues, which is making the size of the data to be small and also maintaining the ultimate performance of the system, a process called instance selection is more preferred.

On the other hand, when the system has so many features in the data set, this can increase the complexity of the whole system. To solve this particular problem is complexity, you need to remove all the irrelevant and redundant features where possible, but you should not remove the features which highly depend on one another as this may lead to errors and inaccuracy. To professional deal with this type of issues, you only need to make specific features, either by changing the existing ones or making new ones (Mitchell, 2017). To achieve this, you need to make sure that you select the right algorithms for the dataset and this will improve the performance of the system and hence making it more concise and accurate.

3.3 Genetic Algorithm

The genetic algorithm is a search-based algorithm that utilizes the machine learning techniques in order to trace the prime or near-optimal solutions for a certain problem. This is what is meant by search as used in the description above rather than using alternative algorithms like simulated annealing or local search. In generic algorithm, it is not all about searching for a specific data set but it is looking for the best answer possible for a particular problem as defined in the algorithm (Witten, Frank, Hall, 2017).

Therefore, the generic algorithm makes use of fitness functions and reward systems so as to be able to distinguish between the better solution and the ones which cannot proceed to the next level.

Just like generic development works in the real world and how species develop, genetic algorithm work in a similar way. Below are the steps followed by the generic algorithm:

- i) It starts in an initial population of individuals, in our case the dataset
- ii) In this case, only one dataset is a solution to the issue at hand
- iii) All the datasets are individually evaluated using some calculation which are made to solve the problem at hand
- iv) Our case is detecting web threats and hence we use the given formula

$$\text{fitness} = (\text{correct detections} / \alpha) - (\text{false positives} / \gamma) - (\text{incorrect detections} / \beta \cdot 8)$$

α The number of possible correct detections

β The number of possible incorrect detections

γ The number of possible false positives

In this algorithm, the fitness algorithm, used in generic algorithm, the correct detections will improve the fitness of a specific dataset and these individual datasets are used for later genetic operators. On the other hand, the incorrect detections and false positives impact the fitness negatively but the incorrect detections having a lesser negative impact than the false positives (Zhang, IGI, 2012). The genetic algorithm is currently being used in web threats detection through the use of the variants of an attack to detect network related attacks. But the main idea about in genetic algorithm is how to generate various individuals and testing if they can fit a certain criterion.

3.4 Support Vector Machines

The main technique used by support vector machine for categorizing data is dividing the data sets into two or more sets with prevalent margin imaginable between the separations which are known as hyperplanes. The reason as to why there is a very large separation between the two margins is to reduce the probabilities of categorization error. After the classifications, the points which lie within the margins are known as support vectors hence the name of the algorithm, support vector machine. The support vectors are thus the points which are used to calculate the hyperplanes and the other points are ignored.

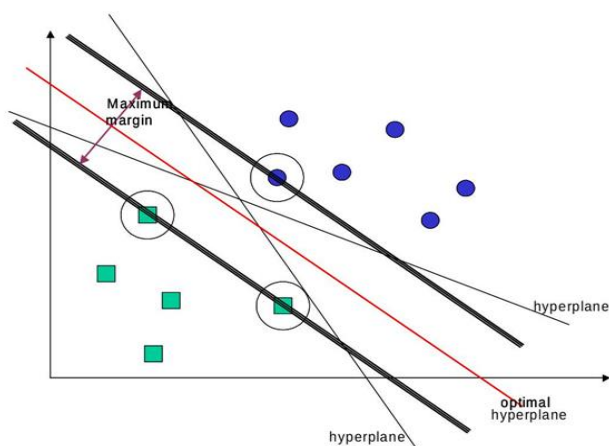


Fig2: Example of a linear separated SVM

The support vectors are usually a small subset of the training vectors and they are determining factors of how great support vector machines are great because it means the speed does not suggestively slow down with a higher number of features (Witten, Frank, Hall, 2017). Moreover, data can be imagined especially where a separation is not simple and this can be overcome by using lenient margins which allow for mis categorizations but in more complex scenarios the data can be mapped to a larger dimensional space to make room for other options of apportioning the data.

When we talk about distorted feature space, it is the one in the same with this larger dimensional space

but the simple linear separations in this altered feature space change into non-linear separation when you breakdown back into the original space (Zhang, IGI, 2012). Support vector machines are appropriate for web threat detection especially the application layer and the other lower layers in the OSI model which deal with network allied attacks such as denial of service, network threats among others. Due to its efficiency, it has the capability detecting web attacks to an overall accuracy of ninety nine percent. And hence making it better compared to the genetic algorithm.

IV. METHODS & PROCEDURES

4.1 System Overview:

Despite the fact that two rather different algorithms will be used the system is designed to operate more or less the same with the genetic algorithm or support vector machine components as loosely coupled modules to avoid having to redesign the system for both approaches. Web requests will be processed through a parser that looks for various aspects related to each of the three possible web attacks. The results are then output and can be used as input for either the genetic algorithm, support vector machine, or potentially another algorithm that could extend the testing. Finally, the testing modules will output the results to a file that can be processed by graphing tools, in this case the R programming language will be used to create graphs that can be used to draw conclusions on the two approaches.

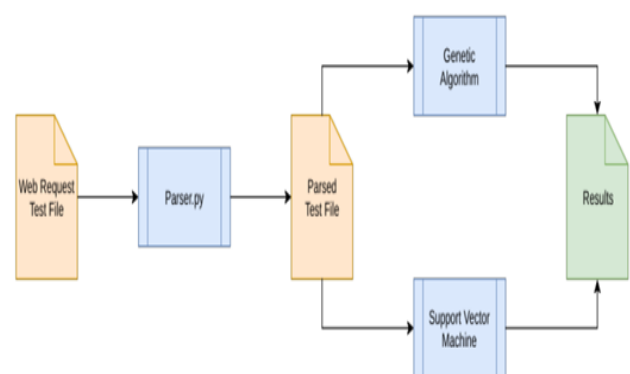


Fig3: Overview of the system

4.2 Gathering Test Data:

All data will be gathered from as close to a real-world scenario as possible. In order to do so, automated enumeration exploiting tools will be used to gather a large sample size of varied attacks. In order to gather SQL injection attacks the popular tool sqlmap will be used, for XSS attacks Grabber and XSSer will be used.

Request Type	Generation Method
SQL Injection	sqlmap [27]
XSS Attack	grabber [28], xsser [29]
RFI Attack	randomized generation
Normal Requests	httpfox [30]

Table1: Breakdown of test data generation

4.3 Parsing the Requests:

With the completed test file(s) containing the proper amount of each attack and normal

requests, the next step is to parse each request into numeric values so that the algorithms

can work on them.

4.4 Genetic Algorithm Based Signature Detection:

The method of the using a genetic algorithm for signature-based detection is largely the same as the proposed and tested system in previous research with a few modifications. One major difference is that instead of allowing the signatures to change to different attack type signatures (ex. SQLi to XSS) we specify what type of attack we want to search for and the algorithm uses that parsed result for every request.

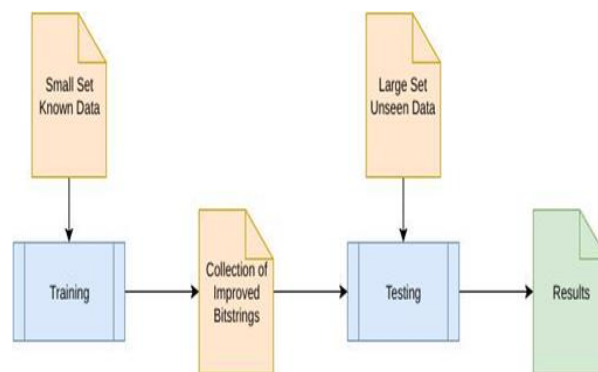


Fig4: Overview of the genetic algorithm system

4.5 Support Vector Machine Detection:

The SVM detection will follow a similar process to the genetic algorithm but instead of changing the parameters of the algorithm, the training data will be changed instead. This is because the parameters for the SVM that will make a difference are automatically optimized by a grid search approach provided by the same library used to implement the SVM in Python.

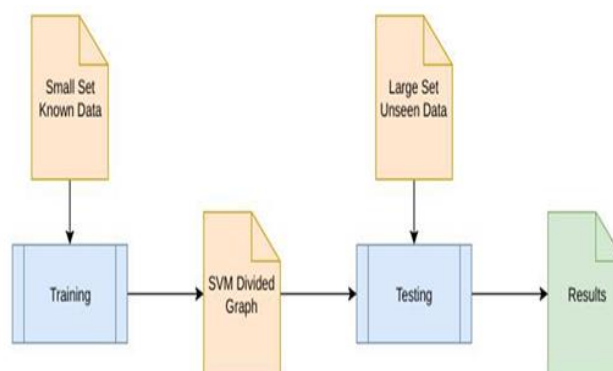


Fig5: Overview of the support vector machine system

V. SIMULATION RESULTS

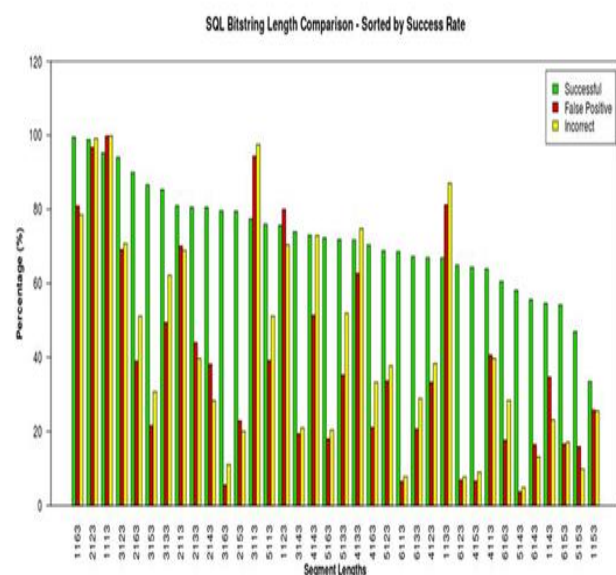
For curtness the focal point of the correlation is on the SQL infusion results over all experiments as it is the most mind-boggling demand out of the three, the other demand types will have their graphical outcomes shown nearby also however the former dialog will in general spotlight on the SQL infusions generally vigorously. A full posting of the content-based outcomes for each separate diagram is incorporated

5.1 Genetic Algorithm:

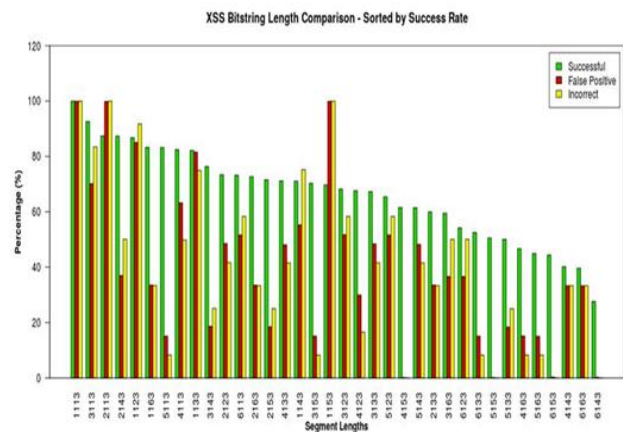
In the accompanying outcomes each test was rehashed multiple times and arrived at the midpoint of to give a superior thought of the run of the mill execution of the hereditary calculation approach portrayed before (Algorithm 4). The parameter that was changed sorts each test, and also the reason and expected consequences of the progressions to the specific parameters is incorporated. Likewise, except for the bit string section length test, every single other test are led utilizing the typical fragment lengths.

5.2 Bit string Segment Length Effects:

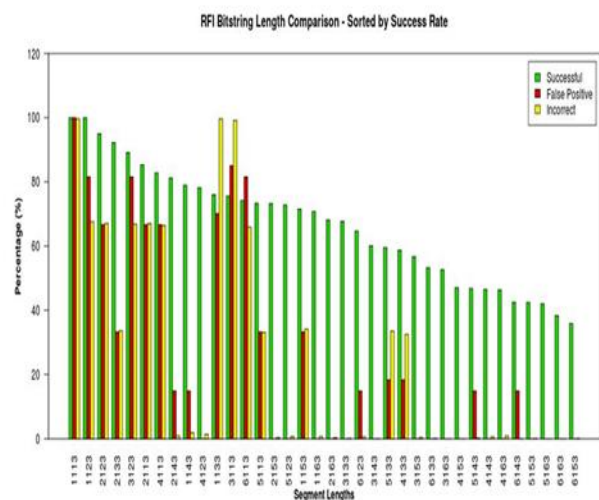
Since the hereditary calculation can identify extra assaults by producing new marks, if the quantity of conceivable mark blends is less because of the fragment lengths, at that point it would probably create these marks that coordinate with the preparation or testing assaults. Anyway it additionally opens up the likelihood of making it less demanding to produce poor marks because of a littler hunt space, and additionally possibly misleadingly increment results because of section over owes



Result1: Effects of Different Segment Lengths on Detecting SQLi



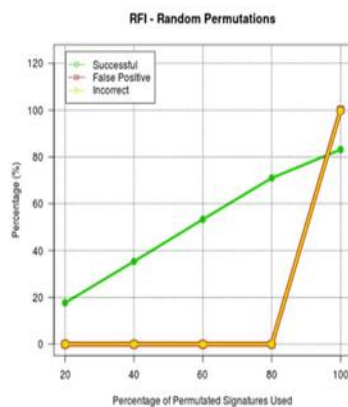
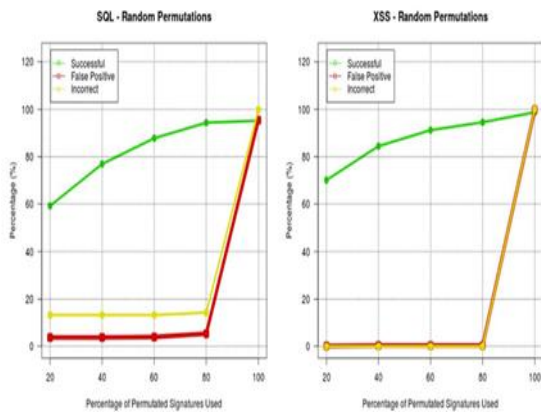
Results2: Effects of Different Segment Lengths on Detecting XSS



Results3: Effects of Different Segment Lengths on Detecting RFI

5.3 Random Permutations with Fitness:

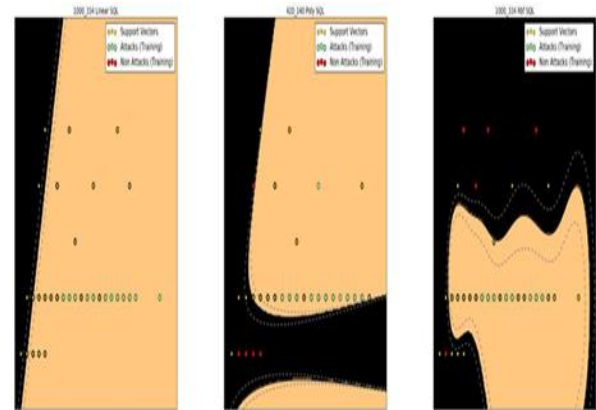
The hereditary calculation approach works thanks to the capability to haphazardly produce new marks utilizing well-being to get rid of the terrible marks, thus it might be exceptionally fascinating to distinction the methodology and primarily manufacturing each conceivable mix and simply utilizing the bit strings that performed well utilizing an identical well-being calculation used within the hereditary calculation. this might conceivably keep one's eyes off from the complexness and calculation time that joins a hereditary calculation



Results4: Permutation of Bitstrings for Detection

5.4 Support Vector Machine:

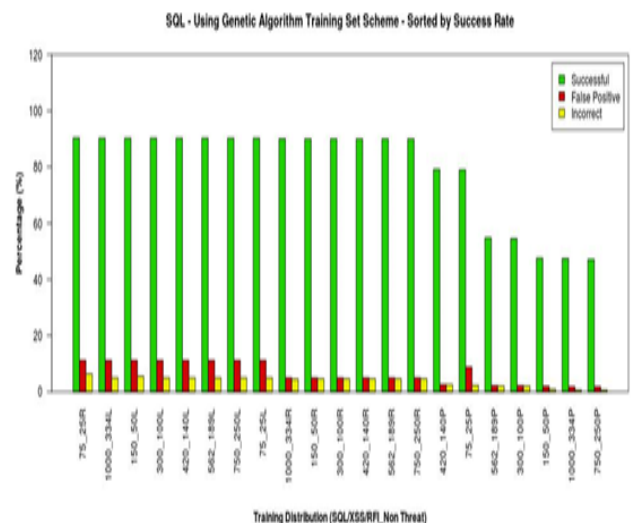
For the testing of the support vector machine it absolutely was not necessary to average along multiple results as there aren't any random parts within the formula so a similar results area unit created whenever. All tests used a similar testing knowledge to verify the coaching method yet because the same coaching knowledge whenever the desired range of requests didn't exceed the number of utilized in the genetic algorithms coaching set. once doing the genetic formula testing it absolutely was doable to live changes by adjusting the algorithm's parameters however within the case of the SVM this can be not the case so instead manipulating the coaching knowledge can provide observations on any changes in performance.



Results5: Classifier Output, Linear: 1000 334, Poly: 420 140, RBF: 1000 334

5.5 Comparison with Genetic Algorithm

Obtaining the remainder results uses a similar proportions of attack varieties for coaching as within the genetic algorithmic program coaching set of half hours for every attack kind and 100% remaining is non-threats. the aim of this check is to form as honest a comparison as doable with the genetic algorithmic program. The output classifier for the simplest performing arts instance of every kernel is additionally enclosed.



Results6: Genetic Algorithm and SVM comparison for SQLi Detection

VI. CONCLUSION

In conclusion, web applications have increasingly become one of the customary platforms for service releases and representing data and information over the worldwide web and they are the major cause of web threats and attacks. The most common application layer web threats include; remote file inclusion attacks, SQL injections, and cross-site scripting. The major ways in place today of fighting web threats and attacks is through the use of genetic algorithms and support vector analysis. Support vector machine uses the technique of categorizing data is dividing the data sets into two or more sets with prevalent margin imaginable between the separations which are known as hyperplanes. On the other hand, genetic algorithm is a search-based algorithm that utilizes the machine learning techniques in order to trace the prime or near-optimal solutions for a certain problem. However, from the analysis, support vector machines is more efficient compared to the genetic algorithm and thus it is more preferred in detecting and preventing web attacks.

REFERENCES

- [1] Gopichand, G., & Saravanaguru, R. A. K. (2016). A Generic Review on Effective Intrusion Detection in Ad hoc Networks. *International Journal of Electrical & Computer Engineering* (2088-8708), 6(4).
- [2] G. Gopichand, R.A.K. Saravanaguru, K. Ramesh Babu, Fully secured intrusion detection system for sensing attacks in MANET, *Journal of Advanced Research in Dynamical and Control Systems*, vol. 10, no. 4 Special Issue, pp. 810-816, 2018
- [3] Gopichand G, Saravanaguru RA.K., Collaborative Packet Dropping Intrusion Detection in MANETs, *Recent Patents on Computer Science* (2019) 12: 1. <https://doi.org/10.2174/2213275912666190618163426>
- [4] Gopichand G., Sankeerth K.S., Parlapalli A, Evaluation of recommendation systems using trust aware metrics, *International Journal of Recent Technology and Engineering*, Volume-7, Issue-6S4, April 2019
- [5] Gopichand G, Vishal Lella, SaiManikantaAvula, Enhancing Performance of Map Reduce Workflow through H2HADOOP: CJBT, *International Journal of Recent Technology and Engineering*, Volume-7, Issue-6S4, April 2019
- [6] Gopichand G, Sailaja G, N. VenkataVinod Kumar, T. Samatha, Digital Signature Verification Using Artificial Neural Networks, *International Journal of Recent Technology and Engineering*, Volume-7 Issue-5S2, January 2019
- [7] Gopichand G, Ra.K.Saravanaguru, .K.RameshBabu, Usage of AODV and AOMDV Protocols in Perceiving Black hole Attacks in a MANET, *International Journal of Pharmacy & Technology*, Volume 8, Issue 4, December 2016
- [8] Mehta M., Rajesh Mamilla, Sunithavenugopal, Gopichand G, Growth and development of start-ups in India - A study with respect to mechanical and production engineering, *International Journal of Mechanical and Production Engineering Research and Development*, Volume : 8-2, April 2019
- [9] Jitesh Shaw, P. M. Durai Raj Vincent, SenthilnathanPalaniappan, *, Arun Kumar Sangaiah, Gopichand G, Intelligent Phishing Detection System Using Feature Analysis, *Journal of Computational and Theoretical Nanoscience* Vol. 15, 2533–2538, 2018
- [10] SenthilnathanPalaniappan, SaiprasadPalli, Gopichand G, SirajudeenAmeerjohn, Siva ShanmugamGopal, Enhanced Handwritten Number Detection Using Kernel Discriminant Analysis (KDA), *Journal of Computational*

- and Theoretical Nanoscience Vol. 15, 2539–2543, 2018
- [11] H R Swathi, Shah Sohini, Surbhi, Gopichand G, Image compression using singular value decomposition, IOP Conference Series: Materials Science and Engineering 263(4).
- [12] Santhi H, Gopichand G, Gayathri P, Automated Smart Parking System using IoT, Journal of Advanced Research in Dynamical & Control Systems, Vol. 10, 09-Special Issue, 2018
- [13] P Gayathri, MayankAgarwal, H Santhi, Gopichand G, Bone Breakage Identification Using Image Processing Techniques, Journal of Advanced Research in Dynamical & Control Systems, Vol. 10, 09-Special Issue, 2018
- [14] Krishna Ganeriwal, G. P., Gopichand, G., &Santhi, H. Data Mining in Social Networks and its Application in Counterterrorism.
- [15] Priyadarsini, M. J. P., Rajini, G. K., Naseera, S., Balaji, S., Reddy, P. S. K., &Gopichand, G. (2006). AUTOMATIC OBJECT RECOGNITION BASED ON EUCLIDEAN DISTANCE RESTRICTED AUTO ENCODER.