

Application of Negative Selection Theory in Intrusion Detection

Yunhui Xing^{1, 2,a}, Zhaowen Lin^{1,2,3,b}, and Yan Ma^{1,c}

¹Network and Information Center, Institute of Network Technology, Beijing University of Posts and Communications, Beijing, 100876, China;

²Science and Technology on Information Transmission and Dissemination in Communication Networks Laboratory;

³National Engineering Laboratory for Mobile Network Security(No. [2013] 2685).

^awww.xingyh@163.com, ^blinzw@bupt.edu.cn, ^cmayan@bupt.edu.cn

Keywords: network security, artificial immune systems (AIS), negative selection algorithm(NSA), intrusion detection.

Abstract. Along with the rapid development of Internet, network security has become one of high-profile issues in the current life. In this paper, the problem of intrusion detection based on artificial immune system model and the key technologies negative selection algorithm(NSA) are discussed in detail, by the reviews of the recent history and present state. Simulation experiments reveals their characteristics of three matching algorithms including r-contiguous-bits, r-adjustable-bits and hamming distance algorithm.

Introduction

Now as the scale of computer network expanding and more and more types of network attacks appeared, information security is becoming more and more prominent. Intrusion detection comes to be a key technology for network security protection system, and turns into one of the hot issues of the computer network research in recent years. Intrusion detection is a kind of network security technique which detects acts that will sabotage security, availability and integrity of system. The most important research branch of intrusion detection is based on artificial immune theory [1].

Nowadays computer network are faced so many threats, which are caused by variety of factors. For example, information loss or disclosure intentionally or unintentionally during the network transmission; critical information (such as user, passwords and etc) stolen; using the Internet or computer resources without authorization; spreading virus (worms or Trojan) [2]. All of above are fatal factors leading to fragile network circumstance. Network security guarantees a normal service of the Internet, however, the traditional detection technology cannot adequately deal with such a situation any longer.

As one of the representatives of interdisciplinary, the artificial immune theory has been more and more important role to play in many fields. It is modeled off the T-cell maturing process[3] that happens in the thymus. T-cells of enormous diversity are first assembled with a pseudo-random genetic rearrangement process and those that recognize self-cells are eliminated before the rest are deployed into the immune system to recognize and attack outside pathogens. The methods based on this process are generally called negative selection algorithms (NSA) with the descriptive word "artificial" omitted[4] in order to improve the efficiency of network intrusion detection.

The Related Research at Home and Abroad

The Network Attack.

As the global informationization happened, a wide variety of network attacks appeared one after another. Network information are vulnerable to these events. In the middle of December, 2009, a malicious network attack known as the "Operation Aurora" launched quietly, in which Gmail service sustained attack. By sending emails to Google employees, the attacker obtained core data in Google company, leading to a huge loss for Google and its users. On the other hand, this event

made Advanced Persistent Threat(APT) a fashionable noun[5]. In 2013, Snowden tipped off the "PRISM", said that a large number of Internet users' personal information were being watched all over the world, after which the "network intelligence" characterized by big data security was raised to the level of national security strategy. In December 2014, Woo Yun platform Vulnerability Report issued a message that a large number of 12306 users' data were published on the Internet. In the attack, the hacker collected over 130,000 user data from 12306 by bumping database, owing to the lack of security mechanism. As a result, many train tickets booked were refund maliciously. The occurrence of these events exposed the huge flaw of network information security, meanwhile to network information security, it is also an unprecedented challenge.

Hacking is one of the most severe threats. Hackers seek for system flaws by various means, and raise an attack to destroy important data in system, causing threats to network security. Isolating inner network from Internet by firewall or others is a way to protect inner network against hacking. Whereas attackers often bypass the firewall by either faking data packets and modifying IP address, or launch a DDoS attack[6]. Under the flow of TCP packets, the server is not able to provide service normally. Against such attacks, TCP flow changes can be analyzed to determine whether attacks happened ,because fairness of network led to fewer and fewer resources to be obtained by TCP data flow that comply with congestion control, but this method has a high rate of false positive.

As a consequence of huge network access point with numerous public information, Illegal invaders destroy the validity and integrity of information in purpose and seizure resources. Also they modify computer information in an unauthorized manner that interfere normal operation or consume excessive system resources, leading to the result that authorized users cannot get access and other operation is delayed. In order to prevent "sniffer" or eavesdropping through unauthorized network, people realized that a useful measure of confidential transmission for sensitive data is needed. For instance, encryption in data link layer-setting of Wi-Fi security access. The link is called safe only in the case that the terminal is considered authorized to perform a given transaction legal entity. Enhancing the trust for each network or taking the authentication mechanism will ensure that a safe, reliable and private communication is provided to these entities on the Internet. National governments and their agencies develop a series of security protocols to ensure user authentication and authorization security, like SSL/TLS, a well-known technique of end to end encryption. However, it has not been widely adopted attributed to historical reasons and organizational inertia in part, ignorance or wrong information in the other part.

No matter how powerful the software are, flaws in the design of the beginning is inevitable, mainly due to the unsafe source code and inappropriate use of the program by users. Different software platforms, hardware and system settings will cause different security vulnerabilities. Rely solely on the technical of code analysis is not appropriate, and now most of software attempt to repair vulnerabilities by patching, which often means a lag of time, leads to serious internal insecurity.

Preventive Measure and Technology.

To ensure the security of network information, a variety of strategies will be taken in practical applications, such as virus protection, firewall, and data encryption.

Intrusion detection technology is the key technology of information security discussed in this article[8]. Intrusion detection system is to collect information from a variety of computer system and network system, as well as the analysis of information on the computer and cyber source.

In general, intrusion detection techniques can be categorized into two groups: misuse detection and anomaly detection. The misuse detection systems use patterns of well-known attacks or weak spots of the systems to identify intrusion. The weakness of misuse detection systems is that it is unable to detect any future(unknown) intrusion until corresponding attack signatures are intruded into the signature database. Anomaly detection methods try to determine whether the deviation is from the established normal usage patterns or not. The critical success of anomaly detection relies on the model of normal behaviors.

Combined with the feature of "no prior knowledge of nonself is required", a intrusion detection is mentioned in this paper. It imported the biological immunology of NSA to ensure more accurate

detection rate. There are several significant differences between the algorithm described here and more conventional approaches to change detection. First, the checking activity can be distributed over many sites with each site having a unique signature; second, the quality of the check can be traded off against the cost of performing check; third, protection is symmetric in the sense that the change detector and protected data set are mutually protective; fourth, the algorithm for generating the change detectors is computationally expensive, although checking is cheap, so it would be difficult to modify a protected file and then alter the detectors in such a way that the modification could not be detected. As with other methods, this method relies on the guarantee that the data to be protected are uncorrupted at the time that the detectors are generated.

Negative Selection Algorithm

NSA is one of the major algorithms developed within AIS and can be used for network security, fault detection, especially, anomaly detection. The NSA developed by Forrest is proposed in the use of a binary calculation model for pattern recognition method [9].

Traditional Negative Selection Algorithm.

The main task of NSA can be viewed generally as the problem of learning to distinguish self(legitimate users, corrupted data, etc.) and nonself(unauthorized users, viruses, etc.). We define self as the string to be protected and nonself as any other string. The string could be a string of bots(and hence, anything that can be represented in a digital computer), a string of assembler instructions, a string of data, etc. However, the method appears to be most relevant for strings that do not change over time, that is the protected strings need to be fairly stable.

A partial match between two string of equal length means that the symbols are identical when beyond a definite matching threshold r . Thus, for any two strings x and y , we say that $\text{match}(x, y)$ is true if x and y agree at least r contiguous locations. The matching rule can be applied to strings defined over any alphabet of symbols. It also illustrates an approximate idea by taking partial rules of NSA. Three common matching rules are described in this paper.

The algorithm has three phases:

Step1 Define collection S of self data: each String that is part of protected data defined over the alphabet $\{0,1\}$ in collection S .

Step2 Generate a collection R of detectors: The ones in random collection that match any of strings in collection S are eliminated. Each detector R is a string that does not match any of the protected data in collection S .(shown in figure 1).

Step3 Monitor the protected data by comparing them with the detector R . If detector is ever activated, a change is known to have occurred. Each detected nonself string matches string of detector R (shown in figure 2).

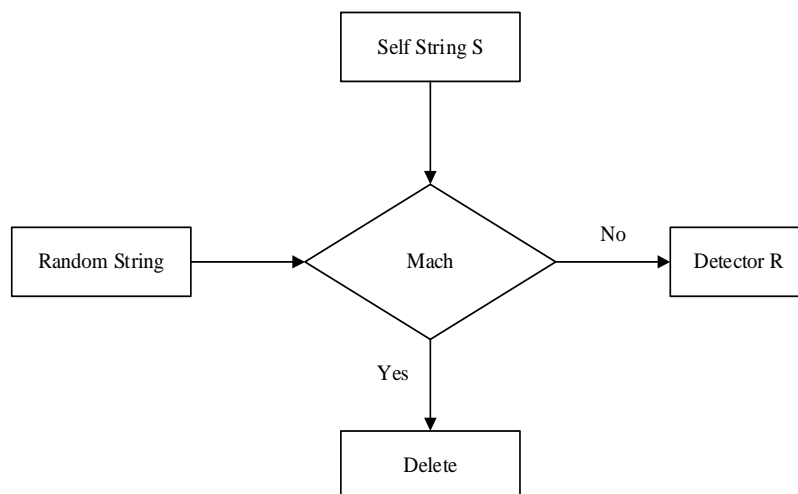


Fig. 1 The way to generate valid detectors

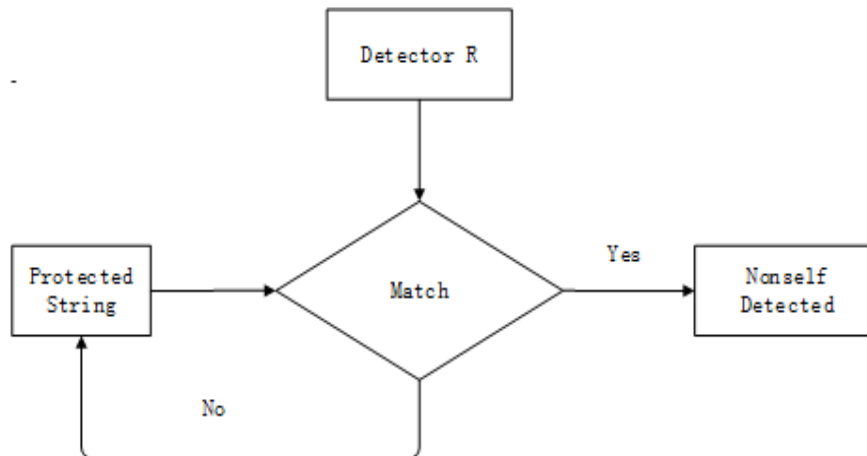


Fig. 2 Monitor protected string for changes

Figure 1 and 2 illustrate how the algorithm works. Each copy of the detection system generates its own unique valid set of detectors once, and then runs the monitoring program regularly to check for changes.

Matching Rules.

R-contiguous-bits matching rule(rcb) is that two strings of equal length are said to match if they are identical in at least r contiguous positions. For example, $x=100111010$, $y=101111001$, x and y defined as binary alphabet $\{0,1\}$ matching at four contiguous locations. Thus, $\text{match}(x, y)$ is false for $r=5$ or greater. Since x and y agree at 4 contiguous locations, $\text{match}(x, y)$ is true for $r=4$ or less. If more nonself strings can be detected, worse specificity it is, conversely, too. It can be seen that the larger r is set, the higher precision and the lower false alarm rate is get and vice versa[10].

R-adjustable-bits matching rule(rab) increased the detection rate by adjusting the matching threshold on the basis of rcb proposed by Zhang Heng[11].The core idea of rab by adjusting the threshold value of this relatively simple approach significantly reduces the number of hole.

Threshold r is set in order to determine whether two strings match in hamming distance matching rule. When the distance between two strings of x and y greater than or equal to the value of threshold r , it can say $\text{match}(x,y)$. It is different from rcb rule that hamming distance rule dose not need to be contiguous.

Hole.

In this paper, partial matching rule specifies the approximation, and the parameter r controls how large the approximation is. As shown in figure 3, it turns out that there may be some strings not in self space, called "holes" for which it is impossible to generate valid detector and cannot be avoid in NSA[12]. For clarity, I have reproduced the example presented there: If S contains two strings s_1 and s_2 that match each other, they may induce two other strings h_1 and h_2 that cannot be detected because any candidate detector would also match either s_1 or s_2 , two strings h_1 and h_2 are defined as "hole".

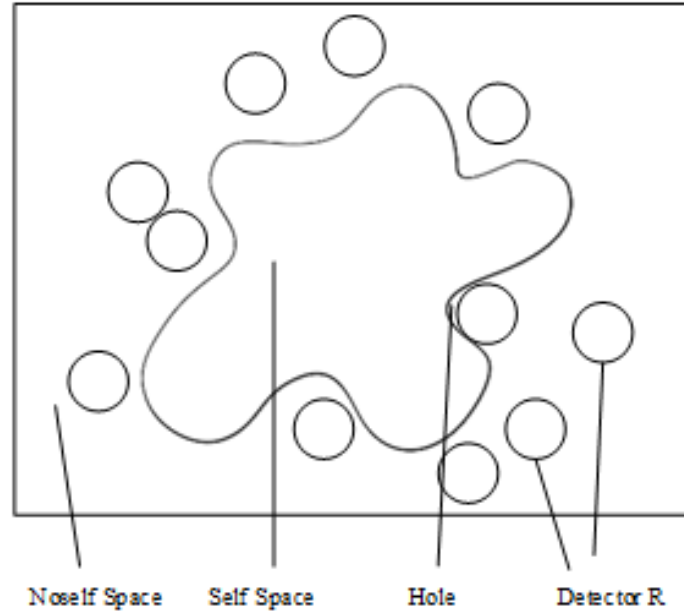


Fig. 3 The schematic diagram of hole

Experiment

Based on the analysis above, simulation experiment has been cast pointing to the three algorithms, including rab, rcb and hamming distance. The experiment uses iris data set as data collection, and analyzes detection distribution, probability of random match between two strings, detection rate and detector quantity in different maximum matching threshold. The following are parameters and definitions:

L length of strings;

N_{R0} the number of initial detector string(before censoring);

N_R the number of detector strings after censoring(size of repertoire);

N_S the number of self strings;

P_M the probability of a match between 2 random strings by rcb:

$$P_M = 2^{-r} \left(\frac{L-r}{2} + 1 \right) \quad (1);$$

P_H the probability of a match between 2 random strings by hamming distance

$$P_H = 2^{-r} \sum_{i=r}^L C_i^r \quad (2);$$

P_S the probability that NR detector detect an intrusion;

f_i a random string does not match any self strings when threshold is r_i :

$$f_i = (1 - P_{Mi})^{N_S} \quad (3);$$

P_i the probability to be a detector when threshold is r_i :

$$P_i = (1 - f_1)(1 - f_2) \dots (1 - f_{i-1}) f_i \quad (4);$$

P_{fM} the probability that NR detector fail to detect an intrusion by rab:

$$P_{fM} = (1 - P_M)^{N_R} \quad (5);$$

P_{fH} the probability that NR detector fail to detect an intrusion by hamming distance:

$$P_{fH} = (1 - P_H)^{N_R} \quad (6);$$

Detection parameters can be specified before the experiment. According to Eq. (1), N_R is independent of N_S when P_f and P_M are fixed. This implies that the size of self data collection has

nothing to do with the increase of detector quantity, and therefore a relatively small size of data collection is required in the experiment.

For rab matching rule, matching threshold gradually increases from $r_1, r_2 \dots$ to r_c . Different detectors produce different detection range, while by rcb, detection range is fixed using different detectors. From Eq. (1), It can be learned that P_M is 2^{-L} while matching threshold grows to L , thus rcb algorithm can be viewed as a special case in rab algorithm.

From Eq. (3) and Eq. (4), these allow us to predict mature detection distribution under different threshold satisfied:

$$N_{Ri} = N_{R0} P_i \tag{7}$$

Combining Eq. (7) and experiment, the distribution of mature detectors is shown in Table 1 when $N_{R0}=2000, N_S=100$.

Table 1 Number of mature detectors under different matching thresholds

Rc	13	14	15	16	17	18	19	20	total
detectors	1312	292	208	99	52	16	11	8	1998
Theoretical 100%	65.57%	27.97%	5.83%	0.6%	0.3%	0%	0%	0%	
experiment 100%	65.65%	14.60%	10.5%	4.95%	2.6%	0.8%	0.55%	0.40%	

It can be concluded that with the increase of the maximum matching threshold, detectors tend to be mature basically, and as a result, the probability of mature of immature detector tends to 1.

Through Eq. (3) and Eq. (4), compare the following three algorithms: hamming distance, rcb and rab matching rules in the same requirements, simulation figure is shown below in Figure 4.

It is known that the matching probability of hamming distance is better than that of rab and rcb, and among which rcb is worse.

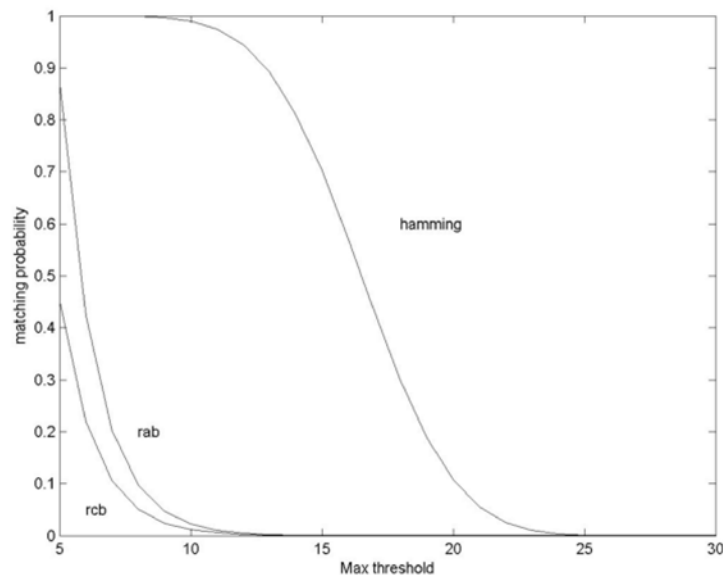


Fig. 4 Probability of matching under rcb, rab and hamming distance

For rab and rcb matching rules, the simulation Figure 5 demonstrates that the detection tends to be flat when the maximum threshold is more than 15. The detection rate is the lowest when $r_c=13$, it happens to be rcb matching rule when $r=13$ in this case. This shows that the detection rate in rab matching rule is higher than it in rcb rule.

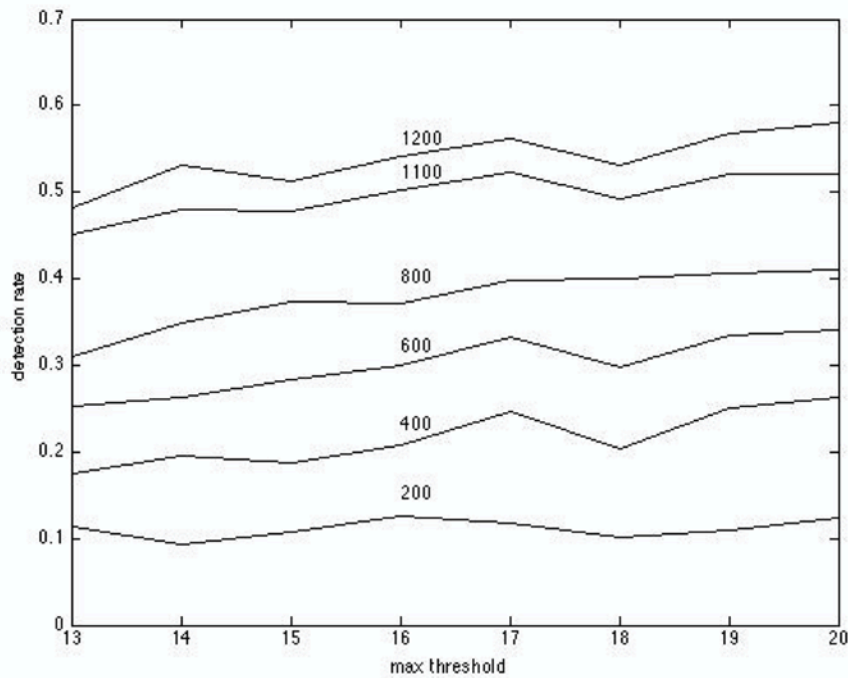


Fig. 5 Detection rate under different max threshold

The core idea is to expand the coverage of each maturity detection by changing the threshold, ultimately reducing the number of hole arising from partial matching. According to Eq. (5) and Eq. (6), It is shown in Table 2 to compare their failure rate, in order to reach the same detection failure rate in the case of the same threshold value, the number of detectors by hamming distance matching rule is far less than rcb matching rule, is more suitable for case of frequently real-time detection.

Table 2 The failure probability between rab and hamming distance

R	13	14	15	16	17	18	19	20
NR	1984	4191	8876	18862	40240	86232	185738	
402470	rcb							
Pf	0.1159	0.1192	0.1219	0.1210	0.1105	0.1190	0.0944	
1.1027								
NR	2	2	2	2	5	7	12	21
hamming								
Pf	0.0582	0.0788	0.0942	0.1154	0.1023	0.1356	0.0944	
0.1018								

Summary

This paper focuses on revisiting various matching rules in NSA to detect nonself. According to simulation experiment, NSA using three different matching rules has been worked out, and the impact of different matching rules for detection has also been shown with the experimental simulation data graphs. The detection rate of Rab algorithm based on rcb matching rule has further improved, and the matching probability of hamming distance rule is much higher than the other two algorithms, more suitable for application in the high requirement of real-time. Currently the theoretical foundation of this method is still developing in the sense that most applications have their own versions of the algorithms, including variability in the data representation, detector representation, detector generation mechanism, etc. A widely accepted standard model or taxonomy is still to be established for this method to mature. More accurate detection method may be found if these cases are further studied and the research on instruction will be strived for further improvement.

Acknowledgments

This work is supported by the National High Technology Research and Development Program of China (863 Program) (Grant No. 2013AA014702), the Fundamental Research Funds for the Central Universities (Grant Nos. 2014PTB-00-04, 2014ZD03-03). In addition, the authors would like to thank the professor Lin and the students in Information Network Center of BUPT for their valuable contribution to recommendations of this paper and the implementation of relevant projects.

References

- [1] Hu R H, Lou P H, Zhao P. A novel approach of detector generation for real-valued negative selection algorithm. *Applied Mechanics & Materials*, 2011, Vols.121-126.
- [2] Zhong Q H, Gou Y C, Jiang Z J, et al. Research and quantitative analysis to factors of network and information security. *Computer Technology and Development*, 2014, 24(2):172-175.
- [3] Zeng J Q, Qin Z G, Tang W W. Anomaly detection using a novel negative selection algorithm. *Journal of Computational and Theoretical Nanoscience (Impact Factor: 1.03)*, 2013, 10(12): 2831-2835.
- [4] Gong M G, Zhang J, Ma J J, et al. An efficient negative selection algorithm with further training for anomaly detection. *Knowledge-Based Systems*, 2012 Vol 30:185-191.
- [5] Virvilis N, Gritzalis D, Apostolopoulos T. Trusted Computing vs. Advanced Persistent Threats: Can a defender win this game? *IEEE 10th International Conference on Ubiquitous Intelligence & Computing and IEEE 10th International Conference*, 2013, 396-403.
- [6] Chen Z, Wen W, Yu D. Detecting SIP flooding attacks on IP Multimedia Subsystem (IMS). *International Conference on Computing, Networking and Communications*, 2012, 154-158.
- [7] He X J, Chomsiri T, Nanda P, et al. Improving cloud network security using the Tree-Rule firewall. *Future Generation Computer Systems*, 2014, 30: 116-126.
- [8] Elfeshawy N A, Faragallah O S. Divided two-part adaptive intrusion detection system. *Wireless Networks*. 2013,19(3):301-321.
- [9] Forrest S, Perelson A S, Allen L, et al. Self-nonsel self discrimination in a computer. *IEEE Symposium on Research in Security and Privacy*, 1994, 202-212.
- [10] Lu TL, Zheng K F, Fu R R, et al. Anomaly detection system with hole coverage optimization based on negative selection algorithm. *Journal on Communications*, 2013, 34(1): 128-135.
- [11] Zhang H, Wu L F, Zhang Y S, et al. An algorithm of r-adjustable negative selection algorithm and its simulation analysis. *Chinese Journal of Computers*, 2005, 28(10): 1614-1619.
- [12] Liu X B, Cai Z X. Properties assessments of holes in anomaly detection systems. *Journal of Central South University (Science and Technology)*, 2009,40(4):986-992.