



Journal of Applied Sciences

ISSN 1812-5654

science
alert

ANSI*net*
an open access publisher
<http://ansinet.com>

An Efficient Analysis of Honeypot Data Based on Markov Chain

¹K.P. Kimou, ²B. Barry, ¹M. Babri, ¹S. Oumtanaga and ¹T.L. Kadjo

¹Laboratory for Informatics and Telecommunications Research,
INPHB, 08 BP 475 Abidjan 08, Cote d'Ivoire

²University Cheikh Anta Diop, BP 5353 Dakar-Fann-Senegal, Dakar, Senegal

Abstract: The aim of this study is to provide a behavioral modeling of hackers for a better predisposition to secure the Internet. Thus, using data collected from the experimental platform of CADHo project, an efficient analysis of these data has been made. This preliminary study has identified the intruders most used ports (80, 135, 139 and 445). Using an approach based on Markov chains, we propose a predictive model characterizing the attack processes on a honeypot and emphasizing the most attacked ports during a given period. By computing the estimated errors, the efficiency of the proposed model is proved. This model is independent of any platform and can be generalized for any number of predetermined ports.

Key words: Cyber criminal, honeypot, honeynet, hacker, intruder

INTRODUCTION

The current protection system against intruders always presents inadequacies. Indeed, the number of incidents and vulnerabilities is on the increase and that has an effect on the productivity of the information processing systems. To improve the protection system a better knowledge on the threats and their intruders are required. To this end, we need methods and data analysis in order to characterize the attacks and evaluate their impacts. Today, the best technology which makes it possible to meet these needs is that based on honeypots.

Once collected, the data on honeypots must be analyzed in order to forewarn or device design for safety tools. Several statistical models have already been used to analyze the data.

On the one hand, for the modeling of the duration between the occurrence of two consecutive attacks, one can quote the Pareto and Exponential distribution or the combined Pareto and Lognormal distribution (Alata *et al.*, 2006; Kaaniche *et al.*, 2006) and on the other hand for the modeling of the attack frequency, the normal distribution has already been used (Oumtanaga *et al.*, 2006).

The purpose of this article is to model the data extracted from the centralized database of the Leurre.com project (Alata *et al.*, 2005). These data are related to the number of observed attacks against the following ports 80, 135, 139 and 445. For that purpose, Markov chains are efficiently used. This work lies within the scope of the CADHo (Collection and Analysis of Data from Honeypots) project (Alata *et al.*, 2005).

A honeypot is an information system whose value lies in its compromising (Alata *et al.*, 2005; Oudot and Glaume, 2006; Kaaniche *et al.*, 2006; Spitzner, 2002). The main goal is to study the attacks against the information processing system and to learn more about their perpetrators with a view to prevention or for a better detection of these incidents. When honeypot platform is deployed, finality is to find optimal countermeasure strategies to protect networks by drawing the attention of the pirates towards the honeypot. For this purpose, this honeypot is called production honeypot, traditionally used by commercial organizations. Also, honeypot platform can be used to observe the pirates' behaviors in order to get information about their attack tools and strategies; this type of platform is called research honeypot. Present study focuses on this latter type of honeypot.

In connection with the research honeypots, there are two types according to the level of interactivity with the pirates. We distinguish honeypots with low interaction and honeypots with strong interaction. The first ones make it possible to collect a maximum of information (like the services concerned, the type of attack, the attacker's operating system, etc.) while offering a minimum of privileges to the attackers. The advantage of honeypots with low interaction is that the attacker cannot use the compromised host as source of further attacks, which is not the case of honeypots with high interaction. These last ones are more generous. A host of the network which offers functional services such as authentication or scripts execution is used.

By implementing a honeypot network (with a firewall and an IDS (Intrusion Detection System)) to simulate a real environment of system and network resources, we obtain a honeynet. We distinguish many kinds of honeynets:

- Physical honeynets which are genuine machines intended to be compromise
- Virtual honeynets which simulate a honeynet on a single machine
- Distributed honeypot systems where the data collected are stored in a central place. The Leurre.com project is one example

The Leurre.com project (Alata *et al.*, 2005) uses several honeypots distributed throughout several countries through the world. Thus, this project aims at developing and giving to the scientific community, a distributed platform of data collection. Our data are linked to that project.

Also, the implementation of a honeypot requires the use of machine emulator or virtual systems as well as a monitoring system. To this end, there is a set of honeypot tools such as honeyed, VMware and Sebek (Lockhart, 2004).

MATERIALS AND METHODS

This study was developed through research projects on honeypots within the LABTIC (Laboratory of Information and Communication Technologies). LABTIC became partner in the project CADHo since 2006. LABTIC has led researches on the analysis of data from honeypots. One of main purposes of this project is to establish a distributed collection of data for attacks analysis (Leurre.com environment) by installing honeypot at each partner's platform. In addition, the project CADHo makes available to its partners, the data collected by the platform. Any team wishing to benefit from these data of the distributed database should become a project partner. The study started since 2007 and data are collected from this platform.

General information on Markov chain: Markov chains are a powerful tool for the modeling of statistical data. They denote a sequence of stochastic events where the probability that a certain future state will occur depends only on the present or immediately preceding state of a variable or system, but not on the events leading up the present state. These chains are efficiently used to describe random processes able to reach a number of final states and the evolution of systems through these states.

Definitions: A Markov chain consists of states and transitions probabilities. Let us consider $(X_t)_{t>0}$ a random process, where t denotes the observation time (t can be expressed in hours, days, week, etc.). It is assumed that during this time, the system can be in one of the following states: $1, 2, \dots, n$, where n is an integer. Let us set E , the finite state space. During the evolution of the system, this last one can be in one of these states with a certain given probability. The probability that the system stays in state i at time t will be $\text{prob}(X_t = i)$ noted by and, $P_{ij}(t)$ called transition probability, represents the probability that, at the given time t , the system being in state i , passes to state j at time $t+1$. These vectors P_{ij} define the transition matrix $P = P_{ij}(t)_{i,j \in E}$ an $n \times n$ matrix:

$$P_{ij} = \text{prob}(X_{t+1} = j / X_t = i) \quad (1)$$

Assumptions on this probability make it possible to obtain the conditions of application of the Markov chains.

Properties: Many properties define Markov chain.

Regularity conditions: A Markov chain with transition matrix P is said to be regular if P^n has all positive entries, from a certain value n . In a regular Markov chain, it is possible to get from any state, any other state in n steps

Independence and homogeneity assumptions: The independence and homogeneity assumptions are defined, respectively as follows:

- The probability $P_{ij}(t)$ that the process goes into the state j at the moment $t+1$ knowing that it is in state i at the moment t is known and independent from the system previous states
- The probability $P_{ij}(t)$ is time independent

The probabilistic state of the system at the moment t is a probability law of the possible states on set E . It is a stochastic line vector $\alpha_{(t)}$ defined by:

$$\alpha_{(t)} = (\text{prob}(X_t = 1) \text{ prob}(X_t = 2) \dots \text{prob}(X_t = n)) \quad (2)$$

$\alpha_{(t)}$ is also called law of marginal distribution and $\alpha_{(0)}$ initial law or initial state of the Markov chain.

Under the independence and homogeneity assumption, we have:

$$\alpha_{(t+1)} = \alpha_{(t)}P \text{ or } \alpha_{(t)} = \alpha_{(0)}P^t \quad (3)$$

The second equation of relation (3) is obtained from the first one by reiterating it to order 0. Thus, the law of

Markov chain is well determined by the data of its matrix of transition P and the initial state $\alpha_{(0)}$.

Convergence of Markov chain: The convergence of a Markov process makes it possible to anticipate the behavior of the system independently from its initial state. Thus, when the conditions of convergence are satisfied, it is shown that at a given time, the law of probability is independent from the initial law. To this end, there are theorems which make it possible to prove under certain conditions, the convergence of Markov chain (Bremaud, 2001; Nuel and Prum, 2007). Among them, we will state the principal theorem in the theory of Markov chain (Graham, 2008). We now state the principal theorem in the theory of the Markov chain: if P is a regular matrix of transition from a Markov chains process, if β is an unspecified vector of states, then:

- $\lim_{n \rightarrow \infty} P^n = S$, where S is a stochastic matrix having its lines all equal
- $\lim_{n \rightarrow \infty} \beta P^n = \alpha$, where α is a permanent stochastic vector (the sum of its coordinates is 1). α is called the vector of balance or stationary distribution

Related works based on Markov chain: Used by Claude Shannon to introduce the concept of entropy in the book (Shannon, 1948), Markov chains are used by others researchers to describe various situations. For examples, Markov chains have allowed the setting up of an efficient arithmetic coding used for data compression (Haixiao *et al.*, 2006). In bioinformatics, these chains have been useful for the modeling of certain nucleotides' properties (Nuel and Prum, 2007).

Since this current decade, Markov chains have become efficient tools used to describe intruders behaviors. Thus, Nong (2000) has used Markov chain model to represent a temporal profile of normal behavior in a computer and network system. This technique

developed by Nong allows detecting the anomalous behaviors of system caused by intruders. To construct the Markov chain, he has identified a set of 284 possible states which result from audit events on a single UNIX-based host. The intrusive activities are efficiently detected by Nong Ye using Markov chain model.

Callegari *et al.* (2008) compared performance of several technique of intrusion detection system based on various types of Markov models. There study improves the mechanism developed by Ju and Vardi (2001) and Jha *et al.* (2001) to describe the use of high order Markov chains to detect masqueraders at the host level.

All these works emphasize the efficient use of Markov chains or other Markov models in intruders' behaviors study.

Therefore, while these former studies take their input from single node or limited networks, our study is part of a large project and uses a great and various data. The honeypot environment is much more adapted for intruders' behaviors study than traditional networks or normal host.

Collected data, reorganization and Markovian modeling:

This section concerns data collected from leurre.com database. The ports with negligible frequency of attack are not part of the final modeling. Only the most attacked ports have been taken into account in the final modeling.

Collected data: The raw data relating to the weekly attack number of ports 80, 135, 139 and 445 and extracted from the centralized database of Leurre.com are presented in Table 1.

The service corresponding to each port is mentioned in the first column of the table. The data of the attacks are observed over 50 periods of a week about over one year.

Reorganization of the data: From previous data, we have the table below, giving for each week, the number of the

Table 1: Results of the number of attacks per port taken from the database of leurre.com

		Observation periods																								
Services	Ports	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11	P12	P13	P14	P15	P16	P17	P18	P19	P20	P21	P22	P23	P24	P25
HTTP	80	0	0	0	7	10	0	14	1	11	7	7	9	4	0	0	0	0	0	17	24	5	0	0	0	
EPMAP	135	0	1	3	77	90	1	170	42	126	108	171	287	97	37	37	0	4	0	3	141	109	51	1	0	0
NETBIOS-SSN	139	0	3	1	33	83	1	112	33	146	97	126	173	70	25	13	0	5	0	5	108	116	77	2	3	0
Microsoft-DS	445	0	0	0	2	1	1	0	0	1	1	0	1	0	0	0	0	3	0	3	200	170	85	3	4	1
		Observation periods																								
Services	Ports	P26	P27	P28	P29	P30	P31	P32	P33	P34	P35	P36	P37	P38	P39	P40	P41	P42	P43	P44	P45	P46	P47	P48	P49	P50
HTTP	80	3	43	107	0	4	61	29	26	20	139	130	6	67	106	67	0	0	21	0	0	0	1	0	0	0
EPMAP	135	49	307	148	0	38	249	203	253	333	380	314	14	195	227	189	0	0	106	0	0	0	0	0	1	0
NETBIOS-SSN	139	44	364	309	1	43	287	326	373	333	487	415	17	262	335	267	0	0	105	0	0	1	0	0	0	0
Microsoft-DS	445	74	402	286	4	43	322	387	351	402	415	375	18	232	263	239	0	0	126	0	0	1	0	0	0	0

Table 2: Chronological tables of the most attacked ports

Periods	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11	P12	P13	P14	P15	P16	P17	P18	P19	P20	P21	P22	P23	P24	P25
Most attacked ports	0	139	135	135	445	135	135	135	139	135	135	135	135	135	135	0	139	0	139	135	139	139	139	139	445
Periods	P26	P27	P28	P29	P30	P31	P32	P33	P34	P35	P36	P37	P38	P39	P40	P41	P42	P43	P44	P45	P46	P47	P48	P49	P50
Most attacked ports	445	445	139	445	445	445	139	445	139	139	445	139	139	139	0	0	445	0	0	139	80	0	135	0	

most attacked port and the frequency of attacks again the service using that port. Two situations can occur to the system during the period of observation:

- No port is attacked
- At least, two ports undergo the same number of attacks

When no port is attacked, the assumption will be made that the port more attacked is the null port noted 0. If two (or more) ports undergo the same number of attacks, we assume that the most-attacked port is the port which will have been attacked the least during the former attack periods (Table 2).

Markovian modeling: In this subsection, Markov chain is used to model the system evolution. First, the initial state is defined. After this, the transition matrix is constructed.

- **Definition of the states of system:** Let us indicate T, the period of observation of the corresponding to one week and by (Σ), the process which consists in observing the ports 80, 135, 139 and 445 each week in order to determine which is the most attacked.

Thus, our system consists in associating to time kT, the current most attacked port, where k>0 describes the set of integer.

For $i \in \{1, 2, 3, 4\}$ let us indicate by P_i the port associated with state i. The analysis of Table 1 shows that (Σ) can be in five states, defined as follow:

- For the state i, with $i \in \{1, 2, 3, 4\}$, P_i denotes the most attacked port. The 4-sequence ports (P_1, P_2, P_3, P_4) equals (80, 135, 139, 445)
- **State 5:** No port is attacked

That defines a random process $(X_t)_{t=kT, k \geq 0}$, which takes its values in E. Indeed, it is impossible when knowing the state of the system at the moment kT, to know the state at moment (k+1)T.

Let us set T = 1 week, we have $(\Sigma) = (X_t)_{k \geq 0}$.

Transition stamp of attack processes: Let us consider the events $X_n = i$ and $X_{n+1} = j$ which, respectively expresses that the system is in state i at the moment n and the

system is in the state j at the moment n+1. Then the event $(X_n = i) \cap (X_{n+1} = j)$ expresses that the system is at state after having been in state i. There is the following definition:

$$\text{prob}(X_{n+1} = j / X_n = i) = \frac{\text{card}((X_n = i) \cap (X_{n+1} = j))}{\text{card}(X_n = i)} \quad (4)$$

By assuming that the probability above does not depend on n, we can define and calculate the transition matrix P of attack process. The data of Table 1 and Eq. 1 and 4 enable us to calculate these elements of P.

Hence,

$$P = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 0 & \frac{8}{13} & \frac{3}{13} & \frac{1}{13} & \frac{1}{13} \\ \frac{1}{16} & \frac{3}{16} & \frac{3}{8} & \frac{1}{4} & \frac{1}{8} \\ 0 & \frac{1}{11} & \frac{4}{11} & \frac{5}{11} & \frac{1}{11} \\ 0 & \frac{1}{8} & \frac{1}{2} & \frac{1}{8} & \frac{1}{4} \end{pmatrix}$$

Thus one can consider the sequence $(X_k)_{k \geq 0}$ as Markov chain with transition matrix P. We can from now, apply the properties of Markov chain previously defined.

- **Convergence of the system:** By calculating P^3 we get

According to the principal criterion of regular process previously quoted, the studied process is a regular Markov chain. Thus, the probability law converges and is independent of the initial law, according to the convergence properties.

We consequently deduce that there exists a finite value such as:

$$\lim_{t \rightarrow \infty} \alpha_{(t)} = \alpha \quad (5)$$

Let us take for initial law that which corresponds at the initial state of the system, e.g., that which does not correspond to any attack of the implied ports. Thus, we get:

$$\alpha_{(0)} = (0 \ 0 \ 0 \ 0 \ 1)$$

RESULTS AND DISCUSSION

Results of the model: We know, according to the study led to section 3, that the probabilistic state of the system is described by the relations (2). Thus, at the end of the first period, we can predict:

$$\alpha_{(1)} = \alpha_{(0)}P$$

$$= (0 \ 0 \ 0 \ 0 \ 1) \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 0 & \frac{8}{13} & \frac{3}{13} & \frac{1}{13} & \frac{1}{13} \\ \frac{1}{16} & \frac{3}{16} & \frac{3}{8} & \frac{1}{4} & \frac{1}{8} \\ 0 & \frac{1}{11} & \frac{4}{11} & \frac{5}{11} & \frac{1}{11} \\ 0 & \frac{1}{8} & \frac{1}{2} & \frac{1}{8} & \frac{1}{4} \end{pmatrix}$$

Thus, after one week, the probability so that port 80 is the most attacked port is null, ports 135 and 445 have 12.5% of chance to be the most attacked ports, port 139 has 50% of chance to be the most attacked port and 25% of chance that no port is attacked. That can be predicted at the end of the second period.

After two weeks, there is 38.68% of chance for port 139 to be the most attacked port.

In a general way, the second relation of (3) enables us to predict the state of the process with nth period. By using the relation (5) and while passing in extreme cases in the first equation of (3) we get:

$$\alpha = \alpha P \Leftrightarrow \alpha(I-P) = 0 \tag{6}$$

In the above relation, I indicates the 5-by-5 identity matrix. Moreover α is a stochastic vector according to the general theorem, therefore, its coordinates α_i with $i \in \{1, 2, 3, 4, 5\}$, verify:

$$\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 + \alpha_5 = 1 \tag{7}$$

where, $\alpha = (\alpha_1 \ \alpha_2 \ \alpha_3 \ \alpha_4 \ \alpha_5)$ $\alpha_i \in \mathbb{R}^+$, $i \in \{1, 2, 3, 4, 5\}$

The relations (6) and (7) form linear system of equations which resolution makes it possible to determine the stationary law. By using the Scilab tool (Kaber, 2002) to solve this system, it follows:

$$\alpha = (0.022 \ 0.267 \ 0.343 \ 0.227 \ 0.141) \tag{8}$$

The port 139 is the most attacked port with a probability of 34.34%.

The pie-chart (Fig. 1) makes it possible to visualize the probabilistic most attacked ports. They correspond to the states occupying the greatest portion.

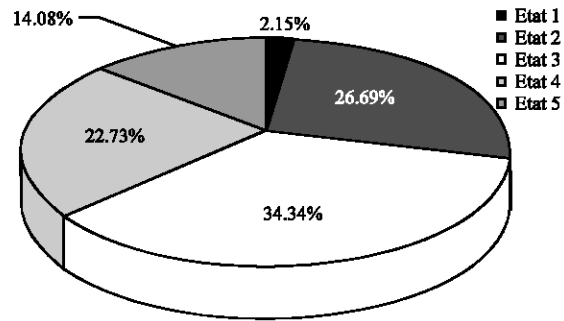


Fig. 1: Analysis of attacks over port

Table 3: Real probabilities for the system to be in its various states at the end of 50 periods

State	1	2	3	4	5
Real probabilities	0.02	0.26	0.32	0.22	0.18

Validation of the model: Here, we will check if the results provided by the model take into account the reality. For this purpose, it is necessary for us to calculate directly from the data of Table 2. The probabilistic state will be calculated at period 50, to ensure convergence (indeed we have seen that starting from a certain period, the probabilistic state of the system becomes stationary).

Direct calculations: The probabilistic state of our system will be calculated in the following way: the occurrence report number where state i occurs in the table during 50 periods on the number of observed periods. If N_i indicates the number of occurrence of state i during 50 periods one has:

$$\text{prob}(X = i) = \frac{N_i}{50} \tag{9}$$

Using the data of Table 2 and formula above, one calculates for each state the probability that the port concerned is attacked the most during 50 periods. The results are presented in the Table 3.

Model validation: By laying out in the same table the results obtained by the model and those obtained by direct calculation at the end of the 50th period, it is possible to appreciate the relevance of the model suggested using an error analysis. By using the results of relation (8) and Table 3, we get the following results of Table 4.

In the Table 4 above, the relative error is calculated for each state by the formula:

$$\Delta = \frac{|P_{exp}^i - \alpha^i|}{P_{exp}^i}$$

Table 4: Estimated errors

States	1	2	3	4	5
Real probabilities	0.020	0.260	0.320	0.220	0.180
Estimated probability	0.021	0.267	0.343	0.227	0.140
Relative error	0.050	0.026	0.072	0.031	0.040

where, P_{exp}^i denote the experimental probability associated with state I.

The vector:

$$(P_{exp}^1 \ P_{exp}^2 \ P_{exp}^3 \ P_{exp}^4 \ P_{exp}^5)$$

is the experimental law of the variable X_n . It is calculated based on the entire sample population (i.e., $n=50$). By using the Table 4, the average error is equal to 0.0438 and the variance is $\sigma = 0.0163$. Thus, it is straightforward that the dispersion effects are weak. The errors are evenly distributed around the average error. We can conclude that the accuracy is good and the model can be considered relevant.

CONCLUSION

This study describes a Markov model to study the evolution of the attacks frequencies against a finished ports based on a honeypot. This model built from the data analysis extracted on CADHo platform project has permitted to distinguish five states from which the system could forward. Moreover the process law has been determined by its transition matrix P and its initial state. The convergence of the chain has also been proved and the stationary law towards which it converges has been determined. This makes it possible to know the probable most ports attacked services at given period. Moreover, one also notes a good adequacy between the results of the model and the data measured. It is significant to note that the model used here neither depends on the platforms from where the data are taken nor to the number of ports concerned.

So, the model can be generalized to any given number of ports.

Our model remains valid for any system having a finished number of ports if the conditions of its application are observed.

REFERENCES

Alata, E., M. Dacier, Y. Deswarte, M. Kaaniche and K. Kortchinsky *et al.*, 2005. Leurré.com: Retour d'expérience sur plusieurs mois d'utilisation d'un pot de miel distribué mondialement. Proceedings of the Symposium Sur La Sécurité des Technologies de l'Information et des Communications, June 1-3, Rennes, France.

Alata, E., M. Dacier, Y. Deswarte, M. Kaaniche and K. Kortchinsky *et al.*, 2006. Collection and analysis of attack dated based on honeypots deployed on the Internet. Advances in Information Security, Springer US. Quality of Protection, pp: 79-91.

Bremaud, P., 2001. Markov Chains: Gibbs field, Monte Carlo Simulations and Queues. 2nd Edn., Springer, NewYork, ISBN: 978-0387985091.

Callegari, C., S. Vaton and M. Pagano, 2008. A new statistical approach to network anomaly detection. Proceedings of the International Symposium on Performance Evaluation of Computer and Telecommunication Systems, Jun. 16-18, Edinburgh, UK, pp: 441-447.

Graham, C., 2008. Chaînes de Markov: Cours. 1st Edn., Dunod, Britain, ISBN-13: 978-2100520831, pp: 274.

Haixiao, C., S. Kulkarni and S. Verdu, 2006. An algorithm for universal lossless compression with side information. IEEE Trans. Inform. Theory, 52: 4008-4016.

Jha, S., K. Tan and R.A. Maxion, 2001. Markov chains, classifiers and intrusion detection. Proceedings of 14th IEEE Computer Security Foundations Workshop, (CSFW'01), IEEE Computer Society, pp: 206-219.

Ju, W.H. and Y. Vardi, 2001. A hybrid high-order markov chain model for computer intrusion detection. J. Comput. Graph. Statistics, 10: 277-295.

Kaaniche, M., Y. Deswarte, E. Alata, M. Dacier and V. Nicomette, 2006. Empirical analysis and statistical modeling of attack processes based on honeypots. Proceedings of the WEEDS 2006-Workshop on Empirical Evaluation of Dependability and Security, June 25-28, Philadelphia, USA., pp: 1-6.

Kaber, S.M., 2002. Introduction à scilab-exercices pratiques corrigés d'algèbre linéaire. Ellipses Marketing, pp: 226.

Lockhart, A., 2004. Network Intrusion Detection Network Security Hack. 2nd Edn., O'Reilly Media Inc., USA., pp: 348-412.

Nong, Y., 2000. A markov chain model of temporal behavior for anomaly detection. Proceedings of the IEEE Workshop on Information Assurance and Security United States Military Academy, June 6-7, West Point, New York, pp: 171-174.

Nuel, G. and B. Prum, 2007. Analyze statistique des séquences biologiques: Modélisation markovienne, alignements et motifs (*Collection bioinformatique*). <http://www.lavoisier.fr/notice/fr418634.html>.

- Oudot, B.L. and V. Glaume, 2003. Global intrusion detection: Prelude hybrid IDS. Technical Report.
- Oumtanaga, S., K.P. Kimou and G.K. Kouadio, 2006. Specification of has model of honeypot attack based on raised data. Proc. World Acad. Sci. Eng. Technol., 17: 207-211.
- Shannon, C.E., 1948. A mathematical theory of communication. Bell Syst. Tech. J., 27: 379-423.
- Spitzner, L., 2002. Honeypots: Tracking Hackers. Addison Wesley, USA., pp: 480.