

Viterbi Algorithm for Detecting DDoS Attacks

Wilson Bongiovanni[‡], Adilson E. Guelfi^{§†}, Elvis Pontes^{*†}, A. A. A. Silva^{*‡}, Fen Zhou[†] and Sergio Takeo Kofuji^{*}

^{*}LSI-POLI, Universidade de Sao Paulo, Brazil, {elvis.pontes, kofuji}@usp.br

[†]CERI-LIA, University of Avignon, Agroparc BP 1228, Avignon, France, fen.zhou@univ-avignon.fr

[‡]IPT, Institute of Technological Research, Sao Paulo, Brazil, wilson.bongiovanni@gmail.com, anderson@ponntes.inf.br

[§]FIPP, Universidade do Oeste Paulista, guelfi@unoeste.br

Abstract—Distributed denial of service attacks aim at making a given computational resource unavailable to users. A substantial portion of commercial Intrusion Detection Systems operates only with detection techniques based on rules for the recognition of pre-established behavioral patterns (called signatures) that can be used to identify these types of attacks. However, the characteristics of these attacks are adaptable, compromising thus the efficiency of IDS mechanisms. Thus, the goal of this paper is to evaluate the feasibility of using the Hidden Markov Model based on Viterbi algorithm to detect distributed denial of service attacks in data communication networks. Two main contributions of this work can be described: the ability to identify anomalous behavior patterns in the data traffic with the Viterbi algorithm, as well as, to obtain feasible levels of accuracy in the detection of distributed denial of service attacks.

I. INTRODUCTION

Nowadays, Distributed Denial of Service (DDoS) attacks affect the availability of computational resources [1, 2]. Those attacks cause financial and institutional losses, violating the code of ethics proposed by the Internet Activities Board (IAB) in [3], or even being classified as violation of law in countries - e.g. in the United States [4] and in the United Kingdom [5].

According to [6], in the first quarter of 2013 the DDoS attacks and the bandwidth consumption of the DDoS attacks have increased 21% and 691% respectively, in comparison with the first quarter of 2012. Intrusion Detection Systems (IDS) can usually be classified as: (1) signature-based detection, which detects intrusions comparing data traffic with a set of predefined rules; and (2) anomaly-based detection, which detects malicious behavior comparing data traffic with a predefined normal behavior modeling [7, 8].

Considering the limitations in recognizing the behavioral changes in data traffic, as well as the occurrence of DDoS attacks in signature-based IDS, other methodologies were already proposed to detect changes in data flows [9, 10]. For instance, Markov models are statistical models for the behavioral analysis of systems, taking into account the observation of their state transitions [11, 12]. An alternative alternative is the use of dynamic programming algorithms [13], e.g. the Viterbi algorithm, together with Markov models.

In this sense, the goal of this paper is to evaluate the feasibility of using the Hidden Markov Model (HMM) based on Viterbi algorithm for detecting DDoS attacks.

The remainder of this paper is organized as follows: the related work on intrusion detection by the use of Markov

models is approached in Section II. Section III introduces the issues that will be approached in this paper - the detection of DDoS attacks. Then, details about the evaluation method used in this paper are reported in Section IV. Section V points out the results obtained by the proposed method. Finally, Section VI summarizes the conclusions of this paper, presenting as well guidelines for possible future works.

II. RELATED WORK

In order to establish an effective plan to detect DDoS attacks, behavioral changes or anomalies in network data traffic should be identified. In [14] an approach based on 2 HMMs was proposed to analyze behavior profiles of network application protocols: the first one considers state transitions through the packet sizes and the second one considers the inter-arrival time of packets. Protocols like FTP (control and data connections), HTTP, HTTP over SSL, SMTP (including outgoing sessions), SSH, Telnet and AOL Instant Messenger (AIM) were monitored using two sources of data flows: packet traces from MIT Lincoln Labs Intrusion Detection Evaluation and TCP sessions from Internet clients in the George Mason University. For the size-based classifier, the best result was achieved with AIM protocol, since HMM correctly classified 78.1% and 81.6% of sequences with block size of 16 bytes and 32 bytes, respectively. For time-based classifier, the best result was also achieved with AIM protocol, since HMM correctly classified 86.4%, 84.2% and 80.2% of sequences with sampling rates of 3, 5 and 7, respectively. For [14] such behavioral changes can be detected by data analysis models.

The main goal in [9] was to propose statistical anomaly detection based on the ARP protocol behavior using a Multi-resolution Hidden Markov Model (MHMM) for the statistical analysis of data collected during the experimental phase. In this work, about 90% of accuracy in detecting anomalies on the ARP protocol data was achieved. However, adjustments in the monitored parameters, such as thresholds of network nodes and sequence sizes of state transitions, increase the accuracy to 99%. Another anomaly detection system for data traffic was also described in [15]. However, this system was based on the implementation of Viterbi algorithm considering commands run on a Solaris 7 operating system. The detection system was able to identify 100% of the occurred attacks. However, although the accuracy was good, the system may

report some differences depending on which dissimilarity distance technique is used.

The work developed by [16] considers the artificial immune system paradigm to detect and deal with intruders, studying its effectiveness for data traffic classification and the techniques for generating anomaly detectors in data traffic to detect DDoS attacks. The system proposed by [16] has shown that the detection times of attacks are considerably significant. Three experiments were carried out with mean detection time of attack close to 1.5 ms for the first one, and more relevant results for the two further experiments, since the attacks were detected, on average, in 0.5 ms and 0.1 ms, respectively.

Table I compares this work with related works. As we can see in Table I, as the main contribution, this work was defined to apply, at the same time, HMM and the Viterbi algorithm to detect specifically the occurrence of DDoS attacks in an updated dataset of computer network data traffic. The updated dataset consists of realistic traffic used in the preparation routines or training of algorithms. In Table I, N/A means non available.

TABLE I
COMPARISON BETWEEN RELATED WORKS.

Paper	HMM	Viterbi	DDoS detection	Up-to-date dataset	2 distance techniques
[14]	Yes	No	No	No	No
[15]	Yes	Yes	No	N/A	Yes
[16]	No	No	No	No	Yes
[9]	Yes	No	No	No	No
our paper	Yes	Yes	Yes	Yes	Yes

III. PROBLEM STATEMENT AND FUNDAMENTAL CONCEPTS

According to [17], most of the commercial IDS work with attack signatures for detecting abnormal data traffic. Although this technique has its effectiveness, the characteristics of DDoS attacks can change (or even new or derived attacks may arise), compromising the performance of a signature-based IDS. Therefore, an additional anomaly-based IDS approach can have good accuracy levels in the intrusion detection process of DDoS attacks, contributing to increase security in computer networks. Fundamental concepts are described in the following subsections.

A. Markov Model and Hidden Markov Model (HMM)

Markov chain is a stochastic process, in which future states depend only on the present state, and are independent of past states. A stochastic process whose state at time t is $X(t)$, for $t > 0$, and whose history of states is given by $x(s)$ for times $s < t$ is a Markov process if [18]:

$$Pr[X(t+h) = y | X(s) = x(s), \forall s \leq t] = \quad (1)$$

$$= Pr[X(t+h) = y | X(t) = x(t), \forall h > 0 \quad (2)$$

Hidden Markov Model (HMM) contains a finite number of unobservable (or hidden) states. As previously defined,

transitions among the states in HMM are also governed by a set of probabilities [11]. As highlighted by [11], HMM are finite state automata that have state transitions according to a probabilistic pattern over time. Thus, considering the inability to predict when the system state transition will occur (i.e. when attack will take place), it is required: (a) the mapping of all the states that the system can adopt, based on the observable variables; and (b) to take into account the probability distribution which can be applied to the states of the observed system.

Representing a HMM, N is the number of model states, $S = S_1, S_2, \dots, S_N$ is the set of individual states, M is the number of observable distinct symbols, $V = V_1, V_2, \dots, V_N$, is the set of individual symbols of the model, $a = a_1, a_2, \dots, a_N$ are the probabilities of the transition of states and $b = b_1, b_2, \dots, b_N$ the output probabilities of the system.

B. Viterbi Algorithm

The Viterbi algorithm is a dynamic programming algorithm for finding the most likely sequence of hidden states, called the Viterbi path, which results in a sequence of observable events related to the HMM (Hidden Markov Model) defined to an autonomous system. Considering the Viterbi path and based on an initial set of probabilities for the system observable parameters, the probabilistic analysis of the system state transitions it feasible, even that the parameters affecting the involved state transitions are only partially observable.

According to [12], the Viterbi algorithm aims to achieve the optimal sequence of state transitions $Q = q_1, q_2, \dots, q_n$, in a given sequence of observations $O = o_1, o_2, \dots, o_n$, having in mind the definition of the best path. In other words, the one that has the highest values in the observed parameters, obtained from:

$$\delta_{t+1}(j) = \max_i [\delta_t(i) a_{ij}] \times b_j(O_{t+1}) \quad (3)$$

To trace the optimal sequence of state transitions, storing the arguments maximized for each instant t is required. Hence, the matrix $\psi_t(j)$ has to be built, and, from the analysis of the sequences of the state transitions, the evaluation of the Viterbi path becomes feasible. Our work applies the following steps, which describe the implementation of Viterbi algorithm proposed in [15]:

Step 1 - Initialization:

$$\delta_t(i) = \pi_i b_i(O_t), 1 \leq i \leq N, \Psi_t(i) = 0 \quad (4)$$

$\delta_t(i)$ is the probability that the observation O_t occurs at time $t = 1$ in the state determined by (i) . The variable $\psi_t(i)$ stores the calculated optimal states.

Step 2 - Recursion:

$$\delta_t(j) = \max_i [\delta_{t-1}(i) a_{ij}] b_j(O_t), 2 \leq t \leq T, 1 \leq j \leq N \quad (5)$$

$$\Psi_t(j) = \operatorname{argmax}_i [\delta_{t-1}(i) a_{ij}] b_j(O_t), 2 \leq t \leq T, 1 \leq j \leq N \quad (6)$$

In this step, $\delta_t(j)$ records accumulate values by the algorithm on the state (i) at time (t) , and $\psi_t(j)$ represents the optimal

state at time $(t - 1)$ which has the lowest amount of state transitions at time (t) .

Step 3 - Computational Calculation:

$$P^* = \max_{S_t \in S} [\delta_T(S)] \quad (7)$$

$$S^*_T = \arg \max_{S_t \in S} [\delta_T(S)] \quad (8)$$

At final time (T) , there are (N) probabilities $\delta_t, t = 1, 2, \dots, N$. The highest probability in the set of analyzed states becomes candidate for inclusion as part of the sequence of optimal states.

Step 4 - Making the Viterbi Path (Backtrack):

$$S^*_t = \Psi_{t-1}(S^*_{t-1}), t = T - 1, T - 2, \dots, 1 \quad (9)$$

The variable $S \times T$ stores the optimal states that were calculated by the algorithm, which, at the end of the calculation, represent the Viterbi path for the set of observable states in traffic, measured in time (T) .

IV. DATASET

For evaluating our proposal (i.e., the feasibility of HMM based on Viterbi algorithm to detect DDoS attacks), we employed the traffic data from the Shiravi's datasets [19]. Table II summarizes the normal and anomalous Data Groups (DG), respectively. According to [19], the Shiravi's dataset intends allowing a more accurate perception of the actual effects of attacks propagated across the network and the respective behavior of the involved nodes, containing traffic data that behave as realistically as possible for normal and anomalous traffic. On the Shiravi's dataset [19], a data sampling and summarization technique is applied based on [20], which depends on the collection of 4 samples per minute from traffic data over 24 hours, representing a total value of 5,760 samples for normal and anomalous DG.

TABLE II
STATISTICAL SUMMARY OF PACKETS COLLECTED BY [19].

Packets	Anomalous DG Jun-15 2010		Normal DG Jun-16 2010	
	Packets	%	Packets	%
Ethernet	35,037,828		24,634,296	
IPv4	35,021,801	99.95 %	24,618,073	99.93 %
IPv6	1,525	0.004 %	1,438	0.06 %
TCP	33,789,362	96.43 %	24,270,717	98.52 %
UDP	1,218,275	3.47 %	345,961	1.4 %
UDpv6	1,447	0.004 %	1,428	0.006 %

V. RESULTS

The Viterbi algorithm is, then, applied over DGs to obtain Viterbi paths. In order to verify if the Viterbi paths obtained for anomalous DG are compatible with the real occurrence of DDoS attacks, Euclidean and Hamming distances are also determined. Then, we calculated the Hamming and Euclidean distances, in order to see whether the model identifies variations in the transitions from states compatible with DDoS attacks. According to [21], the Hamming distance is the

number of different positions in two sequences, and can be represented by:

$$hd(a, b) = \sum_{i=1}^n hd(a_i, b_i), hd(a_i, b_i) = \begin{cases} 0, & \text{if } a_i = b_i \\ 1, & \text{if } a_i \neq b_i \end{cases} \quad (10)$$

The Hamming distances are determined for the sequences of states on the normal and anomalous DG, in order to check how far the sequences are to each other. Here, (a_i) and (b_i) represent the Viterbi paths obtained in the preceding step, and (i) is the number of states included in the calculation step. In analysis of the Hamming distances, the distances from 15:00 (which coincides with the beginning of DDoS attacks) to 21:00 exceeded the threshold value of 0.25, corresponding to the pattern of occurrence of DDoS attacks, as we can see at Figure 1. The same has occurred in the period from 23:00 pm to 24:00 pm (or 00:00) am, respectively. In this case, the threshold value of 0.25 was sufficient to estimate behavioral changes compatible with the occurrence of the DDoS attacks (similar threshold values were obtained by [15] and [16]);

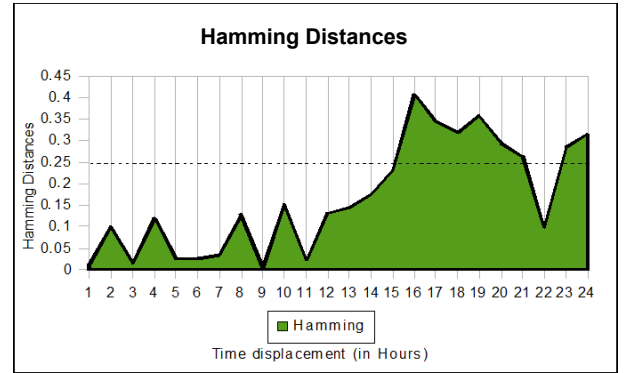


Fig. 1. Results of Hamming Distances.

Subsequently, we performed the calculation of Euclidean distances. Taking into account $x = x_1, x_2, \dots, x_n$ and $y = y_1, y_2, \dots, y_n$ as n -dimensional vectors associated to (a) a Viterbi path from the normal data set, and (b) a sequence states for the Viterbi path calculated for the anomalous data group. The calculation of the Euclidean distance between these two vectors will be the result of the following [16]:

$$ed(x, y) = \sqrt{\sum_{i=1}^N (x_i - y_i)^2}, ed(x, y) \in [0, \infty[\quad (11)$$

Where, N is the number of nodes in the Viterbi path.

The analysis of the Euclidean distance was similar to the analysis of the Hamming distance. However, as depicted by Figure 2, the threshold for the Euclidean distance was 0.8.

We use the same performance metrics as that in [22]. They allow the evaluation of the experimental results that are obtained with the Viterbi algorithm such as: True positive rate (Tpr); False positive rate (Fpr); True negative rate (Tnr);

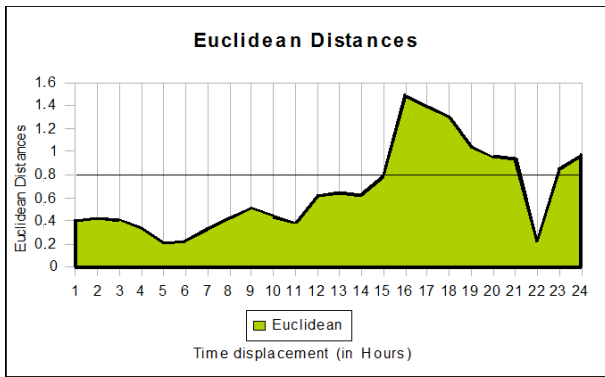


Fig. 2. Results of Euclidean Distances.

False negative rate (Fnr); Accuracy (Ac); Precision (Pr). By applying the same indicators discussed in [22], it is possible to demonstrate the advantages of the proposed method of analysis to determine the occurrence of DDoS attacks, since Tpr for the Hamming Distances is 88.89%, and for the Euclidean Distance is in 77.78% - important to highlight that the Euclidean Distance the Fnr is greater than 22%. Those results demonstrate the behavioral changes, detected by the Viterbi algorithm, is compatible with the occurrence of DDoS attacks. Table III shows such results.

TABLE III
RESULTS CONSIDERING THE INDICATORS PROPOSED BY [22].

Indicators		
Indicator	Hamming distances	Euclidean distances
Tpr	88.89%	77.78%
Tnr	100.00%	100.00%
Fpr	0.00%	0.00%
Fnr	11.11%	22.22%
Ac	95.83%	91.67%
Pr	100.00%	100.00%

VI. CONCLUSION AND FUTURE WORKS

In this paper, we evaluated the feasibility of using the Hidden Markov Model (HMM) based on Viterbi algorithm for detecting DDoS attacks. With validation based on dataset provided by [19], two results for Euclidean and Hamming distances can be highlighted: Tnr = 100%, and Pr = 100%. Thus, the obtained results show HMM based on Viterbi algorithm is preliminary feasible to intrusion detection of DDoS attacks.

As recommendations for future work, we can mention: (a) the use of Viterbi algorithm to analyze larger volumes of data with higher time distances (days, weeks or even months), in order to verify the behavior of Accuracy and Precision indicators; (b) the execution of Viterbi algorithm over datasets that include significant data relating to the IPv6 protocol; (c) the use of a model based on Viterbi algorithm to analyze the on-line network traffic capture, in order to compare response time.

ACKNOWLEDGMENTS

Authors wish to thank the ISCX of the University of New Brunswick, Canada, for sharing the datasets used in this paper. Elvis Pontes is supported by the grant 10581-13-8 - CAPES Foundation, Ministry of Education of Brazil, Brasilia - DF 70040-020, Brazil.

REFERENCES

- [1] S. T. Zargar, J. Joshi, and D. Tipper, "A survey of defense mechanisms against distributed denial of service (DDoS) flooding attacks," *IEEE Comm Surveys & Tutorials*, vol. 15, no. 4, pp. 2046–2069, 2013.
- [2] A. Sperotto, G. Schaffrath, R. Sadre, C. Morariu, A. Pras, and B. Stiller, "An overview of ip flow-based intrusion detection," *IEEE Communications Surveys & Tutorials*, vol. 12, no. 3, pp. 343–356, 2010.
- [3] IAB, "RFC 1087 - Ethics and the Internet," IAB - Internet Activities Board, Tech. Rep., 1989.
- [4] USA, "Title 18 - Crimes and Criminal Procedure Part I - Crimes - Chapter 47 - Fraud and False Statements," <http://www.gpo.gov/fdsys/pkg/USCODE-2010-title18/html/USCODE-2010-title18-partI-chap47-sec1030.htm>, 2010.
- [5] UK, "Office of the Queen's Printer for Scotland - Computer Misuse Act 1990 - Section 3," <http://www.legislation.gov.uk/ukpga/1990/18/section/3>, 1990.
- [6] B. Harris, E. Konikoff, and P. Petersen, "Breaking the ddos attack chain," *Institute for Software Research*, 2013.
- [7] P. Kasinathan, C. Pastrone, M. A. Spirito, and M. Vinkovits, "Denial-of-Service detection in 6LoWPAN based internet of things," in *WiMob*. IEEE, 2013, pp. 600–607.
- [8] K. Scarfone and P. Mell, "Guide to Intrusion Detection and Prevention Systems (IDPS)," *NIST special pub*, vol. 800, no. 2007, p. 94, 2007.
- [9] Y. Yasami, M. Farahmand, and V. Zargari, "An ARP-based anomaly detection algorithm using hidden Markov model in enterprise networks," in *ICSNC*. IEEE, 2007, pp. 69–69.
- [10] I. Forain, A. E. Guelfi, E. Pontes, and A. A. A. Silva, "Intrusion Detection by Time Series Analyzes with ARMAX/GARCH Models," in *SBSeg*, B. C. Society, Ed., dec 2013, pp. 100–113.
- [11] K. Haslum, A. Abraham, and S. Knapkog, "Fuzzy online risk assessment for distributed intrusion prediction and prevention systems," in *UKSIM*. IEEE, 2008, pp. 216–223.
- [12] L. Rabiner, "A tutorial on hidden markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, 1989.
- [13] M. Sniedovich and A. Lew, "Dynamic programming: an overview," *Control and Cybernetics*, vol. 35, no. 3, p. 513, 2006.
- [14] C. Wright, F. Monrose, and G. M. Masson, "HMM profiles for network traffic classification," in *Proceedings of the workshop on Visualization and data mining for computer security*. ACM, 2004, pp. 9–15.
- [15] J.-M. Koo and S.-B. Cho, "Effective intrusion type identification with edit distance for hmm-based anomaly detection system," in *Pattern Recognition and Machine Intelligence*. Springer, 2005, pp. 222–228.
- [16] F. Seredynski and P. Bouvry, "Anomaly detection in TCP/IP networks using immune systems paradigm," *Computer Communications*, vol. 30, no. 4, pp. 740–749, 2007.
- [17] W. Fan, M. Miller, S. Stolfo, W. Lee, and P. Chan, "Using artificial anomalies to detect unknown and known network intrusions," *Knowledge and Information Systems*, vol. 6, no. 5, pp. 507–527, 2004.
- [18] N. Ye, X. Li, Q. Chen, S. M. Emran, and M. Xu, "Probabilistic techniques for intrusion detection based on computer audit data," *IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans*, vol. 31, no. 4, pp. 266–274, 2001.
- [19] A. Shiravi, H. Shiravi, M. Tavallae, and A. A. Ghorbani, "Toward developing a systematic approach to generate benchmark datasets for intrusion detection," *Comp & Sec*, vol. 31, no. 3, pp. 357–374, 2012.
- [20] A. Lazarevic, L. Ertöz, V. Kumar, A. Ozgur, and J. Srivastava, "A comparative study of anomaly detection schemes in network intrusion detection," in *SDM*. SIAM, 2003, pp. 25–36.
- [21] J. Beauquier and Y. Hu, "Intrusion detection based on distance combination," *CESSE07*, 2007.
- [22] S. Kovach, "Detecção de fraudes em transações financeiras via internet em tempo real." Ph.D. dissertation, Universidade de São Paulo, 2012.