



Statistical Learning for Anomaly Detection in Cloud Server Systems: A Multi-Order Markov Chain Framework

Wenyao Sha, Yongxin Zhu , Senior Member, IEEE, Min Chen , Senior Member, IEEE, and Tian Huang, Member, IEEE

Abstract—As a major strategy to ensure the safety of IT infrastructure, anomaly detection plays a more important role in cloud computing platform which hosts the entire applications and data. On top of the classic Markov chain model, we proposed in this paper a feasible multi-order Markov chain based framework for anomaly detection. In this approach, both the high-order Markov chain and multivariate time series are adopted to compose a scheme described in algorithms along with the training procedure in the form of statistical learning framework. To curb time and space complexity, the algorithms are designed and implemented with non-zero value table and logarithm values in initial and transition matrices. For validation, the series of system calls and the corresponding return values are extracted from classic Defense Advanced Research Projects Agency (DARPA) intrusion detection evaluation data set to form a two-dimensional test input set. The testing results show that the multi-order approach is able to produce more effective indicators: in addition to the absolute values given by an individual single-order model, the changes in ranking positions of outputs from different-order ones also correlate closely with abnormal behaviours.

Index Terms— K th-order Markov chain, multivariate time series, anomaly detection, statistical learning

1 INTRODUCTION

CLOUD computing comes with indispensable dependency on networked computer systems. Unfortunately, while every one knows there is no guarantee of its well-being, we tend to simply ignore this painful idea. An increasing number of academical and industrial users are starting to rely solely on cloud computing servers that host entire applications and storage. In fact, these cloud computing and services in the form of distributed and open structure become obvious targets for potential threats. Thus, taking care of both business and personal data, servers expose critical safety and availability issues. Their invulnerability are of major importance to both individuals and the society.

However, during catastrophic disasters such as intrusion, crash or breakdown, the anomalies must be first discovered before any actual remedy could come to its aid. Being recessive at the early stage, such problems would not exhibit distinct traits and often lead to delayed responses and irrecoverable results. Fortunately, a server is the ideal instance whose behavior manifests regularity statistically. This lays the foundation for any anomaly detection algorithm based on

machine-learning or data-mining. All of them adopt the idea of extracting the patterns within the (massive) training set and thus raise the alarm on the deviate ones.

In the literature, anomaly or intrusion detection was implemented in great variety of approaches [1]. These approaches are usually categorized into three groups, i.e. statistical approaches, machine learning approaches and data mining approaches [2], [3]:

In *statistical approaches*, anomaly or intrusion detection systems usually watch behaviors of observed objects to comprise statistical distributions as a set of trained profiles during the training phase. These systems then apply the set of trained profiles by comparing them against a new set of profiles of observed objects during the detection phase. An anomaly or intrusion is detected if these two sets of profiles do not match. In general, any incident whose occurrence frequency goes beyond standard deviations from statistical normal ranges raises an alarm [4].

Machine learning based approaches tend to reduce the super-vised costs during the training phase of statistical approaches by enabling machine learning based systems to learn and improve their performance on their own. These systems are usually designed as a framework that can improve its performance in a loop cycle by adapting its execution strategies according to execution feedbacks, e.g. system call sequence analysis, Bayesian network and Markov model execution results. Neural networks and Hidden Markov Model have been proved to be useful techniques as shown in in [5], [6].

Data mining based approaches exploit unknown rules and patterns by exploring large amounts of data collected either online or off line. Anomaly or intrusion detection systems

- W. Sha, Y. Zhu, and T. Huang are with the School of Microelectronics, Shanghai Jiao Tong University, Shanghai, China.
E-mail: shawenyao@gmail.com, {zhuyongxin, ian_malcolm}@sjtu.edu.cn.
- M. Chen is with the School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan, Hubei, China.
E-mail: minchen@ieee.org.

Manuscript received 15 Oct. 2014; revised 1 Feb. 2015; accepted 22 Feb. 2015.
Date of publication 23 Mar. 2015; date of current version 6 June 2018.

Recommended for acceptance by M. Qiu and S.-Y. Kung.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below.

Digital Object Identifier no. 10.1109/TCC.2015.2415813

...	10:09:10	10:09:30	10:09:32	10:09:33	10:09:34	10:09:40	10:10:00	...
	open()	ioctl()	close()	close()	kill()	fork()	exit()	
	success	success	success	failure	failure	success	success	
		normal				abnormal		

Fig. 1. The time series of system calls.

can be improved with additional inputs in terms of hidden patterns, associations, changes, and events found in data. Common technologies involved in data mining approaches include classification, clustering and outlier detection, and association rule discovery. Application of typical data-mining algorithms such as K-nearest neighbour or clustering could be found in [7].

In the context of this paper, we focus on approaches in categories of statistical and machine learning based approaches. In [8], [9], [10], [11], theories of the classic Markov model are applied so as to detect anomalous patterns in the system, using *the ordering property of events* as proposed in [9]. Ju and Vardi [12] introduces the high-order Markov chain as an extension. Several approaches are then introduced to overcome the formidable cost that seems highly likely to come with it, including the hybrid model in [12] or support vector machine in [13].

In addition to statistical patterns, as originally proposed in [14], the time series of system calls are now by common consent a powerful tool in identifying the nature of a system's behavior. Due to system calls' privileges, a large number of researches of intrusion detection, including [15] and [16] are based on exploiting, modelling, or learning from the audit data of system calls. In 1998, the Cyber Systems and Technology Group of MIT Lincoln Laboratory conducted a seven-week simulation of intrusion in the background of daily usage, and then released all their data named as classic DARPA Intrusion Detection Evaluation Data Set [17].

The major contribution of the paper is our approach based on multi-order Markov chains, which reveals that the combination of mixed-order Markov chains would bring considerably interesting and substantial improvement over any single-order one with fairly reasonable cost. Utilizing not only the multiple order property, this approach effectively suits the application of anomaly detection in addition to its first practice in rainfall modelling in [18]. In our practice, the relative ranking positions between probabilities from multi-order models serve as a new effective indicator for anomalies, which refers to our finding that the ascending order suggests normal, while the descending one exhibits anomalous. Our approach differs from a recent model [19], which exploits mixture of Markov chains by incorporating n-gram transitions to model the normal behavior of users' HTTP requests rather than system calls in underlying servers.

Secondly, we take into accounts a new category of inputs (the return values of system calls) to improve the effectiveness of the multi-order Markov chain based approach and form a two-dimensional model. In the application of anomaly detection, the conventional notion of using system calls to identify a system's behaviour [14], [15], [16] is insufficient in that it does not take into account the resulting status of execution. In a few more recent work [20], [21], return values of system calls were taken into consideration to detect or interpret the anomalies.

In this paper, we further looked into time spans of [17], [20] that were not explored in using our approach based on

multi-order Markov chain scheme. We also provide proof to learning process as part of our scheme framework. The new time span exploration and learning process proof were not covered in [17], [20].

The remainder of the paper is organized as follows. Section 2 formulates the problem and explains the basic theory of the Markov chain model. Section 3 explains our approach to statistical training as well as our high-order models enhanced with multivariate sequence analysis. Section 4 discusses some detailed issues involving algorithm design and complexity. A simplified example of model training and testing is also illustrated in Section 4. Validation results based on the data set is given in Section 5, and a general conclusion is presented in Section 6.

2 PROBLEM FORMULATION AND PRELIMINARIES OF MARKOV CHAIN MODEL

2.1 Problem Formulation

Fig. 1 shows a segment of typical auditing data whose first three lines indicate that the operating system opens or closes a file, reads or writes into the file and creates or terminates a process. On a much larger scale than the system calls executed, we are able to tell the difference between a normal system condition and an abnormal one as shown in the fourth row. However, on the scale as small as the first three lines, how could we learn from the case and differentiate the malicious activities? In practice, these system calls in the first three lines of Fig. 1 might be used by any applications whose intentions cannot be determined as kind or malicious. To handle the uncertainty in anomaly detection due to the single indicator, results from multiple indicators or even multiple models are desired to trigger a proper alarm against prejudicial operations. In this paper, we address the above question with our scheme based on the multi-order Markov chain based model.

To make the forthcoming analysis easier to understand, it is also worth noting that concurrency is ignored on any level. Instead, it is considered as part of the deterministic nature of the auditing mechanism, which decides how every one of the simultaneous system calls gets recorded in a definitive way.

2.2 The Classical Markov Chain

In our approach, we use the discrete-time Markov chain model with finite state space to characterize the stochastic process, and it shares several remarkable similarities with the approaches for anomaly detection as already evidenced by [8], [9], [10], [11], [12], [13].

2.2.1 The Markov Property and Time-Homogeneous Assumption

Theoretically, in order for the model to hold, it is implicitly assumed that all the observation sequences satisfy the

characteristics of both Markov property and time-homogeneous, a reasonable and necessary simplification based on the potential regularity and periodicity of a server system. Consider the state space \mathbf{X} with m discrete states:

$$\mathbf{X} = \{X_1, X_2, \dots, X_m\} \quad (1)$$

and the observation sequence x_n :

$$x_n \in \mathbf{X}, \quad n = 1, 2, \dots, N_x. \quad (2)$$

If a sequence x_n owns Markov property, then it satisfies that the present state results are solely from the last one regardless of transitions the system might have gone through before:

$$\begin{aligned} P(x_n = X_{i_n} | x_{n-1} = X_{i_{n-1}}, \dots, x_2 = X_{i_2}, x_1 = X_{i_1}) \\ = P(x_n = X_{i_n} | x_{n-1} = X_{i_{n-1}}) \text{ for } n \geq 2. \end{aligned} \quad (3)$$

In addition, the property of time-homogeneous (or stationary) ensures that the way how current state affects the next is independent of time:

$$P(x_{n+1} = X_j | x_n = X_i) = P(x_n = X_j | x_{n-1} = X_i). \quad (4)$$

Therefore, every single step of system variation would conform to a unified transition matrix, while keeping it blind to any time-sensitive anomalies.

2.2.2 The Model of Markov Chain

From the training sequence x_n in equation (2), we can derive both the initial probability distribution matrix Q , where each element q_i represents the initial probability of the corresponding state X_i :

$$Q = [q_1 \quad \dots \quad q_i \quad \dots \quad q_m], \quad \text{where } q_i = P(x_n = X_i), \quad (5)$$

and the one-step transition probability distribution matrix P , where each element p_{ij} represents the transition probability from X_i to X_j :

$$P = [p_{ij}]_{m \times m}, \quad \text{where } p_{ij} = P(x_n = X_j | x_{n-1} = X_i). \quad (6)$$

The model is completely specified once P and Q is given, hence we denote it as $\pi(P, Q)$.

3 STATISTICAL LEARNING FRAMEWORK, HIGH-ORDER AND MULTIVARIATE MARKOV CHAIN MODEL

3.1 Our Statistical Learning Framework

To adapt to the usual supervised learning paradigm, for the original state space \mathbf{X} given in (1), we construct the first-order Markov chain transition space \mathbf{X}^1 :

$$\mathbf{X}^1 = \{ [X_1 \ X_1], [X_1 \ X_2], [X_1 \ X_3], \dots, [X_m \ X_m] \}, \quad (7)$$

which consists of every possible transition (totaling m^2) in the first-order Markov process. Similarly, we have the binary state space \mathbf{Y} similar to the fourth row in Fig. 1:

$$\mathbf{Y} = \{ \text{normal}, \text{abnormal} \}. \quad (8)$$

So the entire set of input data \mathbf{D} can be represented as coordinates mapping from space \mathbf{X}^1 to space \mathbf{Y}

$$\begin{aligned} \mathbf{D} = \{ (x_1^1, y_1), (x_2^1, y_2), \dots, (x_{N_x-1}^1, y_{N_x-1}) \}, \\ \text{where } x_n^1 \in \mathbf{X}^1 \text{ and } y_n \in \mathbf{Y} \end{aligned} \quad (9)$$

where the series of x_n^1 is constructed using the original x_n in (2):

$$x_n^1 = [x_n \ x_{n+1}] \in \mathbf{X}^1, \quad n = 1, 2, \dots, N_x - 1. \quad (10)$$

For training purpose, if series x_n (and thus x_n^1) are recorded in a completely anomaly-free environment, all y_n are set as a constant that represents normal system condition, and therefore will be considered irrelevant in the process of training.

3.1.1 Maximum Likelihood Training

In order to learn from the input data set \mathbf{D} , we try to find the model $\pi(P, Q)$ that maximizes the value of

$$P(\mathbf{D} | \pi(P, Q)), \quad (11)$$

which is the probability of data set \mathbf{D} given the model $\pi(P, Q)$ and it can be easily calculated as

$$P(\mathbf{D} | \pi(P, Q)) = \begin{cases} q_{x_1}, & N_x = 1, \\ q_{x_1} \prod_{n=2}^{N_x} P_{x_{n-1}x_n}, & N_x \geq 2, \end{cases} \quad (12)$$

which equals the product of the initial probability and succeeding transition probabilities. To maximize this product, we aim to maximize every multiplier in equation (12) under the constraints of:

$$\sum_{i=1}^m q_i = 1 \quad (13)$$

$$q_i \in [0, 1], \quad i = 1, 2, \dots, m \quad (14)$$

$$\sum_{i=1}^m p_{ij} = 1, \quad j = 1, 2, \dots, m \quad (15)$$

$$p_{ij} \in [0, 1], \quad i = 1, 2, \dots, m, \quad j = 1, 2, \dots, m. \quad (16)$$

From the inequality of arithmetic and geometric means, we know that once the sum of several values are pre-determined, the product of these values will be maximized if and only if all of them are the same, thus:

$$q_i = \frac{\text{The Number of Occurrences of } X_i \text{ in } x_n}{N_x}, \quad (17)$$

$$i = 1, 2, \dots, m$$

$$p_{ij} = \frac{\text{The Number of Occurrences of } [X_i \ X_j] \text{ in } x_n^1}{N_x - 1}, \quad (18)$$

$$i = 1, 2, \dots, m, \quad j = 1, 2, \dots, m.$$

Till now we have already got everything we need for the model $\pi(P, Q)$. Furthermore, we can calculate the probability of any other data set \mathbf{D}' given $\pi(P, Q)$ using the same equation in (12). In a probabilistic sense, the value given in (12) also serves as a quantitative indicator of "anomalousness"

10:09:20	10:09:21	10:09:22	10:09:24	10:09:26
fork()	fork()	kill()	open()	open()
success	success	failure	success	failure

Fig. 2. Normal series.

by measuring how significantly the data set D' used in the equation deviates from the training data D , as will be shown later.

3.2 High-Order and Multivariate Markov Chain Model

As states X_i of Markov chain are intermediate derivatives, there is no specification about what kind of state X_i really is. Therefore, by constructing new states out of the original ones, we can derive a generalized variation that applies to the high-order or multivariate conditions.

3.2.1 Our High-Order Markov Chain Model

Given sequence x_n that satisfies the K th-order Markov property, where the current state depends on not only the previous statuses but also the K preceding ones:

$$\begin{aligned} P(x_n = X_{i_n} | x_{n-1} = X_{i_{n-1}}, \dots, x_2 = X_{i_2}, x_1 = X_{i_1}) \\ = P(x_n = X_{i_n} | x_{n-1} = X_{i_{n-1}}, \dots, x_{n-K} = X_{i_{n-K}}) \quad (19) \\ \text{for } n \geq K + 1, \end{aligned}$$

we construct new sequence x_n^* from x_n :

$$x_n^* = \underbrace{[x_{n-K+1} \ \dots \ x_n]}_K \in \mathbf{X}^*, \quad n = K, K + 1, \dots, N_x. \quad (20)$$

As a sequence directly derived from x_n , x_n^* is the first-order equivalent of x_n which satisfies the classical Markov property defined in equation (3). The new state space \mathbf{X}^* comprises m^K states in total:

$$\mathbf{X}^* = \{X_1^*, \dots, X_i^*, \dots, X_{m^K}^*\}, \quad (21)$$

where each state X_i^* is a certain permutation of the K original states:

$$X_i^* = [X_{\theta_K(i)} \ \dots \ X_{\theta_j(i)} \ \dots \ X_{\theta_1(i)}]. \quad (22)$$

The permutation $\theta_j(i)$ could be defined as a base m number up to K -digits:

$$\begin{aligned} \theta_j(i) = \frac{1}{m^{j-1}} \left\{ \left[i - \sum_{\omega=1}^{j-1} (\theta_\omega(i) m^{\omega-1}) \right] \bmod m^j \right\}, \quad (23) \\ j = 1 \dots K. \end{aligned}$$

Similarly to equation (5), the initial matrix of all m^K initial states is as follows:

$$Q^* = [q_1^* \ \dots \ q_i^* \ \dots \ q_{m^K}^*], \quad (24)$$

where q_i^* represents the initial probability distribution of X_i^* :

$$q_i^* = P(x_n^* = X_i^*) \text{ for } n \geq K. \quad (25)$$

Slightly different from the classical model, the total amount of possible transitions is limited to $m^K \times m$ instead

10:09:20	10:09:21	10:09:22	10:09:24	10:09:26
fork()	fork()	kill()	open()	open()
success	success	failure	success	success

Fig. 3. Test series 1.

of the entire $m^K \times m^K$ ones in the first-order condition, so the new transition matrix would be:

$$P^* = [p_{ij}^*]_{m^K \times m}, \quad (26)$$

where:

$$p_{ij}^* = P(x_n = X_j | x_{n-1}^* = X_i^*). \quad (27)$$

Therefore, the model in the K th-order case is $\pi(P^*, Q^*)$.

Similar to the first-order Markov chain transition space \mathbf{X}^1 in (28), we construct the K th-order Markov chain transition space \mathbf{X}^K :

$$\begin{aligned} \mathbf{X}^K = \{ & \underbrace{[X_1 \ X_1 \ \dots \ X_1 \ X_1 \ X_1]}_K, \\ & \underbrace{[X_1 \ X_1 \ \dots \ X_1 \ X_2 \ X_1]}_K, \\ & \dots, \\ & \underbrace{[X_m \ X_m \ \dots \ X_m \ X_m \ X_m]}_K \}, \quad (28) \end{aligned}$$

which consists of every possible transition (totaling $m^K \times m$) in the K th-order Markov process. The label space \mathbf{Y} remains unchanged. Now input data set D can also be represented as coordinates mapping from space \mathbf{X}^K to space \mathbf{Y}

$$\begin{aligned} D = \{ (x_1^K, y_1), (x_2^K, y_2), \dots, (x_{N_x-K}^K, y_{N_x-K}) \}, \quad (29) \\ \text{where } x_n^K \in \mathbf{X}^K \text{ and } y_n \in \mathbf{Y}, \end{aligned}$$

where the series of x_n^K is constructed using the original x_n in (2) and x_n^* in (20):

$$x_n^K = [x_n^* \ x_{n+1}^*] \in \mathbf{X}^K, \quad n = K, K + 1, \dots, N_x - 1. \quad (30)$$

Correspondingly, equation (12) becomes:

$$P(D | \pi(P^*, Q^*)) = \begin{cases} q_{x_n^K}^*, & N_x = K, \\ q_{x_n^K}^* \prod_{n=K+1}^{N_x} p_{x_{n-1}^* x_n^*}^*, & N_x \geq K + 1. \end{cases} \quad (31)$$

3.2.2 Enhancement with Multivariate Sequences

Given multiple sequences for anomaly detection, it is intuitive to apply the model to each of them in separation. In this way, the hidden correlation among multiple variables in these sequences are lost unfortunately, which might be the very factor that contributes to anomalies. For example, consider the following case, where Fig. 2 is a piece of time series collected in anomaly-free context, and we are evaluating the likelihood of the test series in Figs. 3 and 4 being also normal.

10:09:20	10:09:21	10:09:22	10:09:24	10:09:26
fork()	fork()	kill()	open()	open()
failure	failure	success	failure	sucess

Fig. 4. Test Series 2.

From a stand-alone point of view, all the transitions of system calls and return values are legitimate, as their occurrences have already been observed in the training series. The only difference is that Fig. 3, being almost the same as the training series, may yield a higher likelihood than Fig. 4. However, a closer look at the joint states of system calls and return values tells us that none of the transitions in Fig. 4 are possible. For example, the fork command with the return value of failure can never transit into another fork with the return value of failure in the normal series, which means Fig. 4 actually stands for something more extreme than the transitions along each dimension suggests. In other words, looking at only one dimension of system calls or return values at a time gives us less information than looking at them jointly.

Therefore, combing these variables into a new space would be a practically superior enhancement. Combining multiple inter-related sequences as a multivariate into a single model would be another feasible approach to improve the sensitivity of detection. The example shown in Figs. 2, 3 and 4 is not detectable by major existing methods in the form of statistical pattern detection and time series analysis [8], [9], [10], and [15].

For a S -dimensional sequence, if the original state spaces are noted as $\mathbf{X}_1, \mathbf{X}_2$ to \mathbf{X}_S and each consists of m_1, m_2 or m_S states respectively, the corresponding state space of combination comprises totally m^S states:

$$m^S = \prod_{i=1}^S m_i, \quad (32)$$

and the new state space $\mathbf{X}(S)$ is given by:

$$\mathbf{X}(S) = \{X_1(S), \dots, X_i(S), \dots, X_{m^S}(S)\} \quad (33)$$

each state $X_i(S)$ of which is a combination of the original states:

$$X_i(S) = [X_{i_1} \ \cdots \ X_{i_j} \ \cdots \ X_{i_S}], \quad (34)$$

where $1 \leq i_j \leq m_j$ and $j = 1, 2, \dots, S$,

where each element belongs to one of the original state spaces:

$$X_{i_j} \in \mathbf{X}_j, \quad j = 1, 2, \dots, S. \quad (35)$$

Hence, the multivariate sequences are able to comply with the classical univariate Markov chain model $\pi(P, Q)$ described in previous sections.

Moreover, the variation of high-order and multivariate model could be applied jointly, creating a S -dimensional K th-order Markov chain model $\pi[P^*(S), Q^*(S)]$ with S -dimensional K th-order Markov chain transition space $\mathbf{X}^K(S)$.

Though combination of sequences often results in higher time and space complexity, in our practice of algorithm

design discussed later, both of these generalizations are implemented with complexity reduction techniques as the Multi-order Markov Chain based model, where we utilize the system calls and their return values as a two-dimensional input ($S = 2$) of the first, second and third-order Markov chain ($K = 1, 2, 3$).

4 ALGORITHM DESIGN AND IMPLEMENTATION

To formulate our algorithm, we follow the process of anomaly detection workflow comprising of training stage and test stage. On top of the generic training and testing stages, we further tackle the issues of model order selection, zero-probability event transform and sparse matrix storage etc. Our solutions to these issues will be explained before the complete algorithm is formulated.

4.1 Design Issues

Order of the model in our approach is the trade-off between model complexity and generalizability. System of higher-order tend to neutralize the idealization brought by the assumption of Markov property and thus provides a more sensitive response over anomalous exceptions. In other words, high-order Markov chain generates larger support for frequent patterns within the training set while producing smaller support for rare ones at the same time. On the contrary, high-order models might lead to overfitting for noise or idiosyncrasy of the training sequence, and loses the ability to generalize outside it. In consequence, Markov chain models of up to the third order are utilized in our simulation.

Zero-probability-event issues refer to the fact that due to the non-exhaustive nature of any training sequence, one or more elements in P or Q will remain zero after training, and a single zero in the multipliers would ultimately dominate the whole output value. Here a solution proposed in [8] is implemented by replacing the zeros with a value far smaller than any other element (yet still essentially different from the real zero).

Another typical issue resulting from the non-exhaustiveness of training set is the occurrence of completely new states in test sets. Even after going through all the training data, what we consider as the state space \mathbf{X} could just be a subset of the real one because some states simply might not have occurred at all in the recording of training set. So for application purpose, we add a new artificial state X_{m+1} to the original space \mathbf{X} to represent "everything else". In this way, without compromising any accuracy of calculation, state space \mathbf{X} will be independent of the test set and we will not have to deal with an growing \mathbf{X} as we proceed.

4.2 Algorithm Design

The final algorithms that take both zero-probability events and storage optimization into consideration are presented as Algorithm 1-4. Algorithm 1, *Training*, represents the training stage which processes the training set and generates the initial probability distribution and transition probability distribution matrices, of which the latter one is the major result of the training stage. Algorithm 2, *Testing*, takes the two matrices and the test set in and calculates its probability of occurrence given the model trained. Algorithm 1

and Algorithm 2 also rely on two functions i.e. *IncreaseTransitionMatrix* and *GetTransitionMatrix* to calculate transition probabilities and retrieve values in these matrices.

Algorithm 1. Training()

Input: Training sequence $x_n \in \mathbf{X}$ ($n = 1 \rightarrow N_x$), Number of states m , Order K

Output: Initial probability distribution matrix Q , Transition probability distribution matrix P

- 1: Initialize() // Set Q and P equal to 0
- 2: **for** $n = K \rightarrow N_x$ **do**
- 3: $x_n^* \leftarrow [x_{n-K+1} \quad x_{n-K+2} \quad x_{n-K+3} \quad \dots \quad x_n]$
// Construct new sequence in space \mathbf{X}^*
- 4: **end for**
- 5: **for** $n = K \rightarrow N_x - 1$ **do**
- 6: $x_n^K \leftarrow [x_n^* \quad x_{n+1}^*]$
// Construct new sequence in space \mathbf{X}^K
- 7: **end for**
- 8: **for** $n = K \rightarrow N_x$ **do**
- 9: **if** $x_n = X_i$ **then**
- 10: IncreaseInitialMatrix(i)
// Get the index of x_n and increase q_i^* in Q by 1
- 11: **end if**
- 12: **end for**
- 13: **for** $n = K \rightarrow N_x - 1$ **do**
- 14: **if** $x_n^K = [X_i^* \quad X_j^*]$ **then**
- 15: IncreaseTransitionMatrix(i, j)
// Get the index of x_n^K and increase p_{ij}^* in P by 1
- 16: **end if**
- 17: **end for**
- 18: NormalizeInitialMatrix() // Divide the matrix Q by $N_x - K$
- 19: NormalizeTransitionMatrix()
// Divide the matrix P by $N_x - K - 1$

Algorithm 2. Testing()

Input: Testing sequence $y_n \in \mathbf{X}$ ($n = 1 \rightarrow N_y$), Q, P

Output: Probability of y_n given Q and P

- 1: **for** $n = K \rightarrow N_y$ **do**
- 2: $y_n^* \leftarrow [y_{n-K+1} \quad y_{n-K+2} \quad y_{n-K+3} \quad \dots \quad y_n]$
// Construct new sequence in space \mathbf{X}^*
- 3: **end for**
- 4: **for** $n = K \rightarrow N_y - 1$ **do**
- 5: $y_n^K \leftarrow [y_n^* \quad y_{n+1}^*]$
// Construct new sequence in space \mathbf{X}^K
- 6: **end for**
- 7: **for** $n = K \rightarrow N_y$ **do**
- 8: **if** $y_n = X_i$ **then**
- 9: probability \leftarrow GetInitialMatrix(i) // Get q_i^* in Q
- 10: **end if**
- 11: **end for**
- 12: **for** $n = K \rightarrow N_y - 1$ **do**
- 13: **if** $y_n^K = [X_i^* \quad X_j^*]$ **then**
- 14: probability \leftarrow probability \times GetTransitionMatrix(i, j)
// Multiply by p_{ij}^* in P
- 15: **end if**
- 16: **end for**
- 17: Output : probability

Algorithm 3 and Algorithm 4 are two functions called during training and testing, which either increases a specific element of transition matrix by one or returns its value.

Algorithm 3. IncreaseTransitionMatrix(i, j)

Input: The index of state in transition(i and j)

Output: Transition matrix (stored in boolean values), Nonzero value table

- 1: MAXNONZERO $\leftarrow 1 \times 10^4$
// Define max number of nonzero values
- 2: **if** TransitionMatrixBit[i][j] = 0 **then**
// If the nonzero value does not exist
- 3: CounterOfNonzeroValues \leftarrow CounterOfNonzeroValues + 1
// Increase counter
- 4: **if** CounterOfNonzeroValues > MAXNONZERO
then // If exceeded
- 5: Print : Error! Please increase MAXNONZERO!
// Print message
- 6: exit // Algorithm aborted
- 7: **end if**
- 8: iOfNonzeroValues[CounterOfNonzeroValues] $\leftarrow i$
// Set index of nonzero value
- 9: jOfNonzeroValues[CounterOfNonzeroValues] $\leftarrow j$
// Set index of nonzero value
- 10: NonzeroValues[CounterOfNonzeroValues] $\leftarrow 1$
// Set initial value
- 11: TransitionMatrixBit[i][j] $\leftarrow 1$
// Set sign of the nonzero value
- 12: **else** // If the nonzero value already exists
- 13: **for** $n = 1 \rightarrow$ CounterOfNonzeroValues **do**
- 14: **if** iOfNonzeroValues[n] = i AND jOfNonzeroValues[n] = j
then
- 15: break // Find where the nonzero value is stored
- 16: **end if**
- 17: **end for**
- 18: NonzeroValues[n] \leftarrow NonzeroValues[n] + 1
// Increase the value by 1
- 19: **end if**

Algorithm 4. GetTransitionMatrix(i, j)

Input: The index of state in transition(i and j), Nonzero value table, Transition matrix (stored in binary bits)

Output: Probability of the transition from state i to state j

- 1: ZERO $\leftarrow 1 \times 10^{-5}$
- 2: **if** TransitionMatrixBit[i][j] = 0 **then** // If it is a zero
- 3: **return** ZERO
- 4: **else** // If it is a nonzero value
- 5: **for** $n = 1 \rightarrow$ CounterOfNonzeroValues **do**
- 6: **if** iOfNonzeroValues[n] = i AND jOfNonzeroValues[n] = j
then
- 7: **return** NonzeroValues[n]
// Find where the nonzero value is stored
- 8: **end if**
- 9: **end for**
- 10: **end if**
- 11: **end for**
- 12: **end if**

4.3 Time and Space Complexity

The training stage involves calculating the initial and transition matrices and takes approximately

$$O(m^{K+1}) \quad (36)$$

10:09:20	10:09:21	10:09:22	10:09:24	10:09:26
kill()	fork()	kill()	fork()	open()
success	failure	success	failure	failure
		normal		

Fig. 5. System call series of Virtual Machine 1 (training set).

's time, which grows either polynomially with the amount of states or exponentially with the order of Markov chain. The test stage requires nothing more than a series of multiplication and consumes about

$$O(N_x) \quad (37)$$

's time, which is linearly proportional to the length of the testing data set.

On the other hand, the algorithm requires the matrix P^* and Q^* be accessible at all time during the training and testing stage, so it can consume memory space as large as

$$O(m^{K+1}). \quad (38)$$

4.4 Curbing the Complexity

The training algorithm, acting as a forefront overhead for the testing algorithm, does not affect the real-time performance of anomaly detection. Also, the testing algorithm is simply a serial of continual multiplication which will hardly become the burden of any modern processor. Even in the cases where multiplying are considered unacceptably time-consuming, it would be highly feasible to replace both transition and initial matrices with their logarithm values in the training stage beforehand, causing all multiplication of the test stage substituted by addition, one of the basic commands supported by any instruction set architecture.

As the amount of states and the order of model grow, the memory consumption for storing transition matrix alone could be formidable. However, when taking into account the sparse feature of the matrix (which assumes a significant amount of potentially possible transitions does not occur in reality), a binary bit of memory space is already sufficient for most of its elements. Together with the extra table declared exclusively for saving all nonzero values and their indices (with subjectively pre-determined array size), the revised algorithm nonetheless succeeds in distinctly optimizing space complexity at the cost of limited sacrifice in time complexity. In our practice, the length of the non-zero value table is set to 10^4 and no overflow has occurred, another proof of how such a look-up approach manages to save storage consumption without compromising any accuracy.

4.5 Implementation Illustration with an Example

For clarity, an example of model training and testing is shown step by step as follows. Suppose the training data set

10:09:20	10:09:21	10:09:22	10:09:24	10:09:26
fork()	fork()	kill()	open()	open()
failure	failure	success	failure	success
		normal		

Fig. 6. System call series of Virtual Machine 2 (training set).

TABLE 1
Mapping from $[RV_1 SC_1 RV_2 SC_2]$ to \mathbf{X}

$[RV_1 SC_1 RV_2 SC_2]$	\mathbf{X}
$[a_1 b_1 a_1 b_1]$	X_1
$[a_1 b_1 a_1 b_2]$	X_2
$[a_1 b_1 a_1 b_3]$	X_3
...	...
$[a_2 b_3 a_2 b_1]$	X_{34}
$[a_2 b_3 a_2 b_2]$	X_{35}
$[a_2 b_3 a_2 b_3]$	X_{36}

\mathbf{D} , which are the series of system calls and return values from a two-virtual-machine system in its normal running state, are shown in Figs. 5 and 6, and a two-dimensional ($S = 2$) model of second-order ($K = 2$) is desired.

First, the two-state space A and three-state space B could be obtained and noted as:

$$A = \{success, failure\} = \{a_1, a_2\}, \quad (39)$$

$$B = \{open(), kill(), fork()\} = \{b_1, b_2, b_3\}. \quad (40)$$

For simplicity, the training sequences of return values and system calls derived directly from Figs. 5 and 6 are noted as RV_1, SC_1, RV_2 and SC_2 respectively:

$$RV_1 = a_1, a_2, a_1, a_2, a_2 \quad (41)$$

$$SC_1 = b_2, b_3, b_2, b_3, b_1 \quad (42)$$

$$RV_2 = a_2, a_2, a_1, a_2, a_1 \quad (43)$$

$$SC_2 = b_3, b_3, b_2, b_1, b_1. \quad (44)$$

Now we create the table for mapping from the input multivariate sequence to a univariate one as shown in Table 1, and construct a univariate equivalent sequence x_n (and the implicit state space \mathbf{X}):

$$x_n = \{a_1, b_2, a_2, b_3\}, \{a_2, b_3, a_2, b_3\}, \\ \{a_1, b_2, a_1, b_2\}, \{a_2, b_3, a_2, b_1\}, \\ \{a_2, b_1, a_1, b_1\} = X_{12}, X_{36}, X_8, X_{34}, X_{19} \quad (45)$$

Then, the first-order equivalent sequence x_n^* (and the implicit state space \mathbf{X}^*) is given using the mapping in Table 2:

$$x_n^* = \{X_{12}, X_{36}\}, \{X_{36}, X_8\}, \{X_8, X_{34}\}, \{X_{34}, X_{19}\} \\ = X_{432}^*, X_{1268}^*, X_{286}^*, X_{1207}^*. \quad (46)$$

So,

$$x_n^K = [X_{432}^*, X_8], [X_{1268}^*, X_{34}], [X_{286}^*, X_{19}]. \quad (47)$$

Evidently, the initial matrix Q^* (1 by 6×6) with only four nonzero elements ($q_{286}^* = q_{432}^* = q_{1,207}^* = q_{1,268}^* = 1/4$):

$$Q^* = [\dots \frac{1}{4} \dots \frac{1}{4} \dots \frac{1}{4} \dots \frac{1}{4} \dots]_{1 \times 1,296} \quad (48)$$

TABLE 2
Mapping from $[x_n x_n]$ to X^*

$[x_n x_n]$	X^*
$[X_1 X_1]$	X_1^*
$[X_1 X_2]$	X_2^*
$[X_1 X_3]$	X_3^*
...	...
$[X_{36} X_{34}]$	X_{1294}^*
$[X_{36} X_{35}]$	X_{1295}^*
$[X_{36} X_{36}]$	X_{1296}^*

10:19:30	10:19:31	10:19:32
fork()	kill()	fork()
failure	success	failure

Fig. 7. System call series of Virtual Machine 1 (test set 2).

10:19:30	10:19:31	10:09:32
fork()	kill()	open()
failure	success	failure

Fig. 8. System call series of Virtual Machine 1 (test set 1).

and the transition matrix P^* (36×36 by 36) with only three nonzero elements ($p_{432,8}^* = p_{1,268,34}^* = p_{286,19}^* = 1$):

$$P^* = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & \cdots & \cdots & \cdots & \cdots & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & \cdots & \cdots & 1 & \cdots & \cdots & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & 1 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}_{1,296 \times 36} \quad (49)$$

are established, which means we have the model $\pi(P^*, Q^*)$. Therefore, the probability of every test sequence constructed this way could be determined using equation (31). For instance, given newly-observed series D' shown in Figs. 8 and 9 as the test input sets, again the test sequences derived are noted as $RV'_1, SC'_2, RV'_1, SC'_2$:

$$RV'_1 = a_2, a_1, a_2, \quad (50)$$

$$SC'_1 = b_3, b_2, b_3, \quad (51)$$

$$RV'_2 = a_2, a_1, a_2, \quad (52)$$

$$SC'_2 = b_3, b_2, b_1. \quad (53)$$

Similarly, we have:

$$\begin{aligned} x'_n &= \{a_2, b_3, a_2, b_3\}, \{a_1, b_2, a_1, b_2\}, \{a_2, b_3, a_2, b_1\} \\ &= X_{36}, X_8, X_{34}, \end{aligned} \quad (54)$$

$$x_n^* = \{X_{36}, X_8\}, \{X_8, X_{34}\} = X_{1268}^*, X_{286}^* \quad (55)$$

$$x_n^K = [X_{1268}^*, X_{34}]. \quad (56)$$

10:59:30	10:59:31	10:59:32
open()	fork()	kill()
success	success	failure

Fig. 9. System call series of Virtual Machine 2 (test set 1).

10:59:30	10:59:31	10:59:32
kill()	fork()	open()
failure	success	success

Fig. 10. System call series of Virtual Machine 2 (test set 2).

According to equation (31), the probabilities of series X^* given the model $\pi(P^*, Q^*)$ are:

$$P(D' | \pi(P^*, Q^*)) = q_{1268}^* \times p_{1268,34}^* = 1 \times 1 = 1. \quad (57)$$

On the other hand, for newly-observed series shown in Figs. 9 and 10, by applying the same methodology as test set 1 and 2, their probability of occurrence given the model $\pi(P^*, Q^*)$ equals 0. As a result, the latter test set is expected to be more "abnormal" than the former one since the latter test set has a lower probability of support than the former one does, which means decision can be made according to the probability of support.

In addition, the ordinal indices of the new state space could conform to any mapping other than Tables 1 and 2 as long as it maintains uniqueness.

5 VALIDATION RESULTS AND DISCUSSION

As a real-world application of the multi-order Markov chain anomaly detecting framework, we take in the classic DARPA Intrusion Detection Evaluation Data Set [17] by the Cyber Systems and Technology Group of MIT Lincoln Laboratory to verify our scheme. The DARPA dataset were collected since June 1998 when DARPA conducted a seven-week simulation of TCP attacks (anomalies) in the background stream of normal and regular user operations. Major results for comparison [8], [9], [10] and [15] were achieved with this data set. The rest data sets in relevant literature were either unavailable to the public or even older than DARPA data set. Hence, we just choose DARPA data set. In [11], authors conducted their experiments on genomic data rather than security data sets. In [12], data were collected from the output of the UNIX acct auditing mechanism on a local machine at authors lab, which were not available to the public. In [16] published in 2006, an even earlier data set (CERT synthetic Sendmail data collected at the University of New Mexico (UNM) by Forrest et al. (1996)) was used. In [19], an unpublished data set on authors own web servers was used, which was not available for cross verification.

With the help of Basic Security Module (BSM) provided by Oracle Solaris, system information including event time, types of system call executed, the return value, et al. were audited during the experiment. A typical piece of BSM audit data can be found in Fig. 11.

...

```

header,150,2,ioc1(2),,Fri Jun 05 07:49:06 1998, + 499467935 msec
path,/etc/security/audit control
attribute,100664,root,other,8388608,62781,0
argument,2,0x5401,cmd
argument,3,0xefffed7c,arg
subject,2122,root,other,root,other,257,257,0 0 pascal.eyrie.af.mil
return,failure: Inappropriate ioc1 for device,-1
trailer,150

```

...

Fig. 11. A typical segment of BSM audit data.

5.1 Validation Results Regarding Multi-Order Markov Chain Model

For training set D , we apply Algorithm 1, *Training*, on all the 819,472 data of system call sequence recorded on the first week's Friday, while data from other Fridays are used as input for Algorithm 2, *Testing*, respectively. A more detailed summary of input parameters and execution could be found in Tables 3 and 4. It is also noteworthy that although BSM is able to monitor up to 243 types of system calls, only 57 of them were encountered during the training stage. In this case, treating all other unencountered system calls as a single category would result in a simplified yet equivalent state space comprising totally 58 states.

It is worth noting that the data inputs in Tables 3 and 4 covered 24 hours of the first week's Friday. This amount of data only takes up to 56 seconds to complete our testing detection. If we feed a segment of data covering only half an hour into our detection scheme for testing, it is expected to take around 1 second to complete the testing. The reason behind the quick response in our detection scheme is that kernel calculations involved in the testing phase only include looking up for the probability of the new states and multiplying sets of probabilities together. In other words, our scheme is able to be deployed for on-line anomaly detection.

Calculation results of week 2 to 7 using the multi-order Markov chain model are shown in Fig. 12, where the horizontal axis represents the ordinal number of a subset of testing sequence with 200 data (a sliding window over the series with a width of 200), while the vertical axis represents the negative logarithmic values of the probability of the subset of test series $D_{testing}$ given the model $\pi(P^*, Q^*)$ just trained:

$$-\log \{ P[D_{testing} | \pi(P^*, Q^*)] \}. \quad (58)$$

TABLE 3
Input Parameters of Whole Day

Parameter	Value
Orders	1,2,3
Training set amount	819,472
Test set amount	100,000 × 6
Sliding window width	200
ZERO	1 × 10 ⁻⁵

The blue, red and green curves stand for the first, second and third-order model respectively.

As many figures show, for most of the time, the calculation results are maintained at a relatively low level, following the descending order from the first-order model to higher ones. Take the third-order model in Fig. 12a for instance, the green curve usually oscillates between 20 and 80, which indicates an average one-step transition probability of approximately 0.4 to 0.8 due to the inverse function of negative logarithm as shown in

$$(10^{-20})^{1/200} \approx 0.80, \quad (59)$$

$$(10^{-80})^{1/200} \approx 0.40. \quad (60)$$

In view of the total 58 possible transitions, a probability of 0.4 to 0.8 would by no means be regarded as a "rare" event and thus agrees with the universal rule of the normality having a greater probability. However, at some point, the curve rises sharply, reaching roughly 330, or

$$(10^{-330})^{1/200} \approx 0.02 \quad (61)$$

in terms of average one-step transition probability somewhere in the tail of Fig. 12a, which is significantly smaller than any one-step transition probability encountered before. At the same time of the "spike", an abrupt reverse of the originally descending order from the first to higher-order model has also occurred. Moreover, the sixth Friday shown in Fig. 12e witnesses a much more dramatic situation in which the third-order result soars to 1,000, or an average one-step transition probability of

$$(10^{-1000})^{1/200} = 10^{-5}, \quad (62)$$

which is the exact same value as the pre-defined constant ZERO in Table 3. This means that for these particular

TABLE 4
Multi-Order Execution Summary for Data of Whole Day

Order	1	2	3
Number of states	58	3,364	1,95,112
Number of transitions	3,364	195,112	11,316,496
Typical runtime (training)	3 s	3 s	5 s
Typical runtime (test)	25 s	38 s	56 s

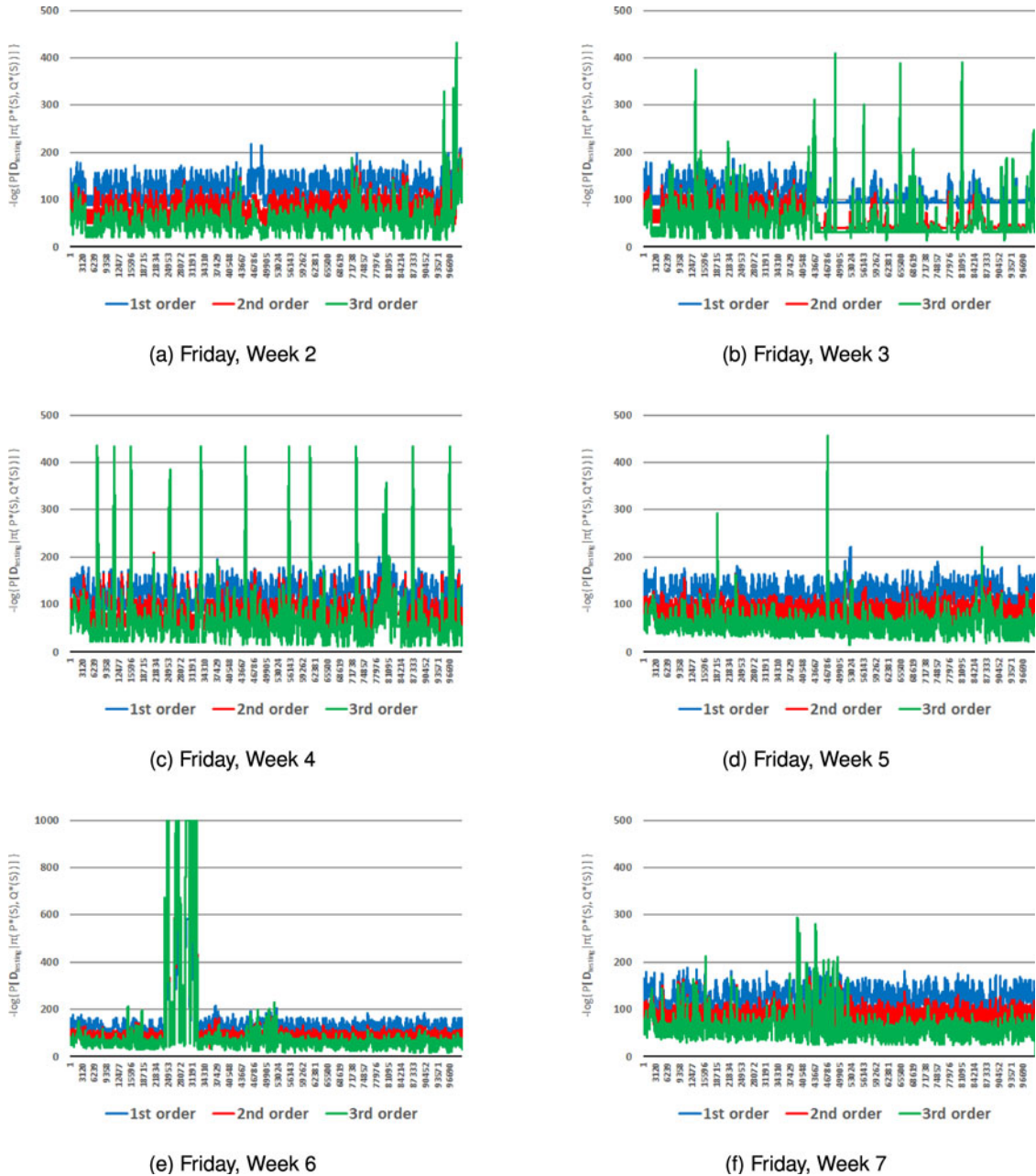


Fig. 12. Calculation results from Friday, Week 3 to 7, where the horizontal axis represents the ordinal number of a sub training sequence with 200 data, and the vertical axis represents the negative logarithmic values of support from the model using the first Friday as the training input. The blue, red and green curves stand for the first, second and third-order model respectively. In other words, the specific steps of generating these results are: 1). Input the data from Friday, Week 1 into Algorithm 1, and by counting the occurrence of every possible state and transition, set up the initial and transition matrix; 2). Input the data from other Fridays into Algorithm 2, and multiply the corresponding elements of the former two matrices to get the conditional probability of the test set given the model; 3). For convenience, the negative logarithmic value are taken as output. Note that all following figures are presented in the same way.

subsets of the testing set, every single transition in them has gone from “rare” or “low probability” to “impossible” or “never happened before”.

According to the record by Lincoln Lab, both of the extreme values depicted in Figs. 12a and 12e coincide with a TCP attack chronologically, whereas the latter one resulted in an unprecedented system crash on July 10th, 1998.

5.2 Validation Results Regarding Multivariate Sequence

Whenever a system call is executed, a certain binary value of either SUCCESS or FAILURE is returned, which can

also be obtained from BSM as shown in Fig. 11. Along with the types of system calls processed above, the return values contribute to forming a kind of simplest two-dimensional time series in a server system. Accordingly, the summary of input parameters and execution can be found in Tables 5 and 6.

The multivariate computing results of the same Fridays in previous figures are shown in Fig. 13. Now the vertical axis has a slight new meaning of the negative logarithmic values of the probability of the subset of test series $\mathbf{D}_{testing}$ given the multi-variate model $\pi[P^*(S), Q^*(S)]$:

TABLE 5
Multivariate Input Parameters

Parameter	Value
Orders	1,2,3
Dimension	2
Training set amount	819,472×2
Test set amount	(100,000×2)×6
Sliding window width	200
ZERO	1×10^{-5}

TABLE 6
Multivariate Execution Summary

Order	1	2	3
Number of states	116	3,364	1,560,896
Number of transitions	13,456	1,560,896	181,063,936
Typical runtime (conversion)	15 s		
Typical runtime (training)	3 s	4 s	5 s
Typical runtime (test)	27 s	44 s	68 s

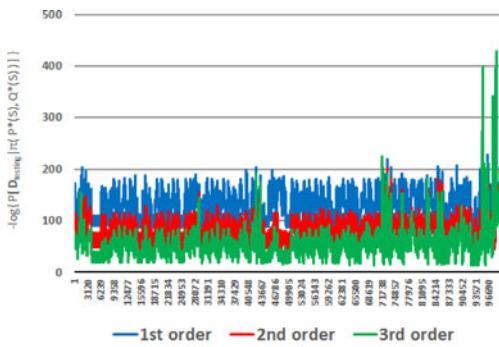
$$-\log \{ P[\mathbf{D}_{testing} | \pi(P^*(S), Q^*(S))] \} \quad (63)$$

instead of equation (58).

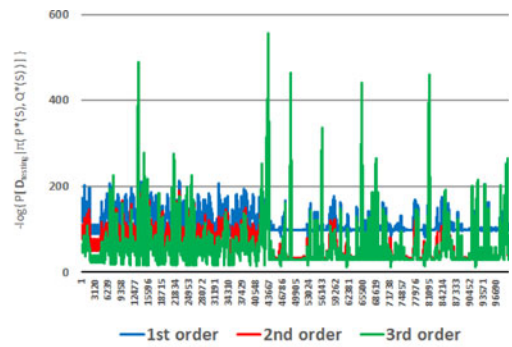
Although somewhat similar at first glance, most of these two sets of results could be easily distinguished in many aspects, e.g., where the univariate model in

Fig. 12a gives 330 to indicate anomalies, the multivariate one generates 400 in Fig. 13a, which further reduces the average one-step transition probability of 0.02 in equation (61) to

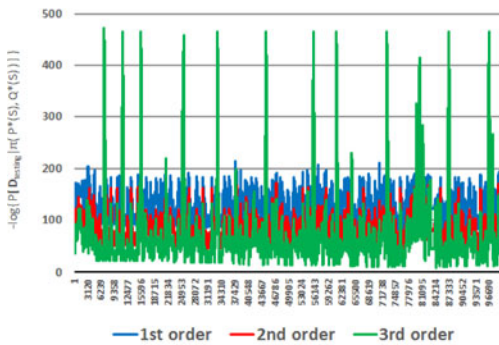
$$(10^{-400})^{1/200} = 0.01. \quad (64)$$



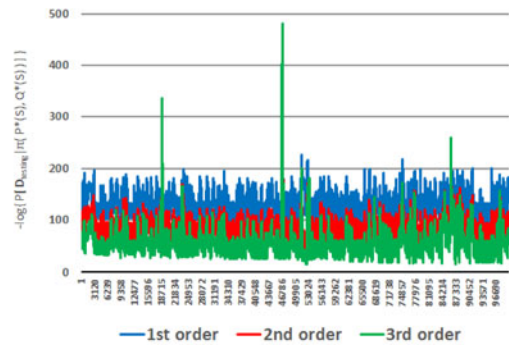
(a) Friday, Week 2 (two-dimensional)



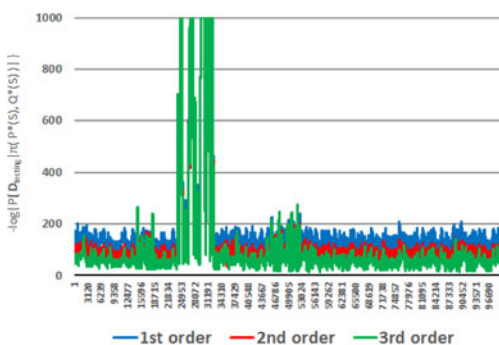
(b) Friday, Week 3 (two-dimensional)



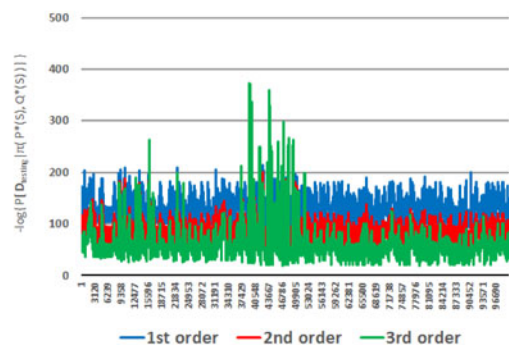
(c) Friday, Week 4 (two-dimensional)



(d) Friday, Week 5 (two-dimensional)



(e) Friday, Week 6 (two-dimensional)



(f) Friday, Week 7 (two-dimensional)

Fig. 13. Friday, Week 3 to 7 (two-dimensional).

The greater values from the multivariate model $\pi[P^*(S), Q^*(S)]$ (or smaller values in terms of average one-step transition probability) naturally lead to a more sensitive anomaly detection considering the fact that when intrusion or attack occurs, it is more likely to have an extensive scale of system calls executed with the return value of FAILURE.

Other multivariate results like the sixth Friday in Fig. 13e fail to yield a better detection, as the curves in Fig. 12e are already at its theoretical maximum of 1,000 (or 10^{-5} as average one-step transition probability) using the univariate approach.

5.3 Discussion

Essentially, the high-order Markov chain model $\pi(P^*, Q^*)$ comes with far more degrees of freedom than the first-order model $\pi(P, Q)$. Therefore, after sufficient training, $\pi(P^*, Q^*)$ is expected to extract “more” information from the training set regarding normal transition patterns, but at the cost of fitting some idiosyncratic noises as well. In the application above, the high-order model $\pi(P^*, Q^*)$ acts like a signal amplifier that increases the level of “normalness” (or “abnormalness”) to a normal (or abnormal) transition. It is this trait that results in the distinct inversion of relative positions from the multi-order models when anomalies occur, and this relation can be useful in detecting anomalies along with the conditional probabilities themselves given by different-order models.

However, it is still perfectly possible to come up with a certain training set that leads to a model with completely opposite properties, which inevitably requires a cautious selection for the training set. After all, it is what is written in the training data that defines the normal and the anomaly in any application. In our study, the unique settings of a server system helps to relax such constraints to some extent.

To summarize, using the *Training* algorithm and *Testing* algorithm, we provide the following rules of thumb:

- 1) Relations between $P[\mathbf{D}_{testing}|\pi(P, Q)]$ and $P[\mathbf{D}_{normal}|\pi(P^*, Q^*)]$ can be regarded a sign of anomaly. For normal (or common) testing data set \mathbf{D}_{normal} , while the first-order model $\pi(P, Q)$ will already assign relatively large value to $P[\mathbf{D}_{normal}|\pi(P, Q)]$, higher-order model $\pi(P^*, Q^*)$ tends to give even larger value to the probability of the corresponding higher-order transitions in space \mathbf{X}^K . On the other hand, for abnormal (or rare) testing data set $\mathbf{D}_{anomaly}$, while the first-order model $\pi(P, Q)$ will already assign relatively small (or ZERO) value to $P[\mathbf{D}_{anomaly}|\pi(P, Q)]$, higher-order model $\pi(P^*, Q^*)$ tends to give even smaller value to the probability of the corresponding higher-order transitions in space \mathbf{X}^K . Putting it altogether, we have:

$$\begin{aligned} P[\mathbf{D}_{normal} | \pi(P, Q)] \\ \leq P[\mathbf{D}_{normal} | \pi(P^*, Q^*)] \end{aligned} \quad (65)$$

$$\begin{aligned} P[\mathbf{D}_{anomaly} | \pi(P, Q)] \\ \geq P[\mathbf{D}_{anomaly} | \pi(P^*, Q^*)]. \end{aligned} \quad (66)$$

- 2) The model $\pi[P^*(S), Q^*(S)]$ that takes into account the multi-variate sequence (with return value being the second dimension for example) can do even better than the multi-order model $\pi(P^*, Q^*)$ alone, or:

$$\begin{aligned} P[\mathbf{D}_{anomaly} | \pi(P^*, Q^*)] \\ \geq P\{\mathbf{D}_{anomaly} | \pi[P^*(S), Q^*(S)]\}. \end{aligned} \quad (67)$$

6 CONCLUSIONS

In this paper, we proposed a multi-order Markov chain based anomaly detection framework. By monitoring the relative relations between results from the different-order models, we provide a new effective indicator of anomalies. In general, due to the regular and periodical behaviors of cloud server systems, if the probability of test set given the lower-order model exceeds that given the higher-order one, it is implied that unusual events might have occurred in the system and further attentions or actions would be necessary.

Besides, combining multi-dimensional inter-related sequences as a multivariate one into a single model would be another feasible approach to improve the sensitivity of detection. As shown previously, the return value series can be a useful complement to the system call series used the traditional practice.

In addition, with both time and space efficiency of the *Training* and *Testing* algorithm, this approach minimizes the possibility of becoming the source of anomalies itself and is fully capable of online (or real-time) detection. The time consumption of the training stage takes no more than 15 seconds for a training set as large as 1.6 million, and for models up to the third-order combined. To further improve efficiency, various ways such as equivalent space construction by including an artificial state, or binary representation for sparse matrix could significantly mitigate the space complexity problem.

Further efforts may involve the time-inhomogeneous Markov chain, which might be able to build transition matrix separately for any specific time period, as the time-homogeneous assumption stated here could seem too restrictive for time-sensitive systems and anomalies. Model aggregations methods may also be relevant in combining multi-order conditional probabilities into one single value for faster decision-making.

ACKNOWLEDGMENTS

This paper is sponsored in part by the Shanghai International Science and Technology Collaboration Program under Grant 13430710400, and National Research Foundation (NRF), Prime Ministers Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CREATE) program by Shanghai Jiao Tong University (SJTU) and National University of Singapore (NUS). Prof. Min Chens work was supported in part by China National Science Foundation under Grants 61300224, the International Science and Technology Collaboration Program (2014DFT10070) funded by China Ministry of Science and Technology (MOST), Hubei Provincial Key Project under grant 2013CFA051, and Program for New Century Excellent

Talents in University (NCET). The authors are also grateful to anonymous reviewers for their valuable constructive comments and suggestions. Yongxin Zhu is the corresponding author.

REFERENCES

- [1] D. E. Denning, "An intrusion-detection model," *IEEE Trans. Softw. Eng.*, vol. SE-3, no. 2, pp. 222–232, Feb. 1987.
- [2] X. D. Hoang, J. Hu, and P. Bertok, "A program-based anomaly intrusion detection scheme using multiple detection engines and fuzzy inference," *J. Netw. Comput. Appl.*, vol. 32, pp. 1219–1228, 2009.
- [3] H. T. Elshoush and I. M. Osman, "Alert correlation in collaborative intelligent intrusion detection systems-A survey," *Appl. Soft Comput.*, vol. 11, pp. 4349–4365, 2011.
- [4] P. G. Bringas and Y. K. Penya, "Next-generation misuse and anomaly prevention system," in *Proc. 10th Int. Conf. Enterprise Inf. Syst.*, 2009, vol. 19, pp. 117–129.
- [5] S. Saravanakumar, A. A. Kumar, S. Anandaraj, and S. Gowtham, "Algorithms based on artificial neural networks for intrusion detection in heavy traffic computer networks," in *Proc. Int. Conf. Adv. Inf. Technol.*, 2011, pp. 6–23.
- [6] C. V. Raman and A. Negi, "A hybrid method to intrusion detection systems using HMM," in *Proc. 2nd Int. Conf. Distrib. Comput. Internet Technol.*, 2005, pp. 389–396.
- [7] Y. Liao and V. R. Vemuri, "Use of K-nearest neighbor classifier for intrusion detection," *Comput. Security*, vol. 21, no. 5, pp. 439–448, 2002.
- [8] N. Ye, "A markov chain model of temporal behavior for anomaly detection," in *Proc. IEEE Syst. Man Cybern. Workshop Inf. Assurance Security*, West Point, NY, USA, Jun. 2000, pp. 171–174.
- [9] N. Ye, X. Li, Q. Chen, S. M. Emran, and M. Xu, "Probabilistic techniques for intrusion detection based on computer audit data," *IEEE Trans. Syst., Man, Cybern. A, Syst. Humans*, vol. 31, no. 4, pp. 266–274, Jul. 2001.
- [10] N. Ye, Y. Zhang, and C. M. Borrer, "Robustness of the Markov-chain model for cyber-attack detection," *IEEE Trans. Reliability*, vol. 53, no. 1, pp. 116–123, Mar. 2004.
- [11] R. Gwadera and M. Atallah, "Markov models for identification of significant episodes," in *Proc. 5th Int. Conf. Data Mining*, 2005, pp. 404–414.
- [12] W.-H. Ju and Y. Vardi, "A hybrid high-order markov chain model for computer intrusion detection," *J. Comput. Graphical Statist.*, vol. 10, no. 2, pp. 277–295, 2001.
- [13] C. Yin, S. Tian, and S. Mua, "High-order Markov kernels for intrusion detection," *Neurocomputing*, vol. 71, no. 16–18, pp. 3247–3252, 2008.
- [14] S. Forrest, A. Steven, Hofmeyr, A. Somayaji, and T. A. Longstaff, "A sense of self for unix processes," in *Proc. IEEE Symp. Security Privacy*, 1996, pp. 120–128.
- [15] G. Tandon and P. Chan, "Learning rules from system call arguments and sequences for anomaly detection," in *Proc. Workshop Data Mining Comput. Security, IEEE Int. Conf. Data Mining*, Nov. 19–22, 2003, pp. 20–29.
- [16] W. Wang, X. Guan, X. Zhang, and L. Yang, "Profiling program behavior for anomaly intrusion detection based on the transition and frequency property of computer audit data," *Comput. Security*, vol. 25, no. 7, pp. 539–550, 2006.
- [17] R. P. Lippmann, D. J. Fried, I. Graf, J. W. Haines, K. R. Kendall, D. McClung, D. Weber, S. E. Webster, D. Wyschogrod, R. K. Cunningham, and M. A. Zissman, "Evaluating intrusion detection systems: The 1998 DARPA off-line intrusion detection evaluation," in *Proc. DARPA Inf. Survivability Conf. Expo.*, 2000, vol. 2, pp. 12–26.
- [18] M. Stowasser, "Modelling rain risk: A multi-order Markov chain model approach," *J. Risk Finance*, vol. 13, no. 1, pp. 45–60, 2011.
- [19] Y. Song, A. D. Keromytis, and S. J. Stolfo, "Spectrogram: A mixture-of-Markov-chains model for anomaly detection in web traffic," presented at the Network and Distributed System Security Symp., San Diego, CA, USA, Feb. 2009.
- [20] W. Sha, Y. Zhu, T. Huang, M. Qiu, Y. Zhu, and Q. Zhang, "A multi-order markov chain based scheme for anomaly detection," in *Proc. IEEE 37th Annu. Comput. Softw. Appl. Conf. Workshops*, 2013, pp. 83–88.
- [21] J.-M. Koo and S.-B. Cho, "Interpreting chance for computer security by viterbi algorithm with edit distance," *New Math. Natural Comput.*, vol. 1, No. 3, pp. 421–433, 2005



Wenyaoyao Sha received the BS degree in microelectronics from Shanghai Jiao Tong University in 2012. His research interest includes machine learning, data mining, algorithm optimizing, asset pricing, and financial forecasting.



Yongxin Zhu (M'00-SM'10) received the BEng degree in EE from Hefei University of Technology, and the MEng degree in CS from Shanghai Jiao Tong University in 1991 and 1994, respectively. He received the PhD degree in CS from National University of Singapore in 2001. He is an associate professor with the School of Microelectronics (merged into School of Electronic Information and Electrical Engineering recently), Shanghai Jiao Tong University, China. His research interest is in computer architectures, embedded systems,

medical electronics, and multimedia. He has authored and coauthored more than 70 English journal and conference papers and 30 Chinese journal papers. He holds more than 10 approved Chinese patents. He is a senior member of the China Computer Federation and a senior member of the IEEE.



Min Chen (M'08-SM'09) is a professor in the School of Computer Science and Technology, Huazhong University of Science and Technology (HUST). He is the chair of IEEE Computer Society (CS) Special Technical Communities (STC) on Big Data. He was an assistant professor in School of Computer Science and Engineering at Seoul National University (SNU) from September 2009 to February 2012. He was the R&D director at Confederal Network Inc. from 2008 to 2009. He was a postdoctoral

fellow in the Department of Electrical and Computer Engineering, University of British Columbia (UBC) for three years. Before joining UBC, he was a postdoctoral fellow at SNU for one and half years. He received Best Paper Award from IEEE ICC 2012, and Best Paper Runner-up Award from QShine 2008. He has more than 180 paper publications, including 85 SCI papers. He serves as an editor or associate editor for *Information Sciences*, *Wireless Communications and Mobile Computing*, *IET Communications*, *IET Networks*, *Wiley I. J. of Security and Communication Networks*, *Journal of Internet Technology*, *KSII Trans. Internet and Information Systems*, *International Journal of Sensor Networks*. He is a managing editor for IJAACS and IJART. He is a guest editor for the *IEEE Network*, *IEEE Wireless Communications Magazine*, etc. He is the co-chair of IEEE ICC 2012-Communications Theory Symposium, and the co-chair of IEEE ICC 2013-Wireless Networks Symposium. He is the general co-chair for IEEE CIT-2012 and Mobimedia 2015. He is the general vice chair for Tridentcom 2014. He is a Keynote speaker for CyberC 2012 and Mobiquitous 2012. He is a TPC member for IEEE INFOCOM 2013 and INFOCOM 2014. His research focuses on Internet of things, big data, machine to machine communications, body area networks, e-healthcare, mobile cloud computing, ad hoc cloudlet, cloud-assisted mobile computing, ubiquitous network and services, and multimedia transmission over wireless network, etc. He is a senior member of the IEEE.



Tian Huang (M'13) received the BS degree in electronics science and technology from Shanghai Jiao Tong University in 2008. He is currently working the PhD degree in electronics science and technology at Shanghai Jiao Tong University. His research interest includes computer architecture, embedded system, multimedia signal processing, parallel processing, and system on chip. He is a member of the IEEE.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.