



**STATISTICAL LEARNING FOR ANOMALY DETECTION IN CLOUD SERVER
SYSTEMS A MULTI-ORDER MARKOV CHAIN FRAMEWORK**

KALYANAM HEMA SATYA SAI BHAGAVAN

PG Scholar, Department of Computer Science,
SVKP & Dr K S Raju Arts & Science College,
Penugonda, W.G.Dt., A.P, India
khssbhagavan1919@gmail.com

B.N.SRINIVASA GUPTA

Associate Professor in Computer Science,
SVKP & Dr K S Raju Arts & Science College,
Penugonda, W.G.Dt., A.P, India
bnsgupta@gmail.com

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

keywords : Markov processes, Training, Cloud computing, Algorithm design and analysis, Mathematical model, Equations, Servers

1. INTRODUCTION

1.1 Introduction:

An increasing number of academics 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 is 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. These approaches are usually categorized into three groups, i.e. statistical approaches, machine learning approaches and data

mining approaches: 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.

Machine learning based approaches tend to reduce the supervision 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.

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 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 the context of this work, we focus on approaches in categories of statistical and machine learning based approaches. 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. Ju and Vardi introduce 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 or support vector machine. In addition to statistical patterns, as originally proposed, 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 and 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. 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 modeling. 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, 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 behavior is insufficient in that it does not take into account the resulting status of execution. In a few more recent works, return values of system calls were taken into consideration to detect or interpret the anomalies

1.2 Purpose:

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 academics 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

1.3 Scope:

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.

1.4 Motivation:

We focus on approaches in categories of statistical and machine learning based approaches. In 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. Ju and Vardi 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 i or support vector machine.

1.5 Overview:

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.

2. LITERATURE SURVEY

A Program based anomaly intrusion detection scheme using multiple detection engines and fuzzy inference

This work proposed a hybrid anomaly intrusion detection scheme using program system calls is proposed. In this scheme, a hidden Markov model (HMM) detection engine and a normal database detection engine have been combined to utilize their respective advantages. A fuzzy-based inference mechanism is used to infer a

soft boundary between anomalous and normal behavior, which is otherwise very difficult to determine when they overlap or are very close. To address the challenging issue of high cost in HMM training, an incremental HMM training with optimal initialization of HMM parameters is suggested.

Alert correlation in collaborative intelligent intrusion detection systems—A survey

As complete prevention of computer attacks is not possible, intrusion detection systems (IDSs) play a very important role in minimizing the damage caused by different computer attacks. There are two intrusion detection methods: namely misuse- and anomaly-based. A collaborative, intelligent intrusion detection system (CIIDS) is proposed to include both methods, since it is concluded from recent research that the performance of an individual detection engine is rarely satisfactory. In particular, two main challenges in current collaborative intrusion detection systems (CIDSs) research are highlighted and reviewed: CIDSs system architectures and alert correlation algorithms. Different CIDSs system, architectures are explained and compared. The use of CIDSs together with



other multiple security systems raises certain issues and challenges in, alert correlation. Several different techniques for alert correlation are discussed. The focus will be on correlation of CIIDS alerts. Computational, Intelligence approaches, together with their applications on IDSs, are reviewed. Methods in soft computing collectively provide understandable and autonomous solutions to IDS problems.

Interpreting chance for computer security by viterbi algorithm with edit distance

This work addresses the importance of chance discovery in computer security. There are various methods to discover chances in computer usage, but they have such drawbacks as discovering only anomalies, not interpreting anomalies in conventional approach in the community of computer security. This work focuses on a role of interpreting the type of anomalies by analyzing the state sequences using Viterbi algorithm and evaluating the distance between the standard model of anomaly type and the state sequence of discovered anomalies. Because the state sequences are not always extracted consistently due to environmental factor, the edit distance is utilized to measure the distance effectively.

A multi-order markov chain based scheme for anomaly detection

This work presents a feasible multi-order Markov chain based scheme for anomaly detection in server systems. In our approach, both the high-order Markov chain and multivariate time series are taken into account, along with the detailed design of training and testing algorithms. To evaluate its effectiveness, the Defense Advanced Research Projects Agency (DARPA) Intrusion Detection Evaluation Data Set is used as stimuli to our model, by which system calls and the corresponding return values form a two-dimensional input set. The calculation result shows that this approach is able to produce several effective indicators of anomalies. In addition to the absolute values given by an individual single-order model, we also notice a novelty unprecedented before, i.e., the changes in ranking positions of outputs from different-order ones also correlate closely with abnormal behaviours.

Spectrogram: A mixture-of-Markov-chains model for anomaly detection in web traffic

This work proposed a Spectrogram, machine learning based statistical anomaly



detection (AD) sensor for defense against web-layer code-injection attacks. These attacks include PHP file inclusion, SQL-injection and cross-site-scripting; memory-layer exploits such as buffer overflows are addressed as well. Statistical AD sensors offer the advantage of being driven by the data that is being protected and not by malware samples captured in the wild. While models using higher order statistics can often improve accuracy, trade-offs with false-positive rates and model efficiency remain a limiting usability factor. This work presents a new model and sensor framework that offers a favorable balance under this constraint and demonstrates improvement over some existing approaches. Spectrogram is a network situated sensor that dynamically assembles packets to reconstruct content flows and learns to recognize legitimate web-layer script input. This work describes an efficient model for this task in the form of a mixture of Markov-chains and derives the corresponding training algorithm.

Modelling rain risk: A multi-order Markov chain model approach

The purpose of this work is to investigate the best frequency description of a chain dependent Markov process for the daily simulation of precipitation. The

influence of the order of the Markov chain model to simulate daily precipitation occurrence is evaluated. A mixed-order model is constructed and compared to a simple first-order model to evaluate the importance of the model order for the pricing of a rainfall index put option. Design/methodology/approach—For the first time a mixed-order Markov chain model is presented where the monthly varying order was chosen based on a Bayesian information criteria analysis of rainfall data for one weather station in the US. The outcome of this model is compared to simpler Markov models and to burn analysis results. Findings—The comparison indicate that there is only a slightly better representation of the rain statistics in the theoretically best mixed-order Markov chain model compared to a simpler first-order model. Clear differences between the daily simulation and the burn method are found when pricing a put option on a rainfall index. All daily simulation models underestimate the volatility of the monthly rainfall amount especially in the summer months. Research limitations/implications—To assesses the robustness and any geographical dependence of the bias in the volatility a systematic analysis could be applied to more weather stations across the US in further studies.



Practical implications—The bias in the volatility has significant influence on the price of the put option considered here and limits the use of such a model for risk analyses, e.g. for an extreme event cover. Originality/value—For the first time a multi-order Markov chain model is applied to price a precipitation derivative.

Evaluating intrusion detection systems: The 1998 DARPA off-line intrusion detection evaluation

An intrusion detection evaluation test bed was developed which generated normal traffic similar to that on a government site containing 100's of users on 1000's of hosts. More than 300 instances of 38 different automated attacks were launched against victim UNIX hosts in seven weeks of training data and two weeks of test data. Six research groups participated in a blind evaluation and results were analyzed for probe, denial of -service (DoS), remote-to-local (R2L), and user to root (U2R) attacks. The best systems detected old attacks included in the training data, at moderate detection rates ranging from 63% to 93% at a false alarm rate of 10 false alarms per day. Detection rates were much worse for new and novel R2L and DoS attacks included only in the test data. The best systems failed

to detect roughly half these new attacks which included damaging access to root-level privileges by remote users.

Profiling program behavior for anomaly intrusion detection based on the transition and frequency property of computer audit data

Intrusion detection is an important technique in the defense-in-depth network security framework. In recent years, it has been a widely studied topic in computer network security. In this paper, we present two methods, namely, the Hidden Markov Models (HMM) method and the Self Organizing Maps (SOM) method, to profile normal program behavior for anomaly intrusion detection based on computer audit data. The HMM method utilizes the transition property of events while SOM method relies on the frequency property of events. Two data sets, CERT synthetic Sendmail system call data collected in the University of New Mexico (UNM) and Live FTP system call data collected in the CNSIS lab of Xi'an Jiaotong University, were used to assess the two methods. Testing results show that the HMM method using the transition property of events produces good detection performance while high computational expense is required both for



training and detection. The HMM method is better than other two methods reported previously in terms of detection accuracy for the same data set. The SOM method considering the frequency property of events, on the other hand, is suitable for real-time intrusion detection because of its capability of processing a large amount of data with low computational overhead

3. EXISTING SYSTEM

In existing, anomaly or intrusion detection was implemented in great variety of approaches. These approaches are usually categorized into three groups, i.e. statistical approaches, machine learning approaches and data mining approaches: 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. Machine learning based approaches

tend to reduce the supervision 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. 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 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.

3.1 Disadvantages:

Anomaly Detection accuracy is too low.

4. PROPOSED SYSTEM

This system proposed 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.

4.1 Advantages:

Anomaly Detection accuracy is high.

5. ARCHITECTURE

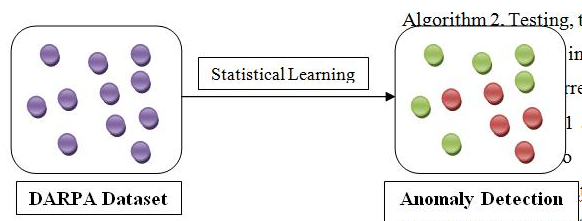


Fig 5.1 System Design

6. IMPLEMENTATION

6.1 Load Dataset:

This module loads DARPA Intrusion Detection Dataset. This dataset has n number of training sequences. Followed by, this proposed work applies training and testing for anomaly detection.

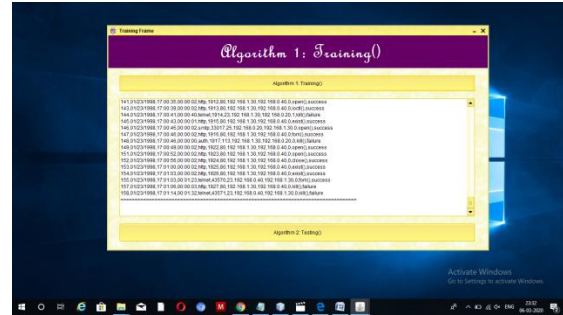
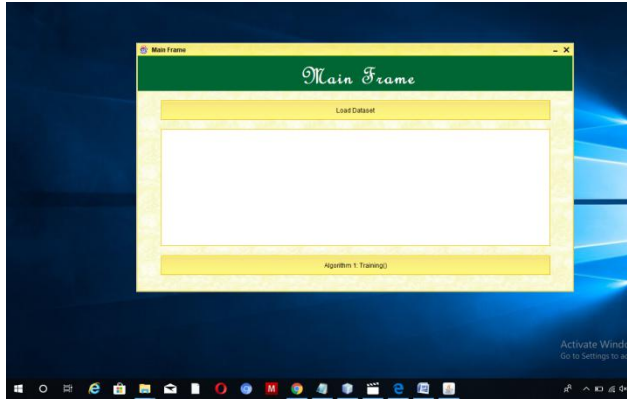
6.2 Training

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.

6.3 Testing

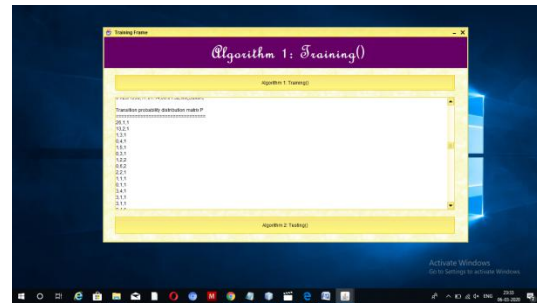
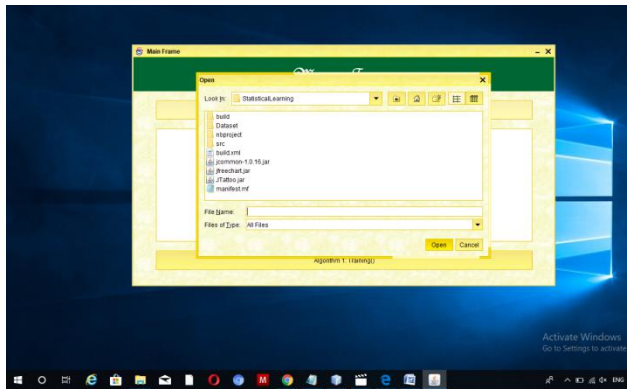
Algorithm 2, Testing, takes the two matrices and the test set in and calculates its probability of occurrence given the model trained. Algorithm 1 and Algorithm 2 also rely on two functions i.e. Increase Transition Matrix and Get Transition Matrix to calculate transition probabilities and retrieve values in these matrices.

7. OUTPUT RESULTS



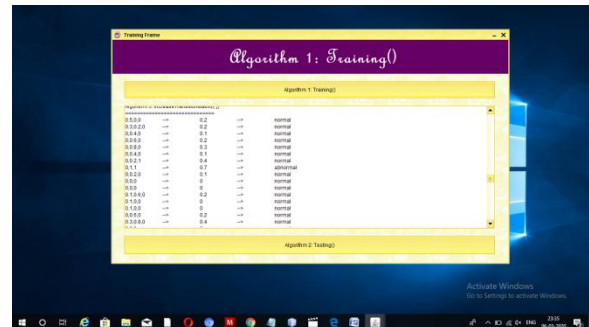
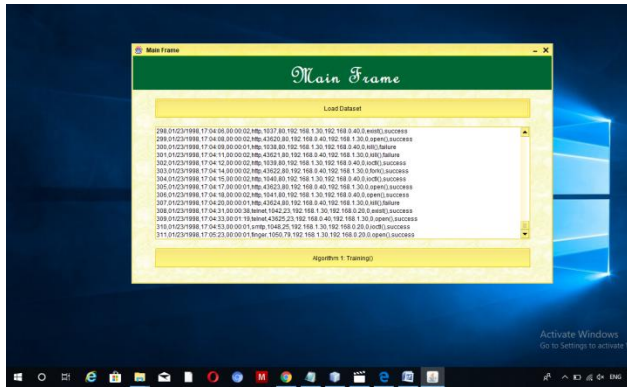
Training Algorithm

Load Data set



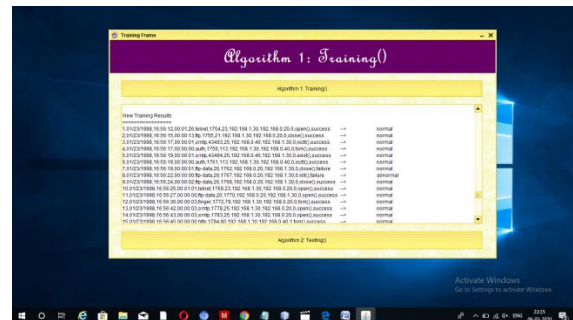
Transition probability distribution matrix p

Open Data set

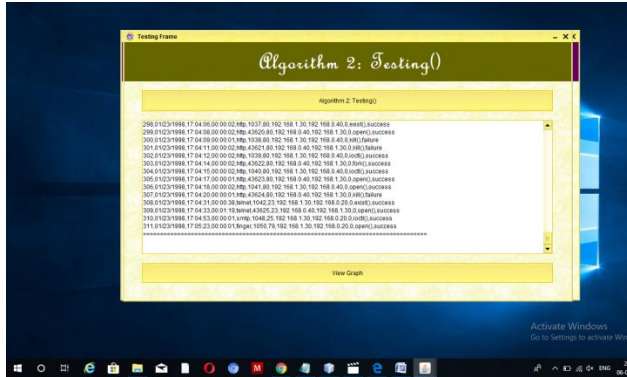


Algorithm -3

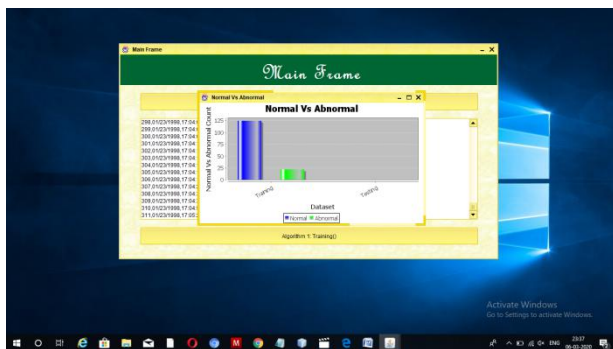
Show data set screen



View Training results.



Testing Algorithm



View Graph.

8.CONCLUSION

This work 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.

9. FUTURE ENHANCEMENT

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.

10 . BIBLIOGRAPHY

[1] 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.

[2] 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.

[3] 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

[4] 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.

[5] 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.

[6] M. Stowasser, “Modelling rain risk: A multi-order Markov chain model approach,” *J. Risk Finance*, vol. 13, no. 1, pp. 45–60, 2011.

[7] R. P. Lippmann, D. J. Fried, I. Graf, J. W. Haines, K. R. Kendall, D. McClung, D.

Weber, S. E. Webster, D. Wyszogrod, 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.

[8] 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.



ABOUT AUTHORS:



Kalyanam Hema Satya Sai Bhagavan is currently pursuing MCA in SVKP & Dr K S Raju Arts & Science College, affiliated to Adikavi Nannaya University, Rajamahendravaram. His research interests include Operation Research, Design and Analysis of Algorithm and Big Data Analytics.



B.N.Srinivasa Gupta is working as Associate Professor in SVKP & Dr K S Raju Arts & Science College, Penugonda , West Godavari District, A.P. He received Master's Degree in Computer Applications from Andhra University. His research interests include Operational Research, Probability and Statistics, Design and Analysis of Algorithm, Big Data Analytics.

