



IJIRCCCE

e-ISSN: 2320-9801 | p-ISSN: 2320-9798



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

Volume 10, Issue 6, June 2022

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

Impact Factor: 8.165



9940 572 462



6381 907 438



ijircce@gmail.com



www.ijircce.com

Modelling and Predicting Cyber Hacking Data Breaches

A. Pranav Sai¹, B. Joshi Naveen², G. Chandra Kiran Reddy³, Md. Shakeel Ahmed⁴

UG Student, Dept. of I.T., Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India ^{1,2,3}

Associate Professor, Dept. of I.T., Vasireddy Venkatadri Institute of Technology, Guntur, Andhra Pradesh, India ⁴

ABSTRACT: Examining cyber incident data sets is an important way to have a better picture of how the threat environment is evolving. This is a very new research issue, thus there are still a lot more studies to be done. A statistical analysis of a breach incidence data set spanning 12 years (2005-2017) of cyber hacking operations, including malware attacks, is presented in this research. We show that, contrary to previous research, both hacking incidence inter-arrival periods and breach sizes should be described by stochastic processes rather than distributions due to autocorrelations. Then, to accommodate the inter-arrival periods and the breach size, we suggest specific stochastic process models. We also show that these models can estimate the time between arrivals and the amount of the breach. We do both qualitative and quantitative trend studies on the data set in order to have a better understanding of the progression of hacking breach incidences. We derive a number of cyber security conclusions, including the fact that while the threat of cyber hacks is increasing in terms of frequency, it is not increasing in terms of the amount of their damage.

KEYWORDS: Cyber incident data sets, Statistical Analysis, Malware Attacks, Breach Size.

I. INTRODUCTION

Data breaches are one of the most devastating cyber incidents. The Privacy Rights Clearinghouse [1] reports 7,730 data breaches between 2005 and 2017, accounting for 9,919,228,821 breached records. The Identity Theft Resource Center and Cyber Scout [2] reports 1,093 data breach incidents in 2016, which is 40% higher than the 780 data breach incidents in 2015. The United States Office of Personnel Management (OPM) [3] reports that the personnel information of 4.2 million current and former Federal government employees and the background investigation records of current, former, and prospective federal employees and contractors (including 21.5 million Social Security Numbers) were stolen in 2015.

The monetary price incurred by data breaches is also substantial. IBM [4] reports that in year 2016, the global average cost for each lost or stolen record containing sensitive or confidential information was \$158. Net Diligence [5] reports that in year 2016, the median number of breached records was 1,339, the median per-record cost was \$39.82, the average breach cost was \$665,000, and the median breach cost was \$60,000. While technological solutions can harden cyber systems against attacks, data breaches continue to be a big problem. This motivates us to characterize the evolution of data breach incidents. This not only will deep our understanding of data breaches, but also shed light on other approaches for mitigating the damage, such as insurance. Many believe that insurance will be useful, but the development of accurate cyber risk metrics to guide the assignment of insurance rates is beyond the reach of the current understanding of data breaches (e.g., the lack of modeling approaches).

II. PROPOSED SYSTEM

In this paper, we make the following three contributions. First, we show that both the hacking breach incident interarrival times (reflecting incident frequency) and breach sizes should be modeled by stochastic processes, rather than by distributions. We find that a particular point process can adequately describe the evolution of the hacking breach incidents inter-arrival times and that a particular ARMA-GARCH model can adequately describe the evolution of the hacking breach sizes, where ARMA is acronym for “Autoregressive and Moving Average” and GARCH is acronym for “Generalized Autoregressive Conditional Heteroskedasticity. “We show that these stochastic process models can predict the inter-arrival times and the breach sizes. To the best of our knowledge, this is the first paper showing that stochastic processes, rather than distributions, should be used to model these cyber threat factors. Second, we discover a positive dependence between the incidents inter-arrival times and the breach sizes, and show that this dependence can be adequately described by a particular copula. We also show that when predicting inter-arrival times

and breach sizes, it is necessary to consider the dependence; otherwise, the prediction results are not accurate. To the best of our knowledge, this is the first work showing the existence of this dependence and the consequence of ignoring it. Third, we conduct both qualitative and quantitative trend analyses of the cyber hacking breach incidents. We find that the situation is indeed getting worse in terms of the incidents inter-arrival time because hacking breach incidents become more and more frequent, but the situation is stabilizing in terms of the incident breach size, indicating that the damage of individual hacking breach incidents will not get much worse. We hope the present study will inspire more investigations, which can offer deep insights into alternate risk mitigation approaches. Such insights are useful to insurance companies, government agencies, and regulators because they need to deeply understand the nature of data breach risks.

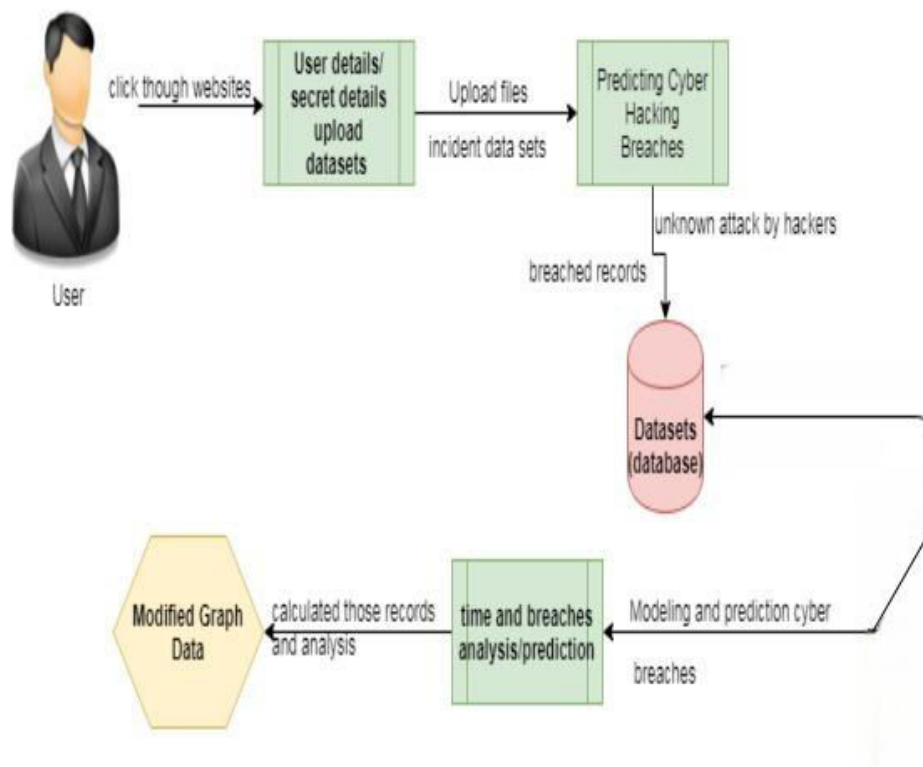
III. METHODOLOGY

SUPPORT VECTOR MACHINE:

Support Vector Machine (SVM) is a supervised machine learning algorithm which can be used for both classification or regression challenges. A support vector machine (SVM) is a supervised machine learning model that uses classification algorithms for two-group classification problems. So you're working on a text classification problem. You're refining your training data, and maybe you've even tried stuff out using Naive Bayes. However, it is mostly used in classification problems. In this algorithm, we plot each data item as a point in n-dimensional space (where n is number of features you have) with the value of each feature being the value of a particular coordinate. Then, we perform classification by finding the hyperplane that differentiate the two classes very well.

The SVM algorithm is implemented in practice using a kernel. The learning of the hyper plane in linear SVM is done by transforming the problem using some linear algebra, which is out of the scope of this introduction to SVM. A powerful insight is that the linear SVM can be rephrased using the inner product of any two given observations, rather than the observations themselves. The inner product between two vectors is the sum of the multiplication of each pair of input values.

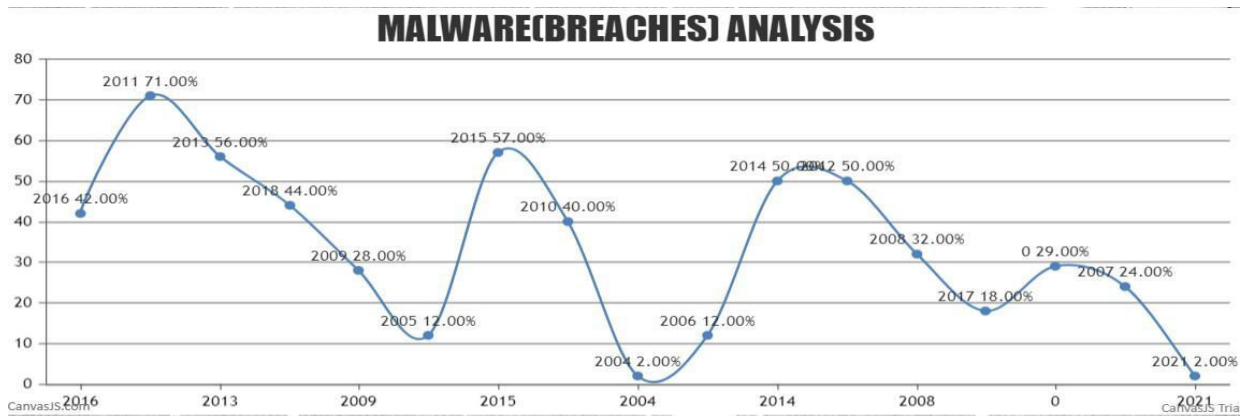
IV. SYSTEM ARCHITECTURE



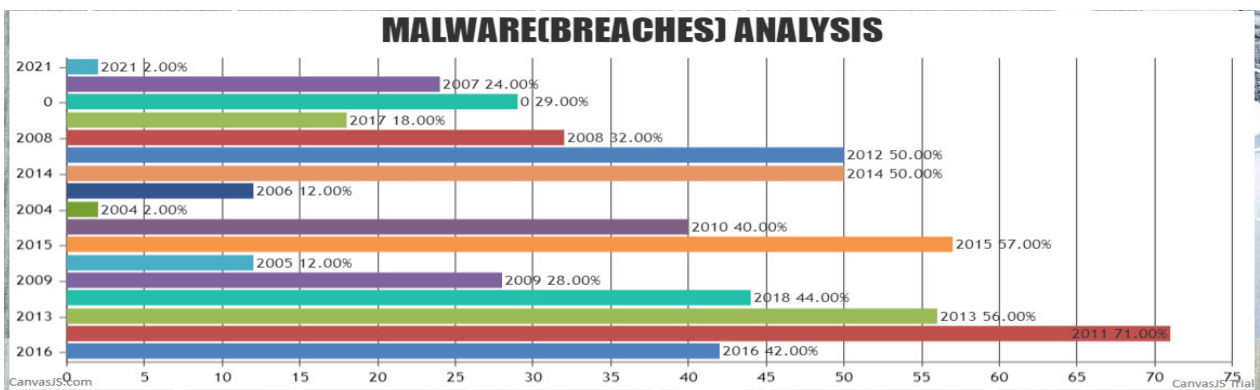
UNMALWARE DATA



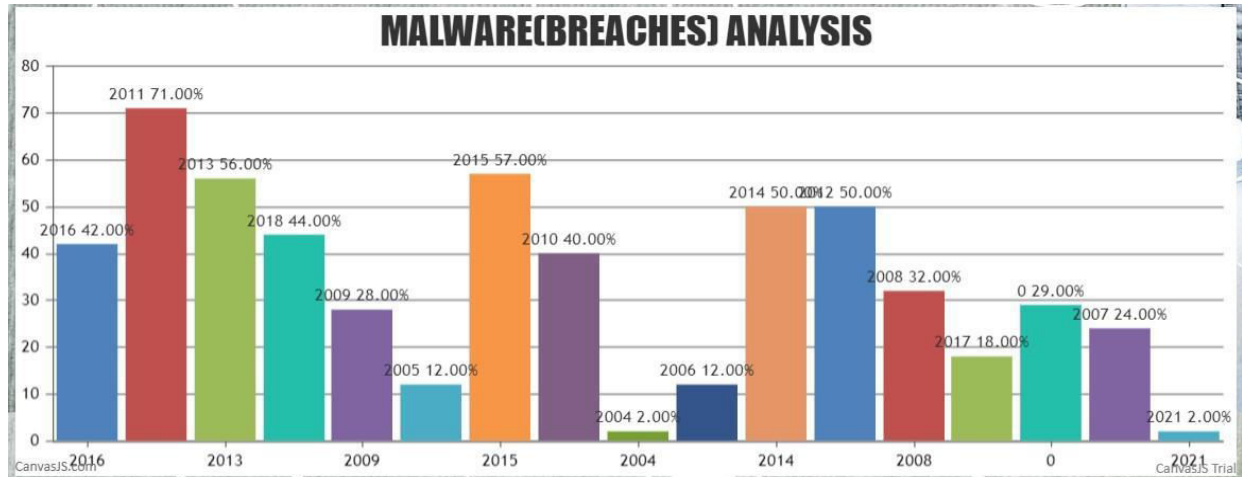
SPLINE CHART ANALYSIS



BAR CHART ANALYSIS



COLOUMN CHART ANALYSIS



VI. CONCLUSION

We analyzed a hacking breach dataset from the points of view of the incidents inter-arrival time and the breach size, and showed that they both should be modeled by stochastic processes rather than distributions. The statistical models developed in this paper show satisfactory fitting and prediction accuracies. In particular, we propose using a copula-based approach to predict the joint probability that an incident with a certain magnitude of breach size will occur during a future period of time. Statistical tests show that the methodologies proposed in this paper are better than those which are presented in the literature, because the latter ignored both the temporal correlations and the dependence between the incidents inter-arrival times and the breach sizes. We conducted qualitative and quantitative analyses to draw further insights. We drew a set of cybersecurity insights, including that the threat of cyber hacking breach incidents is indeed getting worse in terms of their frequency, but not the magnitude of their damage. The methodology presented in this paper can be adopted or adapted to analyze datasets of a similar nature.

VII. FUTURE ENHANCEMENTS

There are many open problems that are left for future research. For example, it is both interesting and challenging to investigate how to predict the extremely large values and how to deal with missing data (i.e., breach incidents that are not reported). It is also worthwhile to estimate the exact occurring times of breach incidents. Finally, more research needs to be conducted towards understanding the predictability of breach incidents (i.e., the upper bound of prediction accuracy).

REFERENCES

1. P. R. Clearinghouse. Privacy Rights Clearinghouse’s Chronology of Data Breaches. Accessed: Nov. 2017. [Online]. Available: <https://www.privacyrights.org/data-breaches>
2. ITR Center. Data Breaches Increase 40 Percent in 2016, Finds New Report From Identity Theft Resource Center and Cybersoul. Accessed: Nov. 2017. [Online]. Available: <http://www.idtheftcenter.org/2016databreaches.html>
3. C. R. Center. Cybersecurity Incidents. Accessed: Nov. 2017. [Online]. Available: <https://www.opm.gov/cybersecurity/cybersecurity-incidents>
4. IBM Security. Accessed: Nov. 2017. [Online]. Available: <https://www.ibm.com/security/data-breach/index.html>
5. Net Diligence. The 2016 Cyber Claims Study. Accessed: Nov. 2017. [Online]. Available: https://netdiligence.com/wp-content/uploads/2016/10/P02_NetDiligence-2016-Cyber-Claims-Study-ONLINE.pdf
6. M. Eling and W. Schnell, “What do we know about cyber risk and cyber risk insurance?” J. Risk Finance, vol. 17, no. 5, pp. 474–491, 2016.



INNO  SPACE
SJIF Scientific Journal Impact Factor

Impact Factor: 8.165

 **doi**[®]
cross **ref**

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

 9940 572 462  6381 907 438  ijircce@gmail.com



www.ijircce.com

Scan to save the contact details