



MPCB : MODELING AND PREDICTING CYBERHACKING BREACHES

Ms. SIDDILA HARSHINI PRIYA¹, Ms. INALA LAHARI PRIYA², Ms. JUTRU BHAGYA NAGA AKHILA³,
Ms. BATHULA VENKATA NAGA SRIVALLI⁴, Ms. CHODISETTI DEEPIKA

1. BTech, Vijaya Institute of Technology For Women, Enikepadu, Vijayawada, Andhra Pradesh, India.
Email : priyaharshini21@gmail.com
2. BTech, Vijaya Institute of Technology For Women, Enikepadu, Vijayawada, Andhra Pradesh, India.
3. BTech, Vijaya Institute of Technology For Women, Enikepadu, Vijayawada, Andhra Pradesh, India.
4. BTech, Vijaya Institute of Technology For Women, Enikepadu, Vijayawada, Andhra Pradesh, India.
5. Assistant Professor, Computer Science and Engineering, Vijaya Institute of Technology For Women, Enikepadu, Vijayawada, Andhra Pradesh, India. Email : deepika3062@gmail.com

ABSTRACT

As part of our study, we will analyse data sets related to cyber incidents in order to better understand how the threat environment has changed over time. There still need to be many studies conducted on this subject because it is so recent. In this research, we present the results of a statistical analysis of a breach incidence data set that covers 12 years (2005–2017) of malware-infected cyber hacking activities. Contrary to findings in the literature, we demonstrate that because used to model them rather than distributions. Then, in order to accommodate the inter-arrival durations and breach sizes, respectively, we provide specific stochastic process models. We also demonstrate that these models are capable of forecasting breach sizes and inter-arrival intervals. We perform both qualitative and quantitative trend studies on the data set to gain deeper insights into the development of hacking breach occurrences. We derive a number of cybersecurity insights, such as the fact that while the threat of cyber intrusions is increasing in frequency, the size of the harm they cause is not.

1 INTRODUCTION

Cyber hacking is an attempt to gain access to a computer's internal network or computing system. Unauthorised access to the network security system is being used for a few illegal purposes. Sensitive, confidential, or otherwise protected data has been accessed without authorization in the data breaches. An assault carried out by cybercriminals utilising one or more computers or networks is known as a cyber attack. A confirmed instance in which sensitive, confidentially protected data has been accessed or revealed without authorization is referred to as a data breach.

A data breach occurs when safe or private/confidential information is accidentally or knowingly disclosed to an unreliable environment. Unintentional information disclosure, data leak, information leakage, and data spill are further words for this phenomena. The incidents range from coordinated hacking operations by "black hats," or people involved in organised crime, political activism, or national governments, to the reckless disposal of outdated computer hardware or data storage media and unhackable sources.

Trade secrets and personal health information may be compromised in data breaches. People who violate privacy regulations run the risk of being humiliated, losing their job or business opportunity, losing their physical safety, and becoming the victim of identity theft. When a cybercriminal successfully accesses a data source and steals private data, there is a data breach. Data breaches are becoming more and more regular, and some of the most recent ones have been the largest ever recorded. This can be done physically by physically accessing a computer or network to take local files, or it can be done remotely by getting over network protection.

One of the most severe cyber events is a data leak. In total, 9,919,228,821 records were compromised between 2005 and 2017, according to the Privacy Rights Clearing house, which cites 7,730 data breaches. 1,093 data breach events were reported in 2016, 40% more than the 780 data breach occurrences in 2015, according to the Identity Theft Resource Centre and Cyber Scout. The United States Office of Personnel Management (OPM) reports that in 2015, the background investigation records of current, former, and potential federal employees and contractors, as well as the personnel information of 4.2 million current and former employees of the federal government, including 21.5 million Social Security



numbers, were stolen. The cost of data breaches in terms of money is likewise high. According to IBM, the average cost per lost or stolen record holding sensitive or confidential information was \$158 globally in 2016. According to NetDiligence, the average number of records that were compromised in 2016 was 1,339, the average cost per compromised record was \$39.82, the average breach cost was \$665,000, and the median breach cost was \$60,000. Despite the fact that technological advancements can fortify cyber systems against attacks, data breaches continue to be a major issue. This encourages us to describe how data breach occurrences have changed over time. Our understanding of data breaches will be improved as a result, and other damage-mitigation strategies like insurance will also become clearer.

Although many people think insurance will be helpful, the creation of precise cyber risk measures to direct the assignment of insurance prices is beyond the capabilities of the current understanding of data breaches (for example, the dearth of modelling methodologies). Researchers have recently begun to model data breach instances. Between 2000 and 2008, Maillart and Sornette looked at the statistical characteristics of lost personal identities in the US. They discovered that the number of breach occurrences significantly rises from 2000 and July 2006 but then stabilises. 2,253 breach instances from 2005 to 2015 were included in the dataset Edwards et al. examined. They discovered that neither the volume of data breaches nor their frequency have grown over time. A dataset that is combined from and correlates to organisational breach incidences between the years 2000 and 2015 was analysed by Wheatley et al. They discovered that while the frequency of major breach occurrences—those involving breaches of more than 50,000 records—occurring to US organisations is not time-dependent, it does show an increasing trend for large breach incidents affecting non-US firms.

The following unanswered concerns served as the impetus for the current study: Are data breaches brought on by cyberattacks increasing, decreasing, or stabilising? A moral response to this query will provide us with a clear understanding of the current state of cyberthreats. Previous research did not provide an answer to this question. The dataset analysed in is more recent, but it contains two types of incidents: negligent breaches (i.e., incidents caused by lost, discarded, stolen devices and other reasons) and malicious breaching. Specifically, the dataset analysed in only covered the time period from 2000 to 2008 and did not always include the breach incidents that are caused by cyber attacks. We did not include negligent breaches in the current study since they are more likely to be caused by human mistake than cyberattacks. The malicious breaches investigated in this study fall into four sub-categories: hacking (including malware), insider, payment card fraud, and unknown. Because of this, this study will only focus on the hacking sub-category (hacked breach dataset), even though the other three sub-categories are also intriguing and should be examined separately.

Data breach is defined as "a security violation in which sensitive, protected, or confidential data is copied, transmitted, viewed, stolen, or used by an individual unauthorised to do so." [1] Financial information, including credit card and debit card numbers, bank account information, personally identifiable information (PII), trade secrets of organisations, and intellectual property, may be compromised in data breaches. The majority of data breaches involve sensitive files, papers, and other unstructured material that has been overexposed and left insecure. [2] Data breaches can be very expensive for enterprises, both in terms of direct costs (cleanup, inquiries, etc.) and indirect costs (damages to victims' reputations, provision of cyber security, etc.).

A total of 227,052,199 individual records containing sensitive personal information were involved in security breaches in the United States between January 2005 and May 2008, according to the nonprofit consumer organisation Privacy Rights Clearinghouse. This number excludes incidents where sensitive data was allegedly not actually exposed. [3] A number of jurisdictions have enacted legislation requiring companies that experience data breaches to notify customers and take other actions to address any potential harms. Any incident that allows unauthorised access to computer data, applications, networks, or devices is referred to as a security breach. As a result, information can be accessed without permission. Usually, it happens when a burglar is able to get past security measures. There is a difference between a security breach and a data breach technically.

A data breach is when a cybercriminal escapes with information, whereas a security breach is essentially a break-in. Consider a burglar; the security breach would be when he scales the window, and the data breach would be when he stole your laptop or purse. Information that is confidential has great value. On the dark web, it is frequently sold; for instance, names and credit card numbers can be purchased and used for fraud or identity theft. It is hardly unexpected that security breaches can result in significant financial losses for businesses. For large organisations, the cost is typically close to \$4



million. It's crucial to differentiate between the definitions of a security incident and a security breach. A malware infestation could be involved in an incident, DDOS attack or a worker leaving a laptop in a cab, but as long as they don't lead to data loss or network access, they are not considered security breaches.

2. LITERATURE SURVEY AND RELATED WORK

In this study, we make predictions about cyber attacks that involve hacking. Data breaches are a continual threat to people's personal and financial security, and they are expensive for the companies that store big amounts of personal data to handle residual IT security risks. However, it is still unclear whether premiums are accurate. As a result, experts in the field and academics have pushed for more reliable and creative pricing schemes for cyber-insurance. By creating a cyber-insurance model in 2011 utilising the newly developed copula methodology, the research bridges this significant gap in the literature. We focus in analysing the attack traffic flows' macroscopic qualities, and in 2015, we distinguish between two key patterns with distinct spatiotemporal properties: deterministic and stochastic. This methodology encourages the use of gray-box models, which take into account the information's statistical features and phenomena. Although our prediction study is based on particular cyber attack data, our methodology may frequently be used to research any data of this kind related to cyber attacks. Incidents of data breaches are on the rise and have had serious financial and legal repercussions for the organisations involved. To determine the elements that could raise or lower the contextual risk of a data breach in 2015, we use the opportunity theory of crime, the institutional anomie theory, and institutional theory. Many people's private information has been made public due to data breaches. According to certain reports, the scope and frequency of knowledge breaches have increased alarmingly. institutions all across the world are being encouraged to deal with what seems to be a worsening scenario in 2016. Cyberattacks have developed into a burden that is endangering the economy, personal privacy, and even national security. Before we can effectively handle the issue, we must have a thorough grasp of cyber threats in 2017 from all angles. The modelling of cybersecurity threats is a crucial yet difficult issue. We introduce the study of modelling multivariate cybersecurity threats in this work. We create the first statistical method, which is based on a Copula-GARCH model that use vine copulas to simulate the multivariate dependence shown by actual data from cyberattacks in 2018. To estimate the inter-arrival time and breach sizes at this moment, we are currently employing a stochastic process model.

3 Implementation Study

Modules:

INSERT DATA

Administrators and authorised users can both upload data resources to databases. To ensure the privacy of the data that is not disclosed without the user's awareness, the data can be uploaded with a key. Based on the information they have provided to the administrator, each user is authorised. Only authorised users are permitted to log in and upload or request files from the system.

ACCESS INFORMATION

The administrators of the database can grant users access to its data. Admin is the sole person with the authority to process accessing information, approve or disapprove users based on their information, and handle uploaded data.

PERMISSIONS PER USER

Data from any resource may be accessed with the administrator's permission alone. Users are given the



opportunity by the administrator to disclose their data prior to data access and to confirm the information they have provided. Users are blocked appropriately if they attempt to access the data in the wrong way. If a user asks to be unblocked, the administrator will unblock them depending on the user's requests and prior behaviour.

Analysing data

Graphs are used to help in data analysis. In order to obtain the best analysis and prediction of the dataset and specified data policies, the collected data are applied to a graph. Through this visual analysis, the dataset can be examined to gain a deeper understanding of the data's specifics.

4 PROPOSED WORK

We provide the following three contributions in this work. First, we demonstrate that stochastic processes, rather than distributions, should be used to describe both the hacker breach incident interarrival times (indicating incidence frequency) and breach sizes. We discover that a specific point process can adequately describe the evolution of hacking breach incidents' inter-arrival times, and that a specific ARIMA model—the abbreviation for "Auto Regressive Integrated Moving Average"—can adequately describe the evolution of hacking breach sizes.

We demonstrate that the inter-arrival periods and breach sizes can be predicted using these stochastic process models. To the best of our knowledge, this is the first paper to demonstrate the need to represent these cyber threat elements using stochastic processes rather than distributions. The second finding is that there is a positive correlation between the breach sizes and the events' inter-arrival intervals. We demonstrate that a certain copula may appropriately capture this correlation. We also demonstrate that the reliance must be taken into account in order to accurately anticipate inter-arrival periods and breach sizes. We believe that this is the first piece of work to demonstrate both the existence of this dependence and the effects of ignoring it.

Third, we undertake trend studies of the cyber hacking breach episodes using both qualitative and quantitative data. Because hacking breach incidents are occurring more frequently, we discover that the situation is indeed getting worse in terms of incidents between arrivals, but that it is improving in terms of incident breach size, suggesting that the harm from individual hacking breach incidents won't get significantly worse. We anticipate that the current study will spark further research that can provide in-depth understanding of other risk mitigation strategies. Because they must thoroughly comprehend the nature of data breach threats, insurance companies, governmental organisations, and regulators can benefit from these insights.

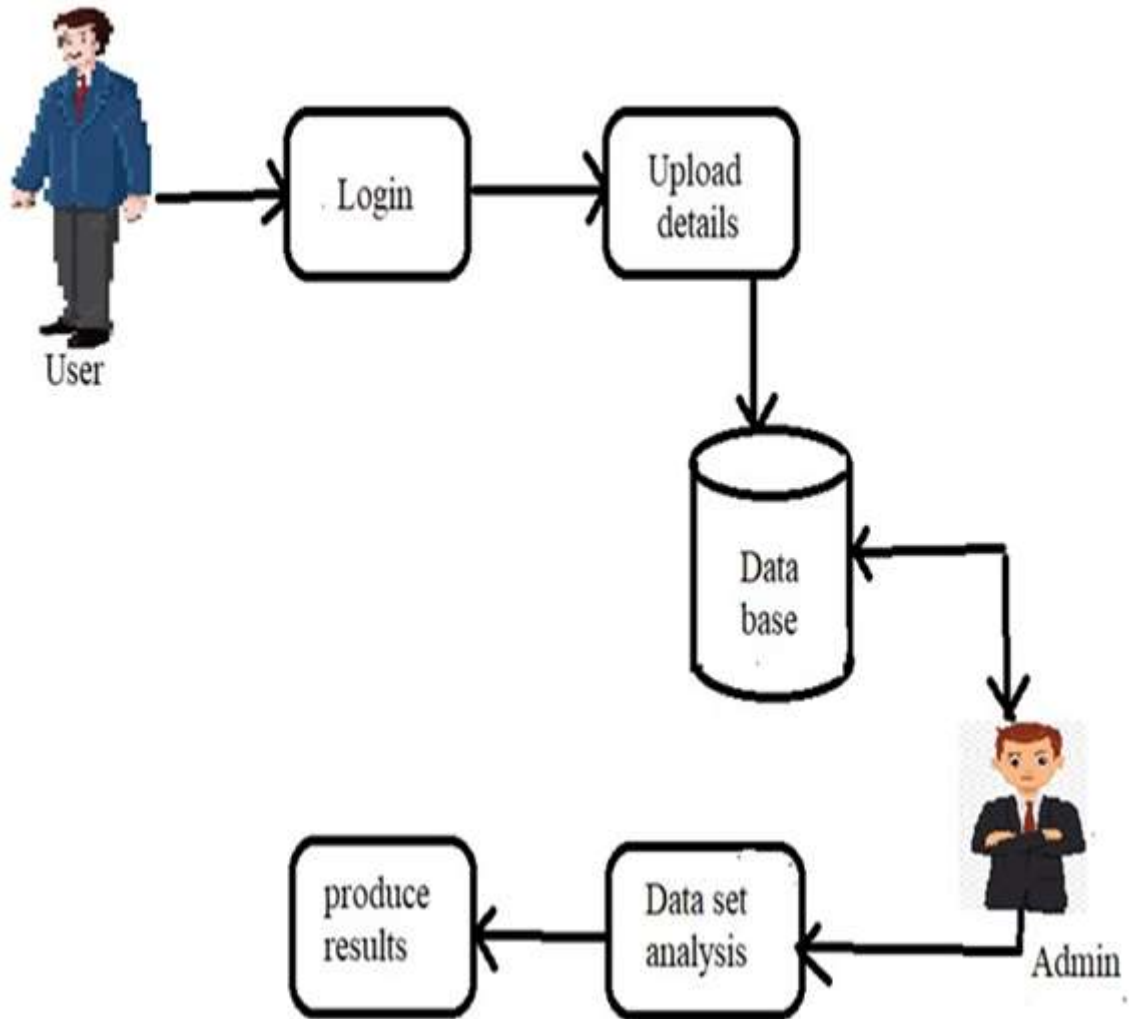


Fig-1: System architecture.

5 METHODOLOGIES

VECTOR SUPPORT MACHINE

A supervised machine learning approach called "Support Vector Machine" (SVM) can be applied to classification and regression problems. However, classification issues are where it's most frequently employed. This algorithm plots every data point as a point in n-dimensional space, where n is the number of features you have and each feature's value is a specific coordinate value. Then, we perform classification by identifying the hyper-plane that effectively distinguishes the



two classes (see the image below for an example). Support Vectors are nothing more than an individual observation's coordinates. The method that best separates the two classes (hyper-plane/line) is support vector machine.

Formally speaking, a support vector machine creates a hyper plane or set of hyper planes in a high- or infinite-dimensional space that may be applied to tasks like outliers detection, regression, and classification. It makes sense that the hyper plane with the greatest distance from the nearest training data point for any class will accomplish good separation, as the higher the margin, the lower the classifier's generalisation error is generally. The sets to discriminate are frequently not linearly separable in that space, despite the fact that the original problem may have been expressed in a finite dimensional space. In order to facilitate the separation, it was suggested that the original finite-dimensional space be mapped into a much higher-dimensional space.

6. RESULTS AND DISCUSSION SCREENSHOTS

Date Made Public	Company	City	State	Type of breach	Type of organization	Total Records	Description of incident	Information Source	Source URL	Year of Breach	Latitude	Longitude	Unnamed: 13	Unnamed: 14
1-13-2005	George Mason University	Fairfax	Virginia	HACK	EDU	32,000	Names, photos, and Social Security numbers of	Databases (DB)	NaN	2005-0	36.948224	-77.306173	NaN	NaN
1-13-2005	University of California, San Diego	San Diego	California	HACK	EDU	3,500	A Hacker breached the security of two Univers...	Databases (DB)	NaN	2005-0	32.715320	-117.157256	NaN	NaN
1-22-2005	University of Northern Colorado	Greeley	Colorado	PHISH	EDU	15,790	A hard drive with list of stolen is contained	Databases (DB)	NaN	2005-0	40.425314	-104.709102	NaN	NaN
2-12-2005	Science Applications International Corp. (SAIC)	San Diego	California	STAR	BSD	45,000	On January 25 thieves broke in into	Databases (DB)	NaN	2005-0	32.715320	-117.157256	NaN	NaN
2-15-2005	Chromehunt	Alpharetta	Georgia	PHISH	BSD	1.83.000	Fraudsters who presented themselves as legit...	Security Breach Letter	NaN	2005-0	34.075310	-84.294090	NaN	NaN

Fig-2: Displaying dataset.



Fig-3: Preprocessing the data.

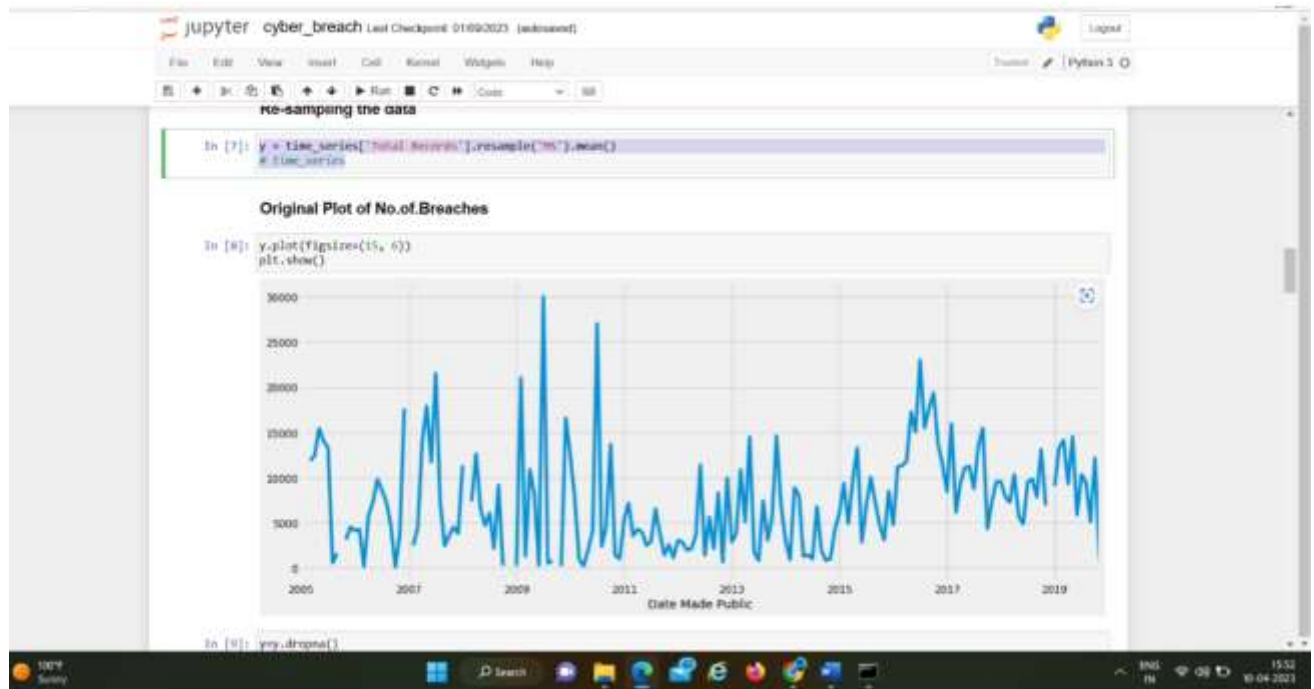


Fig-4: Plotting the number of Breaches.

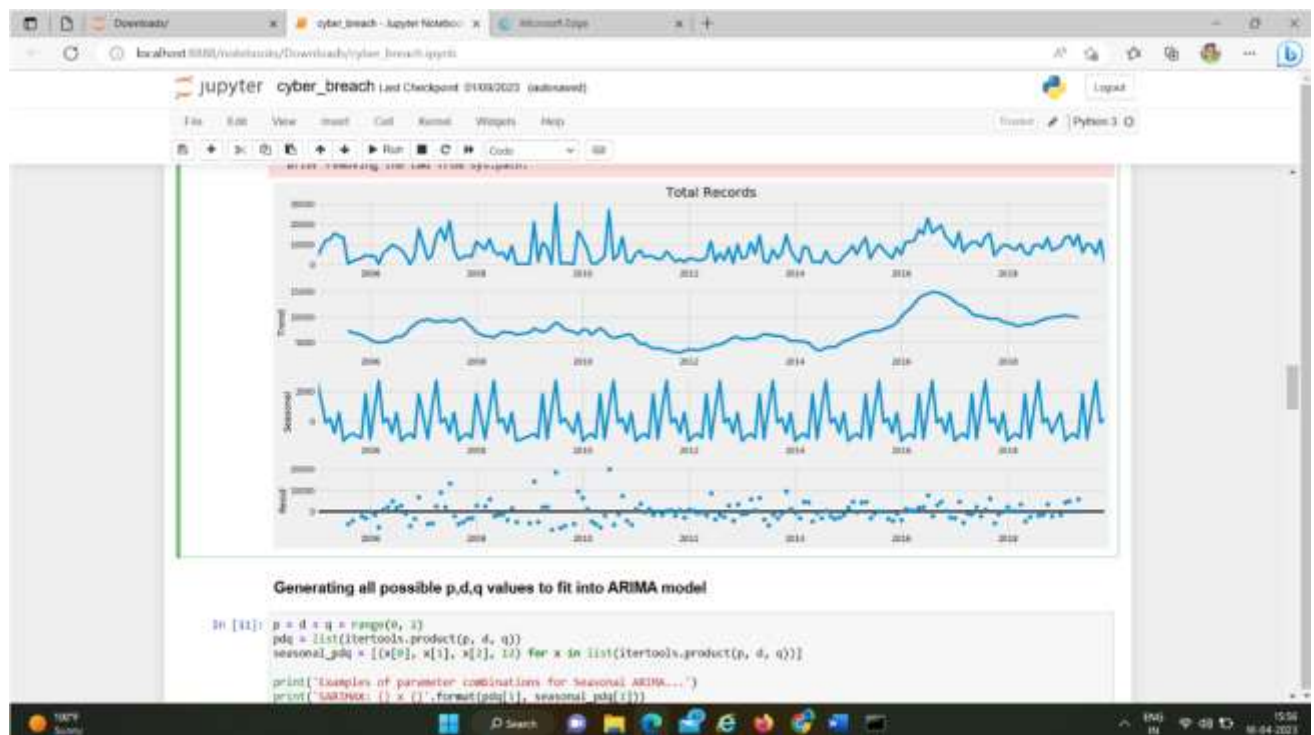


Fig-5: Generating all possible p,d,q values to fill into ARIMA model.

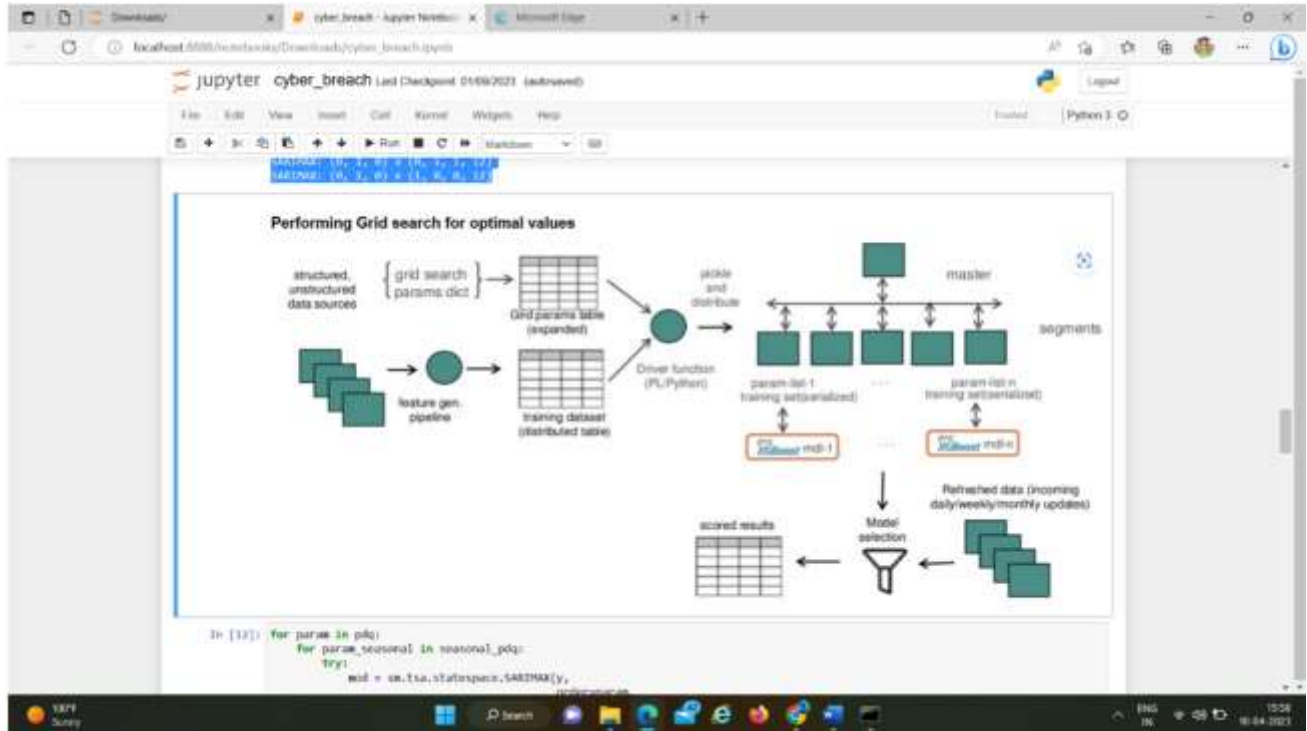


Fig-6: Performing Grid search for optional values.

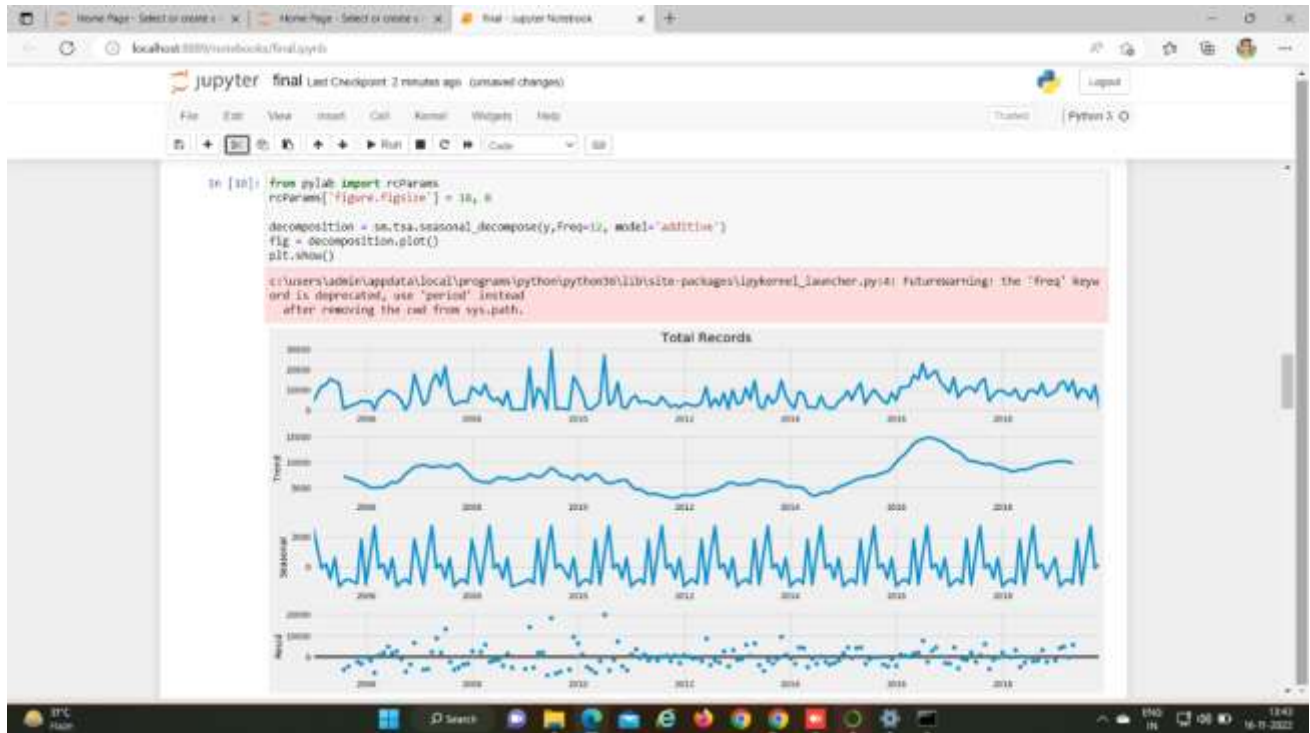


Fig-7: total Records

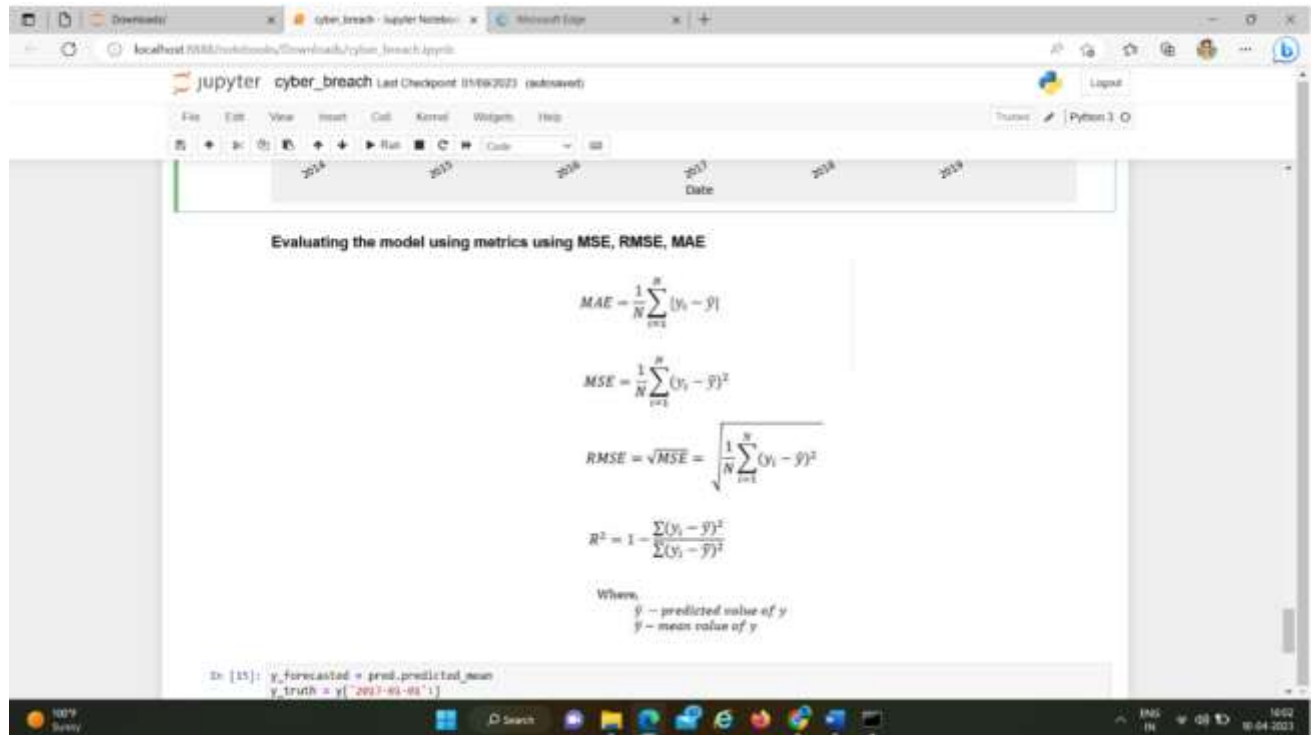


Fig-8: Evaluating the model using metrics using MSE, RMSE, MAE.

7. CONCLUSION AND FUTURE WORK

We looked at a dataset of cyber breaches from the angles of the time between occurrences and the amount of the breach, and we demonstrated that both should be modelled by stochastic processes rather than distributions. The statistical models created in this study exhibit acceptable fitting and prediction accuracies. We specifically suggest employing a copula-based approach to forecast the combined chance that an incident with a specific breach size would take place during the course of the future. According to statistical analyses, the approaches put forth in this research are superior than those found in the literature because the latter neglected temporal correlations as well as the relationship between incident arrival times and breach sizes. To gain additional insights, we carried out qualitative and quantitative assessments. We came up with a number of cybersecurity insights, such as the fact that while the frequency of cyber hacking breach instances is increasing, so is the size of the harm they do. The approach described in this research can be used to analyse datasets with a similar structure.

There are a lot of unresolved issues that require further study. Investigating how to forecast extremely big numbers and how to handle missing data (i.e., breach instances that are not reported) are two examples of investigations that are both intriguing and difficult. Estimating the precise times at which breach occurrences will occur is also important. Lastly, additional research is required to comprehend the predictability of breach situations (i.e., the maximum level of prediction accuracy). Creating a graphical user interface where users may login to contribute data, determine the type of assault involved, and observe the analysis of entered data is part of this project's future work.



8. REFERENCES

1. P. R. Clearinghouse, Privacy Rights Clearinghouse's Chronology of Data Breaches, Nov, 2017, [online] Available: <https://www.privacyrights.org/data-breaches>.
2. *Data Breaches Increase 40 Percent in 2016 Finds New Report From Identity Theft Resource Center and CyberScout*, Nov. 2017, [online] Available: <http://www.idtheftcenter.org/2016databreaches.html>.
3. C. R. Center, Cybersecurity Incidents, Nov. 2017, [online] Available: <https://www.opm.gov/cybersecurity/cybersecurity-incidents>.
4. *IBM Security*, Nov. 2017, [online] Available: <https://www.ibm.com/security/data-breach/index.html>.
5. *The 2016 Cyber Claims Study*, 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.
7. T. Maillart and D. Sornette, "Heavy-tailed distribution of cyber-risks", *Eur. Phys. J. B*, vol. 75, no. 3, pp. 357-364, 2010.
8. R. B. Security, Datalosdb, Nov. 2017, [online] Available: <https://blog.datalosdb.org>.
9. B. Edwards, S. Hofmeyr and S. Forrest, "Hype and heavy tails: A closer look at data breaches", *J. Cybersecur.*, vol. 2, no. 1, pp. 3-14, 2016.
10. S. Wheatley, T. Maillart and D. Sornette, "The extreme risk of personal data breaches and the erosion of privacy", *Eur. Phys. J. B*, vol. 89, no. 1, pp. 7, 2016.
11. P. Embrechts, C. Klüppelberg and T. Mikosch, *Modelling Extremal Events: For Insurance and Finance*, Berlin, Germany:Springer-Verlag, vol. 33, 2013.
12. R. Böhme and G. Kataria, "Models and measures for correlation in cyber-insurance", *Proc. Workshop Econ. Inf. Secur. (WEIS)*, pp. 1-26, 2006.
13. H. Herath and T. Herath, "Copula-based actuarial model for pricing cyber-insurance policies", *Insurance Markets Companies: Anal. Actuarial Comput.*, vol. 2, no. 1, pp. 7-20, 2011.
14. A. Mukhopadhyay, S. Chatterjee, D. Saha, A. Mahanti and S. K. Sadhukhan, "Cyber-risk decision models: To insure it or not?", *Decision Support Syst.*, vol. 56, pp. 11-26, Dec. 2013.
15. M. Xu and L. Hua, *Cybersecurity Insurance: Modeling and Pricing*, 2017, [online]



Available: <https://www.soa.org/research-reports/2017/cybersecurity-insurance>.

16. M. Xu, L. Hua and S. Xu, "A vine copula model for predicting the effectiveness of cyber defense early-warning", *Technometrics*, vol. 59, no. 4, pp. 508-520, 2017.
17. C. Peng, M. Xu, S. Xu and T. Hu, "Modeling multivariate cybersecurity risks", *J. Appl. Stat.*, pp. 1-23, 2018.
18. M. Eling and N. Loperfido, "Data breaches: Goodness of fit pricing and risk measurement", *Insurance Math. Econ.*, vol. 75, pp. 126-136, Jul. 2017.
19. E. Condon, A. He and M. Cukier, "Analysis of computer security incident data using time series models", *Proc. 19th Int. Symp. Softw. Rel. Eng. (ISSRE)*, pp. 77-86, Nov. 2008.
20. Y. Liu et al., "Cloudy with a chance of breach: Forecasting cyber security incidents", *Proc. 24th USENIX Secur. Symp.*, pp. 1009-1024, 2015.



Industrial Engineering Journal

ISSN: 0970-2555

Volume : 52, Issue 4, April : 2023