# BIG DATA NATURAL LANGUAGE PROCESSING AND SECURITY ISSUES

NAMRATHA N

Co Author – DR. SAMITHA KHAIYUM

Master of Computer Applications

Dayananda Sagar College of Engineering

*Abstract - At whatever point we talk about tremendous data, the stress is reliably over the security of the data. Of late the most gotten some answers concerning development is the Natural Language Processing. This new and moving development helps in dealing with the continually tackling security issues which are not completely settled using huge data. Starting with the huge information security issues, this paper oversees keeping an eye regarding the matters related to arrange wellbeing and information security using the Natural Language Processing development. Checking the prominent computerized attacks, for instance, phishing ID and spam disclosure, this paper moreover addresses issues on information confirmation and security like ID of Advanced Persistent Threat (APT) in DNS and shortcoming examination. The target of this paper is to give the diagram of how ordinary language dealing with can be used to address digital protection issues.*

*Keywords—cyber security, Big Data, phishing, spam detection*

## I. INTRODUCTION

Data is generally considered as a significant resource of each organization till date. For what reason is data significant? This inquiry has different answers. Beginning from one's very own data to the significant choices assumed the premise of the data removed from data, data is viewed as the bread and butter to each individual's life.

Data which enters the associations is getting higher and quicker step by step. Because of this increment in the volume and speed of the data, ordinary data sets were insufficient for the capacity of this data. In addition the data is normally unstructured making the conventional frameworks unequipped for handling or analysing it.

This brought about another worldview which was known as the Big data. Big data alludes to the data that is beyond the handling limit of the customary database frameworks. Enormous information is portrayed by the 3 V's which are volume, velocity and variety.

Distributed computing has been changing the IT business by adding versatility to the way in which IT is exhausted, enabling relationship to pay only for the resources and organizations they use. With a ultimate objective to reduce IT capital and operational utilization, relationship of all sizes are using Clouds to give the resources needed to run their applications. Mists vary in a general sense in their specific developments and execution, anyway much of the time give structure, stage, and programming resources as administrations.

The Big Data ascends out of the improvement of conveyed registering and cloud data accumulating. Distributed Computing uses systematized propels, which is wandered from the regular one. Then again we can say that Big Data is the ideal choice for distributed computing, since it needs enormous enrolling power and limit. Big Data and its examination have a considerable amount of benefits which are as per the following:

Supports better decision
making Reduces cost
Better Customer Relationship
Management (CRM)
Improves fraud detection
Saves the time and effort spent on analysis

With several advantages, big data also have few disadvantages attached to it.
They are:

- Security concerns
- Lacks potential stability
- Vulnerable by nature

Among the ones mentioned above the most talked about disadvantage is the big data security.

Since the volume and the speed of information are truly developing, there is a requirement for giving security to the data. Since the data may incorporate individual data, spending subtleties, compensation data, credit archives and so on there is a need to safeguard the protection of each individual or an organization's data.

These security issues of big data can be conquered utilizing the Natural Language Processing (NLP). NLP is a method in the field of software engineering which manages handling human language with the assistance of specific models and algorithms. NLP can be utilized to address the digital protection issues in a superior way.

The organization of the paper begins with the big data security issues followed by NLP as an answer for security issues and how NLP is utilized in data security followed by the conclusion and references

## II. BIG DATA SECURITY ISSUES

Prior to going into the conversation about the security issues, let is first comprehend the distinction among security and privacy. There is a misinterpretation that security and privacy both mean something very similar. Yet, the truth of the matter is that they are two unique points. The comprehension of the distinctions is a significant piece of the conversation. The table beneath gives a reasonable picture in regards to the distinction among security and privacy.

| | Parameter | Privacy | Security |
|---|---|---|---|
| a | Definition | Security is the appropriate use of the client data | Security alludes to the CIA triad i.e. confidentiality, integrity and availability of data |
| b | Offerings | It offers the capacity to choose which data of an individual goes where | It offers the confidence that each choice of an individual is regarded. |
| c | Goal | It typically applies to a client's privileges to shield their data. | The objective of safety framework is to give security to an individual or an organization. |
| d | Practices followed | It is feasible to have poor privacy and great security rehearses. | It is hard to have great protection rehearses without a decent data security program |

Since we have perceived the fundamental distinction among privacy and security, let us talk about the various difficulties to big data security.

There are four unique parts of Big Data security. They incorporate framework security, privacy of data, management data, and integrity and receptive security portrayed as follows:

### 1. INFRASTRUCTURE SECURITY

Infrastructure security manages the utilization of numerous PCs wherein the data are put away, but at the same time are imitated along the group. And further more discovering a design that will guarantee the full accessibility of the infrastructure is the highest need.

### 2. DATA PRIVACY

Data security is the most concerned subject of the organization that utilization Big Data methods. Organization ought not have absolute opportunity to utilize that data without the information on the concerned individual, regardless of whether they need to acquire some profit by the utilization of that data. The primary themes on data security include:

• Access control
• Cryptography
• Confidentiality

### 3. DATA MANAGEMENT

Data Management is worried about how the data must be managed whenever it is put away. It not just tells the best way

to get the data that is put away in the Big Data framework, yet in addition how to share that data without losing its integrity.

## 4. INTEGRITY AND REACTIVE SECURITY

Big Data for the most part gets data from a few roots which may incorporate organized just as unstructured data. This builds the significance of checking the data integrity so it is appropriately utilized without harming the privacy.

Security of a tremendous measure of data can be improved by utilizing procedures like confirmation, validation and encryption.

At the point when we talk about verification and authentication numerous algorithms go to our brains. One such broadly utilized calculation is the advanced mark calculation utilizing RSA. The digital signatures are utilized for checking the realness of reports. Digital signatures are notable for giving validation, respectability and non-disavowal. Despite the fact that it is a successful technique for keeping up credibility and uprightness, there were numerous burdens related with this strategy.

To beat such disadvantage the forthcoming moving innovation which is the natural language handling is utilized. Natural language preparing is the field of software engineering that manages consequently parsing and handling human language. With natural language processing machine can handle a book independent of any language or slang.

## III. NLP AS A SOLUTION TO SECURITY ISSUES

Data researchers apply NLP to defeat the large data security by building a model to all the more likely distinguish and comprehend noxious code. One such methodology towards a NLP procedure explicitly intended for malware examination is a Malicious Language Processing system. The objective of this Malicious Language Processing system is to utilize NLP to address one of the security space's most difficult big data issues via robotizing the recognizable proof of noxious code covered up inside the first code.

## IV. HOW NLP IS USED IN INFORMATION SECURITY

Data put away in any association requires assurance. How is this data gotten? How is NLP utilized in getting this information? The response to this is extremely straightforward. NLP is utilized to recognize pernicious spaces, phishing

assaults, distinguish spam sends and so on. The primary security issues include:

1. FIGHTING SPAM
2. PHISHING IDENTIFICATION
3. SOURCE CODE VULNERABILITY ANALYSIS
4. DOMAIN GENERATION ALGORITHM CLASSIFICATION

Domain Generation Algorithms (DGAs)

They are a class of computations that produce tremendous amounts of territory names. DGAs license malware to make a far reaching number of territories consistently; by a long shot the greater part of them are unregistered. The enormous amounts of unregistered spaces are used to camouflage the selected ones, allowing the tainted botnets to stay away from area and counteraction by imprint or IP-reputation based security disclosure structures.

➤ DGA identification objectives

There are three elements of DGA distinguishing proof, with each subsequent measurement identifying with a climb in earnestness.

✓ If a DGA is recognized, it infers that in any event one of the systems is corrupted by DGA-based malware. The essential objective is to perceive the affected structures and help to hold them back from sullying/harming further.

✓ The following objective is to choose if a given DGA space name is enlisted or not. Thusly, the second fragment of a DGA revelation system is to isolate enrolled spaces from the unregistered ones.

✓ The last and the most essential objective of a DGA recognizable proof structure is to choose if the requesting was compelling with the selected framework and if the contact was made between the sullied system and the host worker. If such development is perceived, some mischief may have recently been done.

➤ DGA working

DGA activity is recognized by getting and taking apart framework packs, as a rule in five general advances.

The 5 stages are as per the following:

Stage 1 – Detect Domain Name Server (DNS) Application

The acknowledgment methodology begins with the DNS request just as response messages. DNS is an essential Internet show, and most firewalls have a procedure to allow dynamic DNS traffic on its held port. It is prominent that the held port for the DNS traffic is 53. Nevertheless, a developer may misuse the finding out about port 53 to send its traffic without adherence to the standard DNS message position.

Stage 2 – Extraction of the Domain Names

At the point when a framework is recognized as DNS, the space names in the DNS question and response messages ought to be isolated.

Stage 3 – Detect any DGA

Examination ought to be performed on the spaces eliminated from DNS messages to choose if they are DGAs or not.

Stage 4 – Detect Registered DGA Domains

In order to recognize if a DGA space name is enrolled, the DNS responses ought to be checked.

Stage 5 – Detect Traffic to Registered DGA Domains

The most existing DGA ID structures focus on recognizing whether a space name is a DGA region; they consistently neglect the most crucial request: Is there any traffic that has been shipped off the enlisted DGA regions? To react to this request, the DGA space revelation should be immovably joined with framework traffic appraisal and the results ought to be resonated back to the traffic survey engine going before any damage is done.

There is another movement which isn't really a piece of area anyway when there is getting together with a neutralizing activity structure, for instance, a Firewall, a standard should be inserted promptly to impede all the traffic to the enrolled areas.

➢ Detection of Advanced Persistent Threat (APT) in DNS

Normally, the APT uses equivalent spaces to that of a standard programming organization (Microsoft, Adobe, Java, etc.) and demands programming revives. The territory names are inventively made attempting to keep up the authenticity of the spaces.

For example: 1inkedin.net

NLP Rank calculation is intended to identify the fake marked areas that regularly fill in as spaces for focused assaults. It utilizes the base alter distance calculation which uses this computation on substrings to check the word separate between the two strings. It assesses the closeness between the two strings by counting the amount of assignments/adjusts needed to change more than one string into another. The operations permitted are:

1. Deletion
2. Substitution
3. Insertion

➢ What does the base alter distance algorithm do?

The minimum alter distance between 2 strings computes the base number of alters it takes (ex. addition, erasure, and replacement) to transform string A into string B.

Algorithm: Minimum edit-distance algorithm

Input: Two strings $A = a_1 \dots a_m$ and $B = b_1 \dots b_n$
Output: The matrix $D = (D_{ij})$

**BEGIN**

$D(0,0)=0$ **for** $i=1$ **to** $m$
**do** $D(i,0)=i$ **for** $j=1$ **to** $n$
**do** $D(0,j)=j$ **for** $i=0$ **to** $m$
**do**
**for** $j=0$ **to** $n$ **do**
$D(i,j) = \min (D(i-1,j)+1, D(i,j-1)+1, D(i-1,j-1)+ c(a_i,b_j))$
**end for end**
**for**
**return** $D(m,n)$

**END**

Example 1:

Introductory Domain: linkedin.com
Objective Domain: linkeddin.com For this model, there is 1 alter because of replacement, making the punishment 1.
Example 2:

Beginning Domain: facebook.com
Objective Domain: faceboook.com For this model, there is 1 alter because of inclusion, making the punishment 1.

Example 3:

Introductory Domain: youtube.in
Objective Domain: yutube.in For this model, there is 1 alter because of cancellation, making the punishment 1.

Identification of Phishing and spam mails

Phishing is a technique used to tempt people to reply with some sensitive information. These spam/phishing messages share a particular model for all expectations and reason like:

✓ A Promising offer enticing the client with cash: Such messages notice the sum, to entice the client who consequently will answer with his/her subtleties with which they endeavor to remove individual data from the client.

✓ Some feeling of crisis: This is a social building strategy where the programmer plays mind diversions with the client depicting some pressing circumstance like the need to exchange cash outside the nation requesting the bank subtleties of the client to move cash into.

In order to apply NLP to this marvel, we need a get-together of documents. This is for regular web data yet when we apply NLP to weaknesses or Malware we should eliminate data by applying express procedures that are static and dynamic to do the examination

Fighting Spam

This territory resembles the previous one anyway is dealt with in an surprising way. Spam distinguishing proof is done very much arranged which is explained as seeks after:

a) Email Input: The email input given to the spam area show.
b) URL Source Check: The moving toward Email's source URL is checked. URL Blacklist Database is an information base that contains each one of the URLs which have been recognized to be a spam in the midst of the planning stage.

c) Threshold Counter: A counter which screens the amount of messages sent throughout some time frame t.

d) Keyword Extractor: This movement secludes the entire message into tokens and sends these tokens as expressions to the classifier.

e) Classifier: It acknowledges the expressions as data and bunches the email into a specific class. This movement absolutely depends upon the past advance as the tokens are made there.

f) NLP Engine: It takes the unclassified Email and its class as a data and after that structures it using the quantifiable NLP approach.

Difficulties of bid data and possible solutions.

Difficulties that enormous information faces in the present situation are going from the plan of preparing framework at the lower layer to investigation implies at the high layer. Here, we have examined great difficulties of huge information:

1) To store huge measure of information Storage
2) To share enormous measure of information Transfer
3) To picture, clean, dissect and search large information issues. Malware and tricks are not, at this point restricted as web.

## V. CONCLUSION

Based on the analysis done during the assortment of data from various sources, one thing is certain that the zero in is just on tracking down a solid answer for the security issues. The methodology may vary starting with one creator then onto the next be that as it may, the objective remaining parts as before. We went over numerous approaches, for example, recognizing the punishment utilizing the base alter distance calculation or utilizing different strategies to characterize a mail as in to if it is a spam or utilizing the machine interpretation strategy to encode the data and so forth that are especially not quite the same as each other. It was likewise seen that each approach

accompanies its own advantages and disadvantages. The choice lies with the client to which calculation to best use in the most exceedingly terrible circumstance. There isn't anything known as the best approach. It is just about utilizing a proper method to defeat the issue in a speedy and viable way.

Digital protection is a mind boggling and a cast subject, whose understanding requires information and aptitude from various disciplines, including yet not restricted to software engineering and data innovation. At last, the inquiry isn't the ticket

the network safety issue can be tackled, but instead how it can be made reasonable.

## VI. REFERENCES

[1]    Abid Mehmood, Iynkaran Natgunanathan, Yong Xiang, Senior Member, IEEE, Guang Hua, Member, IEEE, and Song Guo, Senior Member, IEEE, "Protection of Big Data Privacy", Citation information: DOI 10.1109/ACCESS.2016.2558446

[2]    R.Kalaivani, "Security Perspectives on Deployment of Big Data using Cloud: A Survey", International Journal of Advanced Networking & Applications (IJANA) Volume: 08, Issue: 05 Pages: 5-9 (2017) Special Issue

[3]    K.P.Maheswari[1], P.Ramya[2],S.Nirmala Devi[3] 1,2Department of IT & Networking, 3Department of Computer Science, Subbalakshmi Lakshmipathy College of Science (India) , "Study and analyses of security levels in big data and cloud computing" RTESHM-17

[4]    https://journalofbigdata.springeropen.com/articles/10.1186/s40537-016-0059-y

[5]    P.Joseph Charles1, I.Carol, S.Mahalakshmi J. Clerk Maxwell, "Big Data Security an Overview, A Treatise on Electricity and Magnetism", 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.

[6]    Trupti V. Pathrabe, "Survey on Security Issues of Growing Technology: Big Data", IJIRST (International Journal for Innovative Research in Science and Technology) 2017.

[7]    Deepak Kumar Vishwakarma, "Domain Name Generation Algorithms"-Master thesis-2017

[8]    https://medium.com/@akarshasehwag/you-wake-up-one-nice-morningand-open-up-your-laptop-to-check-the-mails-and-find-out-that-you-wona-7a10471e339b

[9]    K. Ingols, M. Chu, R. Lippmann, S. Webster, and S. Boyer, "Modeling modern network attacks and countermeasures using attack graphs", in Proc. of the 25th Annual Computer Security Applications Conference, Honolulu, HI, USA, Dec. 2009, pp. 117–126.

[10]    Rohit Giyanani, Mukti Desai, "Spam Detection using Natural Language Processing" IOSR(International Oraganization of Scientific Research) Journal of Computer Engineering (IOSR-JCE) e-ISSN: 2278-0661,pISSN: 2278-8727, Volume 16, Issue 5, Ver. IV (Sep – Oct. 2014), PP 116-119.

[11]    https://umbrella.cisco.com/blog/2015/03/05/nlp-apt-dns/

[12]    Rohit Giyanani, Mukti Desai, "Spam Detection using Natural Language Processing", IOSR Journal of Computer Engineering (IOSR-JCE) e-ISSN: 2278-0661,p-ISSN: 2278-8727, Volume 16, Issue 5, Ver. IV (Sep – Oct. 2014), PP 116-119