# Prediction of cyber attacks using data science techniques

J.B.Saraswathi[1], L.Karthika.[2]
[1]PG Scholar, [2]Assistant professor
Department of Computer science and Engineering,
Gojan School of Business and Technology, Redhills, Chennai

**Abstract**— Cyber-attacks end to destroy or virulently manipulate a computing terrain or structure, as well as disrupt data integrity or crack all information. This poses a threat to the organization, maybe performing in data loss. The data from device detectors is collected as big data, which has a wealth of information that can be utilized for targeted assaults. Although being methodologies, models, and algorithms have given the foundation for cyber-attack prognostications, new models and algorithms grounded on data representations other than task-specific ways are needed. Its on-linear information processing armature, on the other hand, can be customized to learn indispensable data representations of network business in order to classify different types of network attacks. In this study, we treat cyber-attack vaticination as a bracket issue, in which networking sectors must use machine literacy approaches to read the type of network assault from a given dataset. The supervised machine learning fashion( SMLT) is used to assay a dataset in order to capture multiple pieces of information, similar as variable identification, uni-variate analysis, bi-variate andmulti-variate analysis, missing value treatments, and so on. A comparison of machine literacy algorithms was conducted to estimate which algorithm is the stylish accurate at prognosticating the types of cyber-attacks. DOS Attack, R2L Attack, U2R Attack, and inquiry Attack are the four types of attacks we classify. The findings reveal that the suggested machine literacy algorithm fashion has the stylish delicacy with entropy computation, perfection, recall, F1 Score, perceptivity, particularity, and entropy.

**Index Terms**— Cyber-attack, DOS Attack, R2L Attack, U2R Attack, inquiry Attack, GUI ,Logistic regression, Decision tree, Random forest
.

————————————————————— ◆ —————————————————————

## 1 INTRODUCTION

### 1.1 Data wisdom.

This is an interdisciplinary field that employs medical procedures, processes, algorithms, and structures to prize knowledge and perceptivity from unshaped and structured information, as well as to track knowledge and practicable perceptivity across a wide range of mileage areas. The International Federation of Bracket Societies( IFCS) has been the top convention to concentrate on records technology since 1996. nevertheless, the description came a work in progress. The term " data period " have come first chased in 2008 thru manner of way ofD.J. Patil, and Jeff Hammerbacher, the colonist leads of data and analytics sweats at LinkedIn and Facebook. In a great deal much lower than a decade, it has crop as one of the most over to date and utmost trending professions with inside the request. This content( Data technology Is the sector of study that combines region understanding, programming chops, and understanding of mathematics and data to prize considerable perceptivity from data). This can be defined as a combination of mathematics, business enterprise wit, tools, algorithms and tool analyzing ways, all of which help us in chancing out the hidden perceptivity or patterns from raw data which can be of use with inside the conformation of huge business.

It is nothing but which questions need to be answered and where to find the answer. They've business sense and logical chops, as well as the capability to prize, clean, and present data. Fact scientists help businesses find, organize, and assay large quantities of unshaped data.
A. needed chops for a data scientist Python, SQL, Scala, Java, R, MATLAB are the programming needed. Natural Language Processing, Bracket, Clustering are Machine literacy ways needed. Tableau, SAS,D3.js, Python, Java, R libraries are the Data Visualization ways needed. MongoDB, Oracle, Microsoft Azure, Cloud era are the Big data platforms needed.

### 1.2 Artificial Intelligence

This( AI) refers to the simulation of mortal intelligence in computers that are programmed to act and move like people. The time period can be used to any device that's well- known for displaying advancements in mortal intellect, similar as mastery and problem- working. Artificial intelligence( AI) refers to intelligence demonstrated by machines rather than mortal or carnal intelligence. Leading AI handbooks define the field as the study of" sensible agents," or any device that's able of perceiving its terrain and acting in ways that increase its chances of achieving its pretensions. Some well- known debtors use the term" synthetic intelligence" to describe machines that mimic" cognitive" features similar as" learning" and" hassle" that people associate with mortal studies, but Important AI experimenters differ with this description. Artificial intelligence( AI) is the simulation of mortal intelligence strategies by computers, especially computer structures. Professional structures, herbal language processing, speech

fissionability, and device vision are exemplifications of AI packages. Superior internet hunt machines, advising systems( as used by Youtube, Amazon, and Netflix), understanding mortal speech( as used by Siri or Alexa), tone- driving motorcars(e.g. Tesla), and contending at the loftiest position in strategic recreation systems are all exemplifications of AI operations( conforming of chess and Go), The AI effect is a marvels that occurs when computers get further able. As machines come more able, scores that need" intelligence" are generally removed from the description of AI. For illustration, optic joe or woman fashionability is occasionally overlooked. Fantastically fine statistical device learning reigned the assiduity in the first numerous times of the twenty-first century, and this system has proven fantastically successful, aiding in the resolution of numerous delicate challenges in business and academia. The colorful subfields of AI exploration are centred on precise solicitations and the operation of precise tools. logic, knowledge representation, planning, mastery, herbal language processing, belief, and the capability to carry and handle bias are all common AI pretensions. Some of the area's long- term pretensions include general intelligence (the capability to break any problem). To address these difficulties, AI experimenters employ strategies similar as seek and fine optimization, formal sense, synthetic neural networks, and statistics, chance, and economics-grounded ways. AI also draws on computer wisdom, psychology, linguistics, gospel, and a variety of other disciplines. The field was renamed after the supposition that mortal intelligence" might be so precisely defined that a device could be erected to imitate it." This heightens philosophical debates on the ethics and studies of creating synthetic beings with mortal- suchlike intellect. Given the age, these issues were examined through myth, fabrication, and gospel. Science fabrication and futurology have also suggested that, due to its enormous capacity and strength, AI could come an empirical trouble to humans. Carriers were scrabbling to vend how their services and products incorporate AI as the hoopla around AI grew. Constantly, what they relate to as AI is clearly one aspect of AI, videlicet device mastery. For designing and training device mastery algorithms. AI necessitates a foundation of specialized tackle and software. Although no single programming language is synonymous with AI, a sprinkle stand out, including Python, R, and Java. In general, AI systems work by collecting large quantities of classified educational data, analyzing the data for correlations and patterns, and also applying those styles to produce prognostications about unborn countries. In this system, a chatbot fed cases of textual content exchanges can study how to produce cultures similar as face- to- face relations, or a print fashionability device can learn how to find and define widgets in images by looking at hundreds of thousands of exemplifications. learning, logic, and tone- correction are three cognitive capabilities that AI programming focuses on. ways for learning This branch of AI programming focuses on gathering data and formulating rules for converting those data into practicable data. The rules, also known as algorithms, give computer bias with step- by step instructions on how to do a specific task. logic ways. This section of AI programming focuses on choosing the stylish collection of rules to achieve a asked result. tone-correction ways.

## 1.3 Machine Learning

This content( Machine literacy) entails anticipating the future from beyond the records. Machine literacy( ML) is a type of artificial intelligence( AI) that allows computer systems to learn without having to be explicitly programmed. Machine literacy focuses on the creation of computer programmes that may change when exposed to new data, as well as the principles of Machine literacy, similar as the design of a simple system learning set of rules using Python. The operation of specialised algorithms is part of the education and vatication process. It feeds the education records to a set of rules, and the set of rules uses the education records to make prognostications on new examination records. Machine literacy can be classified into a many different orders. There are three types of literacy supervised literacy, unsupervised literacy, and  underpinning literacy. The input records and the accompanying labelling are collectively submitted to supervised literacy software, and the study records must be categorised by a person beforehand. There are no groups for unsupervised literacy. It handed a set of guidelines for literacy. The clustering of the enter records must be parented out by this set of Criteria. Eventually, underpinning learning stoutly interacts with its surroundings and receives positive or negative feedback to ameliorate its performance. To discover styles in python that affect in useful perceptivity, data scientists use a variety of different forms of system learning algorithms. At a high position, those unique algorithms can be divided into two groups grounded on how they" study" records to make prognostications supervised and unsupervised literacy.

## 2.EXISTING SYSTEM

They proposed first to produce a contrastive tone- supervised studying to the incongruity discovery hassle of attributed networks. CoLa, is particularly includes 3 factors contrastive illustration brace slice, GNN- primarily grounded completely contrastive studying interpretation, and multi round slice- primarily grounded completely anomaly rating calculation. Their interpretation captures the connection among every knot and its neighbouring shape and makes use of an anomaly- associated thing to educate the contrastive studying interpretation. We agree with that the proposed frame opens a brand new possibility to extend tone supervised studying and contrastive studying to an adding number of graph anomaly discovery operations. The multiround anticipated conditions via way of means of the contrastive studying interpretation are in addition used to assess the abnormality of every knot with statistical estimation. The training member and the conclusion member. In the training member, the contrastive studying interpretation is educated with tried illustration dyads in an unmonitored fashion. After that the
incongruity standing for every knot is acquired withinside the conclusion member.

## 3.PROPOSED SYSTEM

The proposed model is to construct a machine literacy model for anomaly discovery. Anomaly discovery is an critical system for spotting fraud conditioning, suspicious conditioning, community intrusion, and different peculiar occasions that could have awful significance still are hard to descry. The contrivance studying
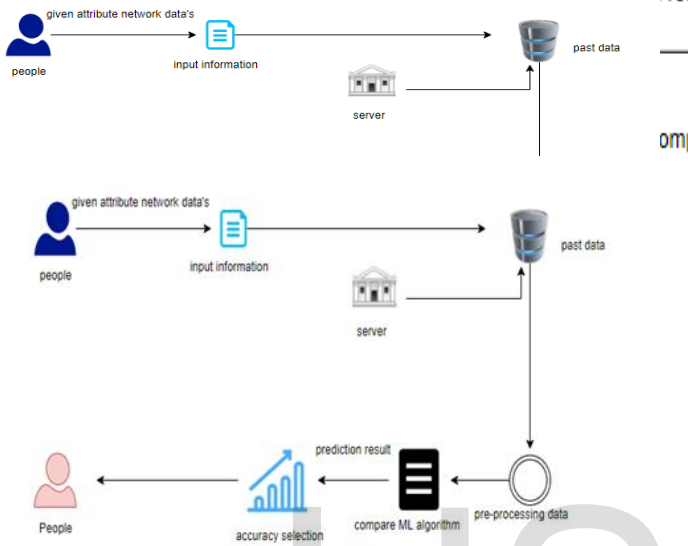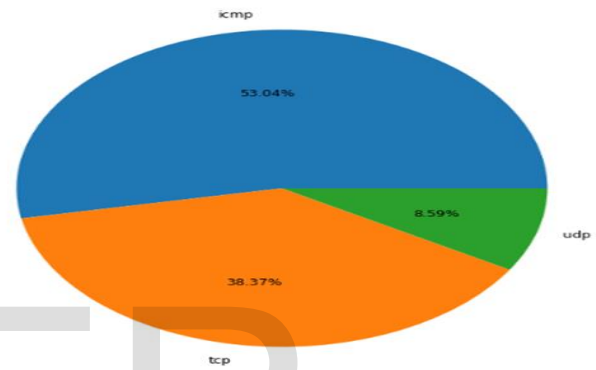
given attribute network data's

input information

people

server

past data



given attribute network data's

input information

people

server

past data

prediction result

People

accuracy selection

compare ML algorithm

pre-processing data



protocol_type (%) (Per Count)

icmp

53.04%

8.59%

udp

38.37%

tcp

protocol_type

## 5.MODULE DESCRIPTION

### 5.1. Variable Identification Process / data validation process.

Validation techniques in machine learning are used to get the error rate of the Machine Learning (ML) model, which can be considered as close to the true error rate of the dataset. If the data volume is large enough to be representative of the population, you may not need the validation techniques. However, in real-world scenarios, to work with samples of data that may not be a true representative of the population of given dataset. To finding the missing value, duplicate value and description of data type whether it is float variable or integer. The sample of data used to provide an unbiased evaluation o f a model fit on the training dataset while tuning model hyper parameters. The evaluation becomes more biased as skill on the validation dataset is incorporated into the model configuration. The validation set is used to evaluate a given model, but th is is for frequent evaluation. It as machine learning engineers uses this data to fine-tune the model hyper parameters. Data collection, data analysis, and the process of addressing data content, quality, and structure can add up to a time -consuming to do list. During the process of data identification, it helps to understand your data and its properties; this knowledge will help you choose which algorithm to use to build your model. For example, time series data can be analyzed by regression algorithms; classification algorithms can be used to analyze discrete data. (For example to show the data type format of given dataset) Importing the library packages with loading given dataset.

### 5.2. Performance measurements of DoS attacks

In computing, a denial-of-service attack (DoS attack) is a cyber-attack in which the perpetrator seeks to make a machine or network resource unavailable to its intended users by temporarily or indefinitely disrupting services of a host connected to the internet. Denial of service is typically accomplished by flooding the targeted machine or resource with superfluous requests in an attempt to overload systems and prevent some or all legitimate requests from being fulfilled. In a distributed denial-of-service attack (DDoS attack), the incoming traffic flooding the victim originates from many different sources. This effectively makes it impossible to stop the attack simply by blocking a single source. A DoS or DDoS attack is analogous to a group of people crowding the entry door of a shop, making it hard for legitimate customers to enter, disrupting trade. A distributed denial-of-service (DDoS) is a large-scale DoS attack where the perpetrator uses more than one unique IP Address, often thousands of them. A distributed denial of service attack typically involves more than around 3–5 nodes on different networks; fewer nodes may qualify as a DoS attack but is not a DDoS attack. Since the incoming traffic flooding the victim originates from different sources, it may be impossible to stop the attack simply by using ingress filtering. It also makes it difficult to distinguish legitimate user traffic from attack traffic when spread across multiple points of origin. As an alternative or augmentation of a DDoS, attacks may involve forging of IP sender addresses further complicating identifying and defeating the attack. An application layer DDoS attack (sometimes referred to as

layer 7 DDoS attack) is a form of DDoS attack where attackers target application layer processes. The attack over-exercises specific functions or features of a website with the intention to disable those functions or features. This application -layer attack is different from an entire network attack, and is often used against financial institutions to distract IT and security personnel from security breaches.

**Classifier report of support vector classifier results:**

| | Precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.79 | 0.97 | 0.87 | 681 |
| 1 | 0.68 | 0.19 | 0.38 | 119 |
| Accuracy | | | 0.78 | 988 |
| Macro | | 0.73 | 0.58 | 900 |

## 5.3. Performance measurements of R2L attacks

Now-a-days, it is very important to maintain a high-level security to ensure safe and trusted communication of information between various organizations. But secured data communication over internet and

**Classification report of Decision Tree Results:**

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.64 | 0.90 | 0.74 | 535 |
| 1 | 0.62 | 0.25 | 0.35 | 365 |
| accuracy | | | 0.63 | 900 |

**Classification report of Decision Tree Results:**

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.64 | 0.90 | 0.74 | 535 |
| 1 | 0.62 | 0.25 | 0.35 | 365 |
| accuracy | | | 0.63 | 900 |

## 5.4. Performance measurements of U2R attacks

Remote to local attack (r2l) has been widely known to be launched by an attacker to gain unauthorized access to a victim machine in the entire network. Similarly, user to root attack (u2r) is usually launched for illegally obtaining the root's privileges when legally accessing a local machine. Buffer overflow is the most common of U2R attacks. This class begins by gaining access to a normal user while sniffing around for passwords to gain access as a root user to a computer resource. Detection of these attacks and prevention of computers from it is a major research topic for researchers throughout the world

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.64 | 0.90 | 0.74 | 535 |
| 1 | 0.62 | 0.25 | 0.35 | 365 |

**Classification report of Decision Tree Results:**

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.64 | 0.90 | 0.74 | 535 |
| 1 | 0.62 | 0.25 | 0.35 | 365 |
| accuracy | | | 0.63 | 900 |

## 5.5. Performance measurements of Probe attacks

Probing attacks are an invasive method for bypassing security measures by observing the physical silicon implementation of a chip. As an invasive attack, one directly accesses the internal wires and connections of a targeted device and extracts sensitive information. A probe is an attack which is deliberately crafted so that its target detects and reports it with a recognizable "fingerprint" in the report. The attacker then uses the collaborative infrastructure to learn the detector's location and defensive capabilities from this report. This is an attack where the attacker attempts to gather information about the target machine or the network, to map out the network. Information about target may reveal useful information such as open ports, its IP address, hostname, and operating system. Network Probe is the ultimate network monitor and protocol analyzer to monitor network traffic in real-time, and will help you find the sources of any network slow-downs in a matter of seconds.

## 5.6 Performance measurements of overall network attacks

Increasingly, attacks are executed in multiple steps, making them harder to detect. Such complex attacks require that defenders recognize the separate stages of an attack, possibly carried out over a longer period, as belonging to the same attack. Complex attacks can be divided into exploration and exploitation phases. Explorat ion involves identifying vulnerabilities and scanning and testing a system. It is how an attacker gathers information about the system. Exploitation involves gaining and maintaining access. At this stage, the attacker applies the know-how gathered during the exploration stage. An example of a complex attack that combines exploration and exploitation is a sequence of a phishing attack, followed by an exfiltration attack. First, attackers will attempt to collect information on the organization they intend to attack, e.g., names of key employees. Then, they will craft a targeted phishing attack. The phishing attack allows the attackers to gain access to the user's system and install malware. The purpose of the malware could be to extract files from the user's machine or to use the user's machine as an attack vector to attack other machines in the organization's network. A phishing attack is usually carried out by sending an email purporting to come from a trusted source and tricking its receiver to click on a URL that results in installing malware on the user's system. This malware then creates a backdoor into the user's system for staging a more complex attack. Phishing attacks can be recognized both by the types of keywords used in the email (as with a spam ema il), as well as by the characteristics of URLs included in the message. Features that have been used successfully to detect phishing attacks include URLs that include IP addresses, the age of a linked-to domain, and a mis match

between anchor and text of a link.

## 6.7. GUI based prediction results of Network attacks

GUI means Graphical User Interface. It is the common user Interface that includes Graphical representation like buttons and icons, and communication can be performed by interacting with these icons rather than the usual text-based or command based communication. A common example of a GUI is Microsoft operating systems. The graphical user interface (GUI) is a form of user interface that allows users to interact with electronic devices through graphical icons and audio indicator such as primary notation, instead of text-based user interfaces, typed command labels or text navigation. GUIs were introduced in reaction to the perceived steep learning curve of command-line interfaces (CLIs) which require commands to be typed on a computer keyboard. The actions in a GUI are usually performed through direct manipulation of the graphical elements. Beyond computers, GUIs are used in many handheld mobile devices such as MP3 players, portable media players, gaming devices, smartphones and smaller household, office and industrial controls. The term GUI tends not to be applied to other lower-display resolution types of interfaces, such as video games (where head-up display (HUD) is preferred), or not including flat screens, like volumetric displays because the term is restricted to the scope of two-dimensional display screens able to describe generic information, in the tradition of the computer science research at the Xerox Palo Alto Research Centre. Graphical user interface (GUI) wrappers find a way around the command-line interface versions (CLI) of (typically) Linuxand Unix-like software applications and their text-based user interfaces or typed command labels. While command-line or text-based applications allow users to run a program non-interactively, GUI wrappers atop them avoid the steep learning curve of the command-line, which requires commands to be typed on the keyboard. By starting a GUI wrapper, users can intuitively interact with, start, stop, and change its working parameters, through graphical icons and visual indicators of a desktop environment, for example. Applications may also provide both interfaces, and when they do the GUI is usually a WIMP wrapper around the command-line version. This is especially common with applications designed for Unix like operating systems. The latter used to be implemented first because it allowed the developers to focus exclusively on their product's functionality without bothering about interface details such as designing icons and placing buttons. Designing programs this way also allows users to run the program in a shell script.

## 7.RESEARCH AND DISCUSSION

The logical process began with data sanctification and processing, followed by exploratory evaluation, and eventually interpretation creation and evaluation. The advanced the quality delicacy on the public examination set, the better the delicacy standing. This can be discovered by comparing each set of rules with the type of all community attacks for unborn vaticination issues by discovering quality connections. This leads to a number of new perceptivity into diagnosing the community assault of each new connection. To give a vaticination interpretation that uses synthetic intelligence to ameliorate over mortal delicacy and give early discovery capabilities.learning technique is useful in

developing prediction models that can helps to network sectors reduce the long process of diagnosis and eradicate any human error.

## 9.FUTURE WORK

The network area must automate the discovery of packet transfer attacks from an eligibility perspective( in real time) grounded on relationship detail.
To automate this process by displaying the vaticination through web software or computer device software. To ameliorate the oils

## 8.CONCLUSION

The analytical process started from data cleaning and processing, missing value, exploratory analysis and finally model building and evaluation. The best accuracy on public test set is higher accuracy score will be find out by comparing each algorithm with type of all network attacks for future prediction results by finding best connections. This brings some of the following insights about diagnose the network attack of each new connection. To presented a prediction model with the aid of artificial intelligence to improve over human accuracy and provide with the scope of early detection. It can be inferred from this model that; area analysis and use of machine similar that they can be used in an Artificial Intelligence environment.

## 10.REFERENCES

[1] Preetish Ranjan and Abhishek Vaish "Apriori viterbri model for prior detection of socio-Technical Attacks in social Network" In 2014 IEEE international conference.

[2] Mohamad Syahir Abdullah;Anazida Zainal;Mohd Aizaini Maarof;Mohamad Nizam Kassim "Cyber-Attack Features for Detecting Cyber Threat Incidents from Online News" In 2018 IEEE international conference

[3] Xiaoyong Yuan, Pan He, Qile Zhu and Xiaolin Li"Adversarial Examples: Attacks and Defenses for Deep Learning" In 2019 IEEE international conference

[4] Wenying Xu and Guoqiang Hu "Distributed Secure Cooperative Control Under Denial-of-Service Attacks From Multiple Adversaries" In 2019 IEEE international conference

[5] Ziwen Sun , Shuguo Zhang" Modeling of Security Risk for Industrial Cyber-Physics System under Cyber-attacks" In 2021 IEEE international conference

[6] Z. Liu, C. Chen, X. Yang, J. Zhou, X. Li, and L. Song, "Heterogeneous graph neural networks for malicious account detection," in Proc. 27th ACM Int. Conf. Inf. Knowl. Manage., Oct. 2018, pp. 2077–2085.

[7] L. Tang and H. Liu, "Relational learning via latent social dimensions," in Proc. 15th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD), 2009, pp. 817–826 .

[8] Y. Zhang et al., "Your style your identity: Leveraging writing and photography styles for drug trafficker identification in darknet markets over attributed heterogeneous information network," in Proc. World Wide Web Conf. (WWW), 2019, pp. 3448–3454.

[9] R . Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, and J. Leskovec, "Graph convolutional neural networks for Web-scale recommender systems," in Proc. 24th ACM SIGKDD Int. Conf. Knowl.Discovery Data Mining, Jul. 2018, pp. 974–983.

[10] W. Fan et al., "Graph neural networks for social recommenda-

tion," in Proc. World Wide Web Conf. (WWW), 2019, pp. 417–426.

[11]  T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in Proc. Int. Conf. Learn. Represent., 2017,pp. 1–14.

[12]  P. Veli˘ckovi´c, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," in Proc. Int.