

CYBER HACKING BREACHES PREDICTION USING MACHINE LEARNING

*Major project report submitted
in partial fulfillment of the requirement for award of the degree of*

**Bachelor of Technology
in
Computer Science & Engineering**

By

**NALLURI KARTHIK 20UECS0659 VTU15337
KOTA MANISH 20UECS0501 VTU17731**

*Under the guidance of
Dr. P.J Beslin pajila M.E ., Ph.D .,
ASSISTANT PROFESSOR - SENIOR GRADE*



**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
SCHOOL OF COMPUTING**

**VEL TECH RANGARAJAN DR. SAGUNTHALA R&D INSTITUTE OF
SCIENCE & TECHNOLOGY**

(Deemed to be University Estd u/s 3 of UGC Act, 1956)

**Accredited by NAAC with A++ Grade
CHENNAI 600 062, TAMILNADU, INDIA**

May, 2024

CYBER HACKING BREACHES PREDICTION USING MACHINE LEARNING

*Major project report submitted
in partial fulfillment of the requirement for award of the degree of*

**Bachelor of Technology
in
Computer Science & Engineering**

By

**NALLURI KARTHIK 20UECS0659 VTU15337
KOTA MANISH 20UECS0501 VTU17731**

*Under the guidance of
Dr. P.J Beslin pajila M.E ., Ph.D .,
ASSISTANT PROFESSOR - SENIOR GRADE*



**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
SCHOOL OF COMPUTING**

**VEL TECH RANGARAJAN DR. SAGUNTHALA R&D INSTITUTE OF
SCIENCE & TECHNOLOGY**

(Deemed to be University Estd u/s 3 of UGC Act, 1956)

**Accredited by NAAC with A++ Grade
CHENNAI 600 062, TAMILNADU, INDIA**

May, 2024

CERTIFICATE

It is certified that the work contained in the project report title “CYBER HACKING BREACHES PREDICTION USING MACHINE LEARNING” by “NALLURI KARTHIK 20UECS0659, KOTA MANISH 20UECS0501, ” has been carried out under my supervision and that this work has not been submitted elsewhere for a degree.

Signature of Supervisor

Computer Science & Engineering

School of Computing

Vel Tech Rangarajan Dr. Sagunthala R&D

Institute of Science & Technology

May, 2024

Signature of Professor In-charge

Computer Science & Engineering

School of Computing

Vel Tech Rangarajan Dr. Sagunthala R&D

Institute of Science & Technology

May, 2024

DECLARATION

We declare that this written submission represents our ideas in our own words and where other's ideas or words have been included, we have adequately cited and referenced the original sources. We also declare that we have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in our submission. We understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

NALLURI KARTHIK

Date: / /

KOTA MANISH

Date: / /

APPROVAL SHEET

This project report entitled “CYBER HACKING BREACHES PREDICTION USING MACHINE LEARNING” by NALLURI KARTHIK 20UECS0659, KOTA MANISH 20UECS0501 is approved for the degree of B.Tech in Computer Science & Engineering.

Examiners

Supervisor

Dr. P.J.Beslin Pajila ,M.E., Ph.D.,
ASSISTANT PROFESSOR.

Date: / /

Place:

ACKNOWLEDGEMENT

We express our deepest gratitude to our respected **Founder Chancellor and President Col. Prof. Dr. R. RANGARAJAN B.E. (EEE), B.E. (MECH), M.S (AUTO),D.Sc., Foundress President Dr. R. SAGUNTHALA RANGARAJAN M.B.B.S.** Chairperson Managing Trustee and Vice President.

We are very much grateful to our beloved **Vice Chancellor Prof. S. SALIVAHANAN**, for providing us with an environment to complete our project successfully.

We record indebtedness to our **Professor & Dean, Department of Computer Science & Engineering, School of Computing, Dr. V. SRINIVASA RAO, M.Tech., Ph.D.,** for immense care and encouragement towards us throughout the course of this project.

We are thankful to our **Head, Department of Computer Science & Engineering, Dr.M.S. MURALI DHAR, M.E., Ph.D.,** for providing immense support in all our endeavors.

We also take this opportunity to express a deep sense of gratitude to our Internal Supervisor **Dr. P.J BESLIN PAJILA M.E ., Ph.D .,** for her cordial support, valuable information and guidance, she helped us in completing this project through various stages.

A special thanks to our **Project Coordinators Mr. V. ASHOK KUMAR, M.Tech., Ms. C. SHYAMALA KUMARI, M.E.,** for their valuable guidance and support throughout the course of the project.

We thank our department faculty, supporting staff and friends for their help and guidance to complete this project.

NALLURI KARTHIK	20UECS0659
KOTA MANISH	20UECS0501

ABSTRACT

The rapid evolution of cyber threats necessitates proactive measures to anticipate and mitigate potential breaches. In this study, we employ machine learning algorithms including Decision Tree, Random Forest, AdaBoost, Logistic Regression, KNN, and SVC to forecast cyber hacking breaches. By analyzing historical data and identifying patterns indicative of potential security breaches, these algorithms enable proactive decision-making and resource allocation for cybersecurity defenses. Through rigorous evaluation and comparison of these algorithms, we aim to determine the most effective approach for predicting cyber hacking breaches. This research contributes to enhancing cybersecurity strategies by providing insights into the predictive capabilities of various machine learning techniques. By leveraging predictive analytics, organizations can strengthen their defenses, detect vulnerabilities, and preemptively thwart cyber threats, ultimately safeguarding sensitive information and preserving operational integrity in an increasingly digital landscape.

Cyber-physical systems (cps) have made significant progress in many dynamic applications due to the integration between physical processes, computational resources, and communication capabilities. However, cyber-attacks are a major threat to these systems. Unlike faults that occurs by accidents cyber-physical systems, cyber-attacks occur intelligently and stealthy. Some of these attacks which are called deception attacks, inject false data from sensors or controllers, and also by compromising with some cyber components, corrupt data, or enter misinformation into the system. If the system is unaware of the existence of these attacks, it won't be able to detect them, and performance may be disrupted or disabled altogether. Therefore, it is necessary to adapt algorithms to identify these types of attacks in these systems. It should be noted that the data generated in these systems is produced in very large number, with so much variety, and high speed, so it is important to use machine learning algorithms to facilitate the analysis and evaluation of data and to identify hidden patterns.

Keywords: Decision Tree, Random Forest, AdaBoost, Logistic Regression, KNN, SVC.

LIST OF FIGURES

4.1	Architecture Diagram	14
4.2	Data Flow Diagram	15
4.3	Class Diagram	16
4.4	Sequence Diagram	17
4.5	Activity Diagram	18
4.6	Collabration Diagram	19
5.1	Dataset Input Diagram	25
5.2	Graph Representing Types of Breaches in Output Diagram . . .	26
5.3	Input for Unit Testing	29
5.4	Result of Black Box Testing	29
5.5	Test Image	30
6.1	Random Forest output Metrics	35
8.1	Offer letter of Nalluri Karthik	39
8.2	Offer letter of Kota Manish	40
9.1	Plagarism Report	42
10.1	Poster Presentation	49

LIST OF TABLES

6.1	Comparing Efficiency Metrics of Existing System and Proposed System	32
-----	--	-----------

LIST OF ACRONYMS AND ABBREVIATIONS

CAM	Comprehensive Assessment Mode
CSRM	Cyber Security Risk Management
CPS	Cyber Physical Systems
DOS	Denial Of Service
GDPR	General Data Protection Regulation
HIPAA	Health Insurance Portability and Ac Countability Act
ISMS	International Standard for Information Security Management Systems
IDS	Intrusion Detection System
KNN	K-Nearest Neighbor
NCS	Networked Control System
ROC	RATE Of Investment
SOC	Security Operations Centers
SIEM	Security Information And Event Management
SVM	Support Vector Machine
TTP	Tactic, Technique and Procedure
VCBD	VERIS Community Dataset

TABLE OF CONTENTS

	Page.No
ABSTRACT	v
LIST OF FIGURES	vi
LIST OF TABLES	vii
LIST OF ACRONYMS AND ABBREVIATIONS	viii
1 INTRODUCTION	1
1.1 Introduction	1
1.2 Aim of the Project	2
1.3 Project Domain	2
1.4 Scope of the Project	3
2 LITERATURE REVIEW	4
3 PROJECT DESCRIPTION	8
3.1 Existing System	8
3.1.1 Disadvantages	8
3.2 Proposed System	9
3.2.1 Advantages	9
3.3 Feasibility Study	10
3.3.1 Economic Feasibility	10
3.3.2 Technical Feasibility	11
3.3.3 Social Feasibility	11
3.4 System Specification	12
3.4.1 Hardware Specification	12
3.4.2 Software Specification	12
3.4.3 Standards and Policies	12
4 METHODOLOGY	14
4.1 General Architecture	14

4.2	Design Phase	15
4.2.1	Data Flow Diagram	15
4.2.2	Class Diagram	16
4.2.3	Sequence Diagram	17
4.2.4	Activity Diagram	18
4.2.5	Collabration Diagram	19
4.3	Algorithm & Pseudo Code	20
4.3.1	Random Forest Algorithm	20
4.3.2	Pseudo Code	21
4.4	Module Description:	23
4.4.1	Data Preprocessing	23
4.4.2	Feature Selection	23
4.4.3	Model Training	23
4.4.4	Prediction	23
4.4.5	Monitoring and Maintenance	24
4.4.6	Visualization	24
5	IMPLEMENTATION AND TESTING	25
5.1	Input and Output	25
5.1.1	Input Design	25
5.1.2	Output Design	26
5.2	Testing	27
5.3	Types of Testing	27
5.3.1	Unit Testing	27
5.3.2	Black Box Testing	29
5.3.3	Test Result	30
6	RESULTS AND DISCUSSIONS	31
6.1	Efficiency of the Proposed System	31
6.2	Comparison of Existing and Proposed System	32
6.3	Sample Code	33
7	CONCLUSION AND FUTURE ENHANCEMENTS	36
7.1	Conclusion	36
7.2	Future Enhancements	37

8	INDUSTRY DETAILS	38
8.1	Industry name	38
8.1.1	Duration of Internship (From Date - To Date)	38
8.1.2	Duration of Internship in months	38
8.1.3	Industry Address	38
8.2	Internship Offer Letter	39
9	PLAGIARISM REPORT	42
10	SOURCE CODE & POSTER PRESENTATION	43
10.1	Source Code	43
10.2	Poster Presentation	49
	References	49

Chapter 1

INTRODUCTION

1.1 Introduction

The advent of the digital age has ushered in unprecedented opportunities for connectivity and innovation. However, it has also given rise to a growing menace—cyber hacking breaches. In recent years, cyberattacks have become increasingly sophisticated and devastating, posing a significant threat to individuals, businesses, and even nations.

This project aims to comprehensively investigate and analyze cyber hacking breach that occurred in the past year. We will delve into the methods, motivations, and impacts of these breaches to gain a deeper understanding of the evolving landscape of cyber threats. By examining a range of high-profile cases, we intend to identify common vulnerabilities and attack vectors.

Cyber hacking breaches are driven by various motivations, including financial gain, cyber espionage, hacktivism, and personal curiosity. Criminals seek profit through data theft and fraud, while state actors engage in espionage. Hacktivists pursue political or social causes, and some hackers explore vulnerabilities for personal notoriety, highlighting diverse threats in cybersecurity. Additionally, we will explore the ethical, legal, and regulatory aspects surrounding cyberattacks and data breaches.

This research not only serves as a valuable resource for cybersecurity professionals but also contributes to raising awareness among individuals and organizations about the importance of robust digital security. By shedding light on the ever-evolving world of cyber hacking breaches, we aim to empower stakeholders to fortify their defenses and safeguard their digital assets in an increasingly interconnected world.

1.2 Aim of the Project

To illustrate the escalating pace of cyber threats necessitates preemptive measures to anticipate and counter potential breaches. This study aims to develop predictive models using machine learning algorithms such as Decision Tree, Random Forest, AdaBoost and KNN to forecast cyber hacking breaches. By leveraging historical data and identifying patterns indicative of security risks, the research seeks to enhance cybersecurity strategies and enable proactive defense mechanisms against evolving cyber threats.

1.3 Project Domain

Using actual cybercrime data, the initial phase is to predict a cybercrime strategy, and the accuracy results are then compared. The second is to examine if the information at hand can be used to predict cybercrime perpetrators. It is employed to hide a system's data. Information theft is caused by sensitive and highly confidential data as well as poor management. The hackers' techniques might be found in two different methods. One is to move through with legal action, get in touch with the victim, and let them know about the violations. The organizations should be aware of the sorts, trends, and patterns of assaults for the purpose of enabling them to monitor the system.

A study on the consequences of these kinds of attacks in an effort for managing the prevention of occurring the breaches. And provide comprehensive study of the breaches that have occurred by the various organizations and financial effect. Because of improvements in information technology, prices for memory and storage devices, and the expansion of the digital economy, businesses and governmental organizations now acquire more data every day.

Businesses and organizations have the threat of data attacks because of the collecting of personal data on their computers. Computer networks are used in manufacturing, healthcare, research. This information is transferring every second through network. These attacks are used for profit and destroy the important information and use that for own need which rises the risk of data. Hybrid based detection is used to detect the high false positive rates and low false positive rates. Anomaly based detection analyse the behaviour of traffic, where signature based detection has the previous attacks records and able to detect the possibilities.

1.4 Scope of the Project

The project scope encompasses assessing existing cybersecurity measures, identifying vulnerabilities, implementing security enhancements, conducting user awareness training, and establishing continuous monitoring protocols. It aims to fortify digital systems and data protection against cyber hacking breaches. Develop a robust machine learning model to analyze historical cyber hacking incidents, aiming to predict and identify potential future breaches based on patterns and trends.

Enhance cybersecurity preparedness by leveraging predictive analytics, enabling organizations to proactively address vulnerabilities and mitigate risks associated with cyber hacking threats. Evaluate the effectiveness of the developed machine learning system in real-time breach prediction, providing a valuable tool for preemptive cybersecurity measures and reducing the impact of potential cyber breaches.

Chapter 2

LITERATURE REVIEW

[1] Kwon et al.(2013), explored the security issue in the state estimation problem is investigated for a networked control system (NCS). The communication channels between the sensors and the remote estimator in the NCS are vulnerable to attacks from malicious adversaries. The false data injection attacks are considered. The aim of this paper to find the so-called insecurity conditions under which the estimation system is insecure in the sense that there exist malicious attacks that can bypass the anomaly detector but still lead to unbounded estimation errors. In particular, a new necessary and sufficient condition for the insecurity is derived in the case that all communication channels are compromised by the adversary .

[2] Pajic et al.(2017), investigated the significant increase in the number of security-related incidents in control systems. These include high-profile attacks in a wide range of application domains, from attacks on critical infrastructure, as in the case of the Maroochy Water breach , and industrial systems (such as the StuxNet virus attack on an industrial supervisory control and data acquisition system , and the German Steel Mill cyberattack) , to attacks on modern vehicles . Even high-assurance military systems were shown to be vulnerable to attacks, as illustrated in the highly publicized downing of the RQ-170 Sentinel U.S. drone . These incidents have greatly raised awareness of the need for security in cyberphysical systems (CPSs), which feature tight coupling of computation and communication substrates with sensing and actuation components.

[3] Sheng et al.(2012), explored the Embedded computational resources in autonomous robotic vehicles are becoming more abundant and have enabled improved operational effectiveness of cooperative robotic systems in civilian and military applications. Compared to autonomous robotic vehicles that operate single tasks, cooperative teamwork has greater efficiency and operational capability. Multirobotic vehicle systems have many potential applications, such as platooning of vehicles in urban transportation, the operation of the multiple robots, autonomous underwater vehi

cles, and formation of aircrafts in military affairs . The study of group behaviors for multirobot systems is the main objective of the work. Group cooperative behavior signifies that individuals in the group share a common objective and action according to the interest of the whole group. Group cooperation can be efficient if individuals in the group coordinate their actions well.

[4] Zeng et al.(2014),discussed the problem of reaching a consensus among all the agents in the networked control systems (NCS) in the presence of misbehaving agents. A reputation-based resilient distributed control algorithm is first proposed for the leader-follower consensus network. The proposed algorithm embeds a resilience mechanism that includes four phases (detection, mitigation, identification, and up date), into the control process in a distributed manner. At each phase, every agent only uses local and one-hop neighbors' information to identify and isolate the mis behaving agents, and even compensate their effect on the system.

[5] Hongtao et al.(2017), focused on resilient control of networked control systems (NCSs) under the denial of service (DoS) attacks which is characterized by a Markov process. Firstly, the packets dropout are modeled as Markov process according to the gamebetweenattack strategies and defense strategies. Then, an NCS under such game results is modeled as a Markovian jump linear system and four theorems are proved for the system stability analysis and controller design. Finally, a numerical example is used to illustrative the application of these theorems. Networked control systems (NCSs) have received an increasing attention in the past decades. Now, NCSs have been widely applied in industrial processes, electric power networks, intelligent transportation and so on. With the growing of the NCSs, network, as a critical element in an NCS, is vulnerable to cyber threats which can menace the control systems.

[6] M. Eling et al.(2016), defined the literature by Eling and Schnell (2016) explores the intricate landscape of cyber risk and cyber risk insurance, providing valuable insights into this rapidly evolving field. The authors delve into the realization of cyber risk, offering a comprehensive analysis that spans issues such as identification, assessment, and mitigation strategies. The study, featured in the Journal of Risk Finance, is particularly noteworthy for its emphasis on the dynamic nature of cyber threats and the corresponding challenges in developing effective insurance mecha

nisms. Eling and Schnell contribute significantly to the scholarly discourse by addressing key facets of cyber risk management, ultimately shedding light on the complexities associated with safeguarding organizations against cyber threats.

[7] Mandal et al.(2020), focused considering the different aspects of social events, responses and their relations to further improve the classification of the social sentiment. The proposed method covers not only the response due to major social events but also predicting and generating alert for situations of significant social importance. The approach has made use of Twitter datasets and performed aspect based sentiment analysis on the obtained text data. It is shown to outperform the state-of-the-art methods.

[8] Poyraz et al. (2020) investigated various factors that can affect the monetary impact of data breaches on companies. This paper introduces a model for the total cost of a mega data breach based on a data set created from multiple sources that categorises stolen data for U.S. residents as personally identifiable information (PII) and sensitive personally identifiable information (SPII). They use a rigorous step wise regression analysis that includes polynomial and factorial multilevel effects of the independent variables. There are three significant findings. First, our model finds a significant relation between total data breach cost and revenue, the total amount of PII and SPII, and class action lawsuits.

[9] Kure et al.(2021), presented an effective cybersecurity risk management (CSRM) practice using assets criticality, predication of risk types and evaluating the effectiveness of existing controls. They follow a number of techniques for the proposed unified approach including fuzzy set theory for the asset criticality, machine learning classifiers for the risk predication and comprehensive assessment model (CAM) for evaluating the effectiveness of the existing controls. The proposed approach considers relevant CSRM concepts such as asset, threat actor, attack pattern, tactic, technique and procedure (TTP), and controls and maps these concepts with the VERIS community dataset (VCDB) features for the risk predication. The experimental results reveal that using the fuzzy set theory in assessing assets criticality supports stakeholder for an effective risk management practice ..

[10] R. R. Subramanian et al.(2021), designed a model by using machine learning to defend a website from security breaches. The primary aim of this research work is to create a machine learning model, which trains in Realtime and monitors the website or a system and trains from the state-of-art attacks. The proposed model has created a web application using Django, which takes the data from multiple sources such as Amazon, Flipkart, Snapdeal, and Shop clues, which shows the data that is safe to obtain from the website. Then, the data will be sorted on our page and then it will be made secured and illegal for the external people to access the data from our website and the proposed model will monitor the website 24/7. The model is trained daily and it generates predictions from the several of datasets available and from the previous state-of-the-art attacks. This model will be trained from the existing datasets and the history of attacks and breaches on our website.

Chapter 3

PROJECT DESCRIPTION

3.1 Existing System

In the existing system, implementation of machine learning algorithms is bit complex to build due to the lack of information about the data visualization. Mathematical calculations are used in existing system for SVM and Logistic Regression model building this may take the lot of time and complexity. To overcome all this, And use machine learning packages available in the scikit-learn library.

The existing system collects data from various sources such as network logs, system event logs, firewall logs, intrusion detection system (IDS) alerts, and other cyber security telemetry data sources. Data preprocessing techniques are applied to clean, normalize, and transform the raw data into a format suitable for machine learning analysis. This may involve handling missing values, removing noise, and encoding categorical variables.

3.1.1 Disadvantages

1. Limited Accuracy: Machine learning predictions for cyber hacking breaches may suffer from limited accuracy due to the evolving nature of hacking techniques and strategies.
2. Over-reliance on Historical Data: Dependence on historical data for machine learning models can lead to a disadvantage as it may not fully capture emerging and novel hacking patterns.
3. Adversarial Attacks: Machine learning models are vulnerable to adversarial attacks, where malicious actors intentionally manipulate input data to deceive prediction systems and compromise their effectiveness.
4. Resource Intensiveness: Implementing and maintaining machine learning solutions for cyber breach prediction can be resource-intensive, requiring significant computational power and expertise.

3.2 Proposed System

This study proposes a predictive system for cyber hacking breaches using machine learning algorithms such as Decision Tree, Random Forest, AdaBoost and KNN. Leveraging historical data analysis and pattern recognition, the system aims to forecast potential security breaches, enabling proactive decision making and resource allocation for cybersecurity defenses. Through rigorous evaluation and comparison of these algorithms, the system seeks to identify the most effective approach for predicting cyber threats, contributing to enhanced cybersecurity strategies and preemptive threat mitigation in digital environments.

Random Forest is a powerful machine learning algorithm that assembles multiple decision trees to make accurate predictions. Each tree is trained on a subset of data and votes on the final prediction, resulting in improved accuracy and robustness. It mitigates overfitting and handles complex relationships in data by averaging predictions from different trees. Random Forest is versatile, handling classification and regression tasks effectively.

3.2.1 Advantages

1. **Early Detection:** Machine learning enables early identification of potential cyber hacking breaches, allowing proactive measures to be implemented before significant damage occurs.
2. **Adaptive Analysis:** ML algorithms adapt to evolving cyber threats, providing continuous and dynamic analysis for more accurate prediction of hacking attempts.
3. **Pattern Recognition:** Machine learning excels in recognizing patterns within vast datasets, enhancing the ability to detect subtle indicators of potential cyber breaches.
4. **Improved Incident Response:** ML-driven predictions empower organizations with timely insights, facilitating faster and more effective incident response strategies.
5. **Enhanced Risk Mitigation:** By leveraging machine learning, organizations can enhance their risk mitigation efforts, preemptively addressing vulnerabilities and minimizing the impact of cyber hacking breaches.

3.3 Feasibility Study

Feasibility study is an important step in the development of any project, These are the requirements that the end user specifically demands as basic facilities that the system should offer. All these functionalities need to be necessarily incorporated into the system as a part of the contract. These are represented or stated in the form of input to be given to the system, the operation performed and the output expected. They are basically the requirements stated by the user which one can see directly in the final product, unlike the non-functional requirements. Examples of functional requirements:

- 1) Authentication of user whenever he/she logs into the system
- 2) System shutdown in case of a cyber-attack
- 3) A verification email is sent to user whenever he/she register for the first time on some software system. It involves evaluating the viability of the project in terms of its technical, economic, and social feasibility.

3.3.1 Economic Feasibility

An economic feasibility study involves a thorough assessment of its financial viability and potential returns on investment (ROI). This evaluation begins by estimating both the initial investment required and the recurring expenses associated with data acquisition, infrastructure setup, and personnel. Subsequently, potential revenue streams are identified, including cost savings resulting from preventing cyber hacking breaches and potential revenue generation from offering cybersecurity services or licensing the predictive model. By comparing the expected benefits to the costs and conducting ROI analysis using financial metrics like NPV, IRR, and payback period, the project's economic performance is evaluated over time. Additionally, a comprehensive risk analysis is conducted to identify and mitigate potential risks that could impact the project's economic feasibility. Based on these findings, conclusions and recommendations are made regarding whether to proceed with the project, modify its scope, or explore alternative approaches. Through this economic feasibility study, stakeholders can make well-informed decisions about investing in the project to predict cyber hacking breaches using machine learning.

3.3.2 Technical Feasibility

Achieving technical feasibility the project's practicality from a technological standpoint. It begins by evaluating the availability and quality of relevant datasets containing historical cybersecurity incidents, network logs, and system event logs. Additionally, the technical requirements for infrastructure, hardware, and software needed for data storage, preprocessing, model training, and deployment are analyzed. Consideration is given to scalability, compatibility with existing systems, and the computational resources required. The feasibility of implementing machine learning algorithms for breach prediction, such as logistic regression, random forests, or neural networks, is also examined. Furthermore, the project's adherence to legal and ethical guidelines regarding data privacy, security, and bias mitigation is assessed. By conducting a comprehensive technical feasibility study, stakeholders can determine whether the necessary resources, technology, and expertise are available to successfully implement the project and achieve its objectives.

3.3.3 Social Feasibility

A social feasibility study involves evaluating its acceptance, impact, and implications within the broader social context. It begins by assessing the societal need for improved cybersecurity measures to protect individuals, organizations, and critical infrastructure from cyber threats. Stakeholder perspectives, including those of end users, cybersecurity experts, policymakers, and the general public, are considered to understand their attitudes, concerns, and expectations regarding the project. Additionally, potential social benefits, such as enhanced trust and confidence in online security, are identified, along with any potential risks or negative consequences, such as privacy concerns or unintended biases in the machine learning models. The project's alignment with societal values, ethical principles, and legal regulations is evaluated to ensure that it promotes fairness, transparency, and accountability in cybersecurity practices. By engaging stakeholders in dialogue, addressing their needs and concerns, and fostering collaboration and trust, the project can contribute to a more resilient and secure digital environment that benefits society as a whole. Through this social feasibility study, stakeholders can assess the project's potential social impact and its ability to garner support and acceptance from the community.

3.4 System Specification

3.4.1 Hardware Specification

- Processor: I5 and above generation
- RAM: 4GB and above

3.4.2 Software Specification

- Operating System: Windows 7/8/10
- Server side Script: HTML, CSS, Bootstrap JS
- Programming Language: Python
- Libraries: Flask, Pandas, Mysql.connector, Os, Smtplib, Numpy
- IDE/Workbench: PyCharm
- Technology: Python 3.6+
- Server Deployment: Xampp Server
- Database: MySQL

3.4.3 Standards and Policies

Standards and policies to ensure ethical, legal, and responsible implementation. One critical aspect is compliance with data privacy regulations such as the General Data Protection Regulation (GDPR) and the Health Insurance Portability and Accountability Act (HIPAA), which govern the collection, processing, and protection of personal and sensitive data. Additionally, adherence to industry-specific standards like ISO/IEC 27001 for information security management systems is essential to ensure the confidentiality, integrity, and availability of data throughout the project lifecycle. Ethical guidelines for AI and machine learning, such as those outlined by organizations like the IEEE and the Partnership on AI, provide principles for fairness, transparency, accountability, and privacy in algorithmic decision-making.

- Standards Organization: International standard for information security management systems (ISMS).
- Division Name: Security Management
- Section Name: Cybersecurity Framework
- Designator of Legally binding Document: ISO/IEC 27001 .
- Jupyter Notebook-Standards Used: ISO3166-1:2018
- Python-Standard Used: ISO6160:1979

Anaconda Prompt: Anaconda prompt is a type of command line interface which explicitly deals with the ML modules and the navigator is available in all Windows, Linux, and MacOS. The Anaconda prompt has many integrated development environments that make coding easier. The user interface can also be implemented in Python.

Standard Used: ISO/IEC 27001

Jupyter: It's like an open-source web application that allows us to share and create documents that contain live code, equations, visualizations, and narrative text. It can be used for data cleaning and transformation, numerical simulation, statistical modeling, data visualization, and machine learning.

Standard Used: ISO/IEC 27001

Chapter 4

METHODOLOGY

4.1 General Architecture

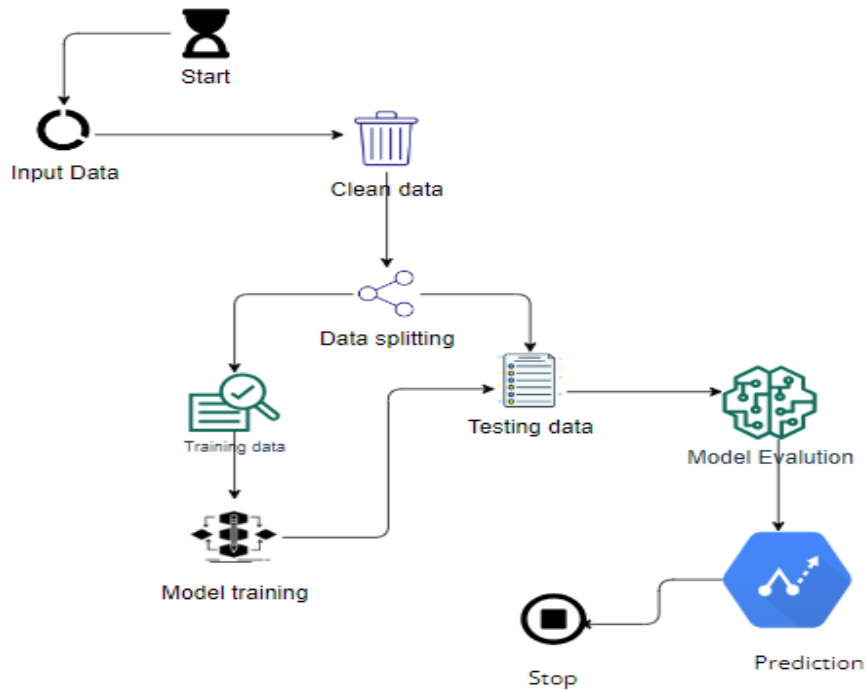


Figure 4.1: Architecture Diagram

Figure 4.1 outlines the system's overall workflow, starting with data acquisition illustrating the end-to-end process of predicting cyber hacking breaches using machine learning, from data collection and preprocessing to model deployment and continuous improvement. It emphasizes the integration of machine learning capabilities into existing cybersecurity infrastructure to enhance threat detection and response capabilities.

4.2 Design Phase

4.2.1 Data Flow Diagram

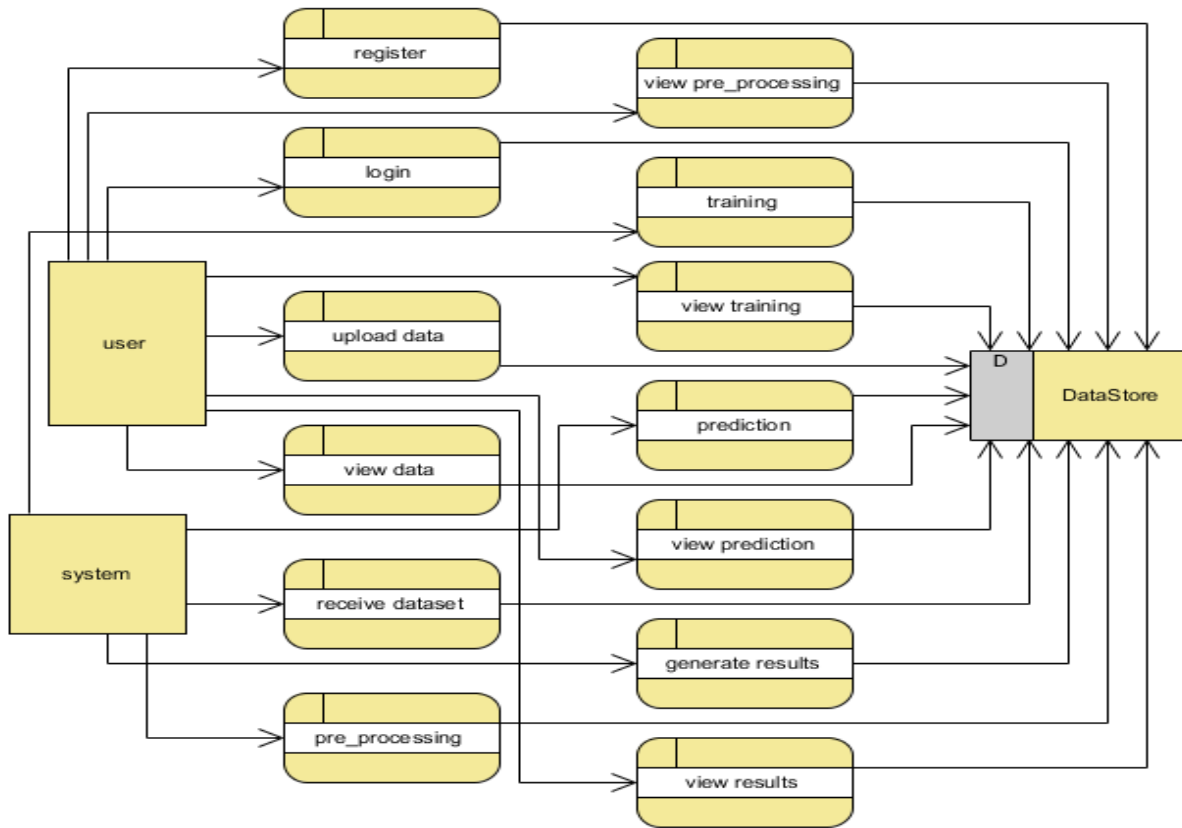


Figure 4.2: Data Flow Diagram

Figure 4.2 shows the flow of data, delineating several sequential steps. Initially, the system receives input in the form of a CSV file comprising data collection and pre-processing to model deployment and continuous improvement. It highlights the key stages and interactions involved in predicting cyber hacking breaches using machine learning. Raw data from the sources undergo preprocessing to clean, transform, and prepare it for analysis. This step involves tasks such as handling missing values, removing noise, normalizing features, and encoding categorical variables. Feature engineering extracts relevant features from the preprocessed data to represent patterns indicative of potential hacking breaches. These features could include network traffic features, system activity features, user behavior features. Deployed models are integrated into the existing cybersecurity infrastructure for real-time or batch processing of data streams. Integration with security information and event management (SIEM) systems, security operations centers (SOCs), or threat intelligence platforms facilitates automated alerting and incident response based on model predictions.

4.2.2 Class Diagram

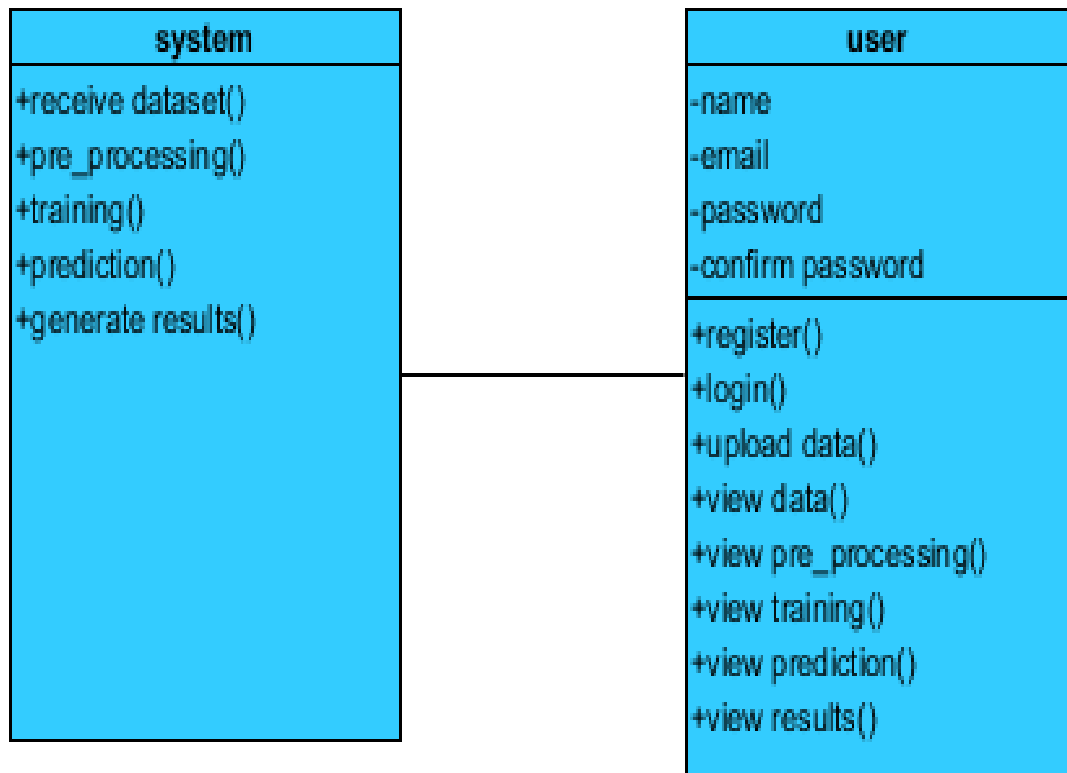


Figure 4.3: Class Diagram

Figure 4.3 illustrates the Class diagram a simplified view of the project's architecture, focusing on the key components responsible for data processing, machine learning, evaluation, visualization, feedback loop management, integration, and compliance. Actual implementations may involve more classes and additional complexities based on specific requirements and design considerations. Static structure diagram that describes the structure of a system by showing the system's classes, their attributes, operations (or methods), and the relationships among the classes. It explains which class contains information DataProcessingSystem class represents the system responsible for preprocessing data and extracting features from raw data sources. MachineLearningModel class encapsulates the machine learning model, including training, prediction, and model storage. ModelEvaluation class handles the evaluation of the machine learning model's performance using various metrics. DataVisualization class provides methods for visualizing the results of model predictions and evaluation. Feedback Loop class manages the feedback loop for monitoring model performance and updating the model with new data.

4.2.3 Sequence Diagram

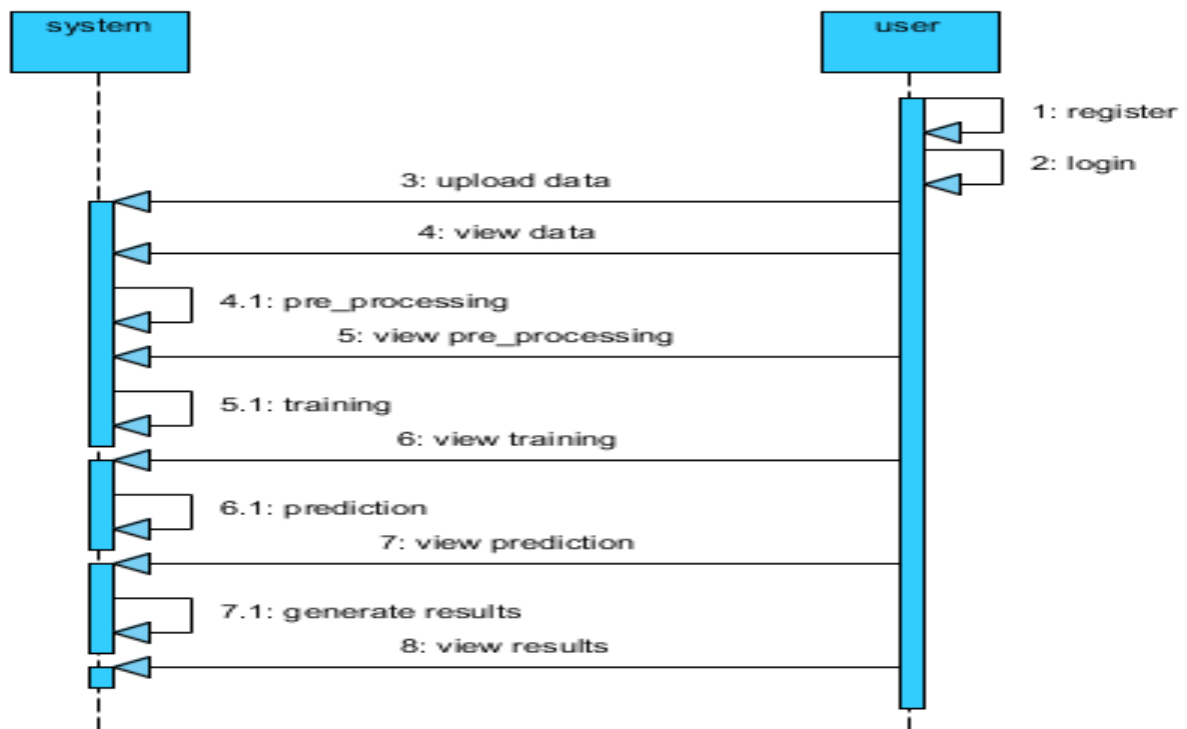


Figure 4.4: Sequence Diagram

Figure 4.4 illustrates the flow of interactions between different components of the project, including data retrieval, preprocessing, feature engineering, model prediction, evaluation, visualization, and feedback loop management. It highlights the key steps involved in predicting cyber hacking breaches using machine learning and the interactions between these steps. Actual implementations may involve more complex interactions and additional steps based on specific requirements and design considerations. It is a construct of a Message Sequence Chart. Sequence diagrams are sometimes called event diagrams, event scenarios, and timing diagrams

4.2.4 Activity Diagram

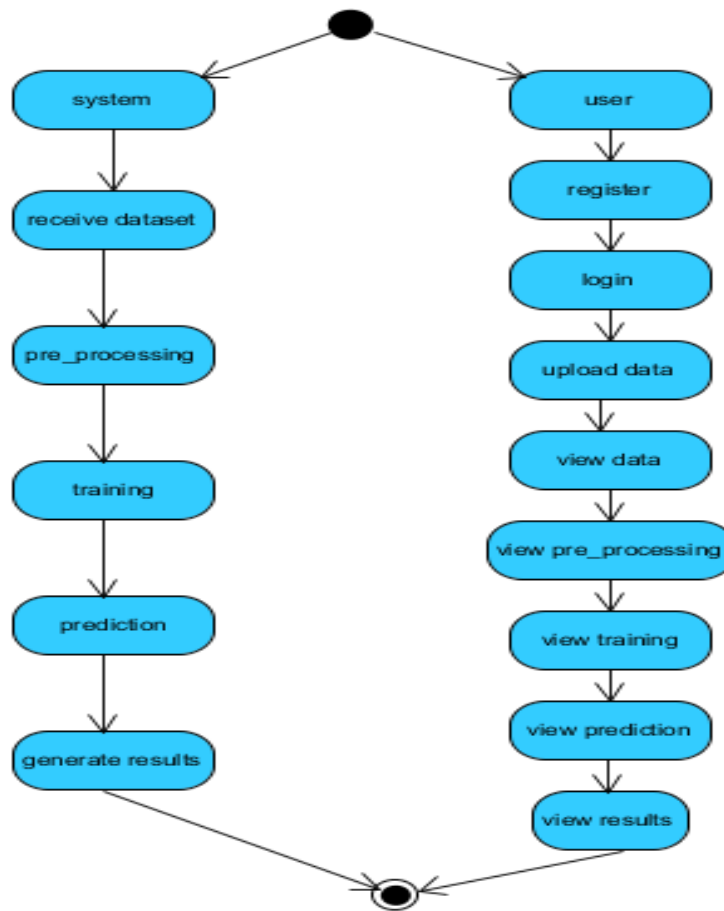


Figure 4.5: Activity Diagram

Figure 4.5 shows the overall flow of control initial phase involves gathering information about the target system or network. Machine learning may be utilized to analyze vast amounts of data quickly and identify potential vulnerabilities or weaknesses. hackers extract sensitive data from the compromised system or network. Machine learning techniques can help obfuscate data transfers, making it harder for security systems to detect and block the unauthorized transfer of data. to cover their tracks to evade detection and attribution. Machine learning can be employed to generate fake or misleading activity, modify logs, or employ other deception techniques

4.2.5 Collaboration Diagram

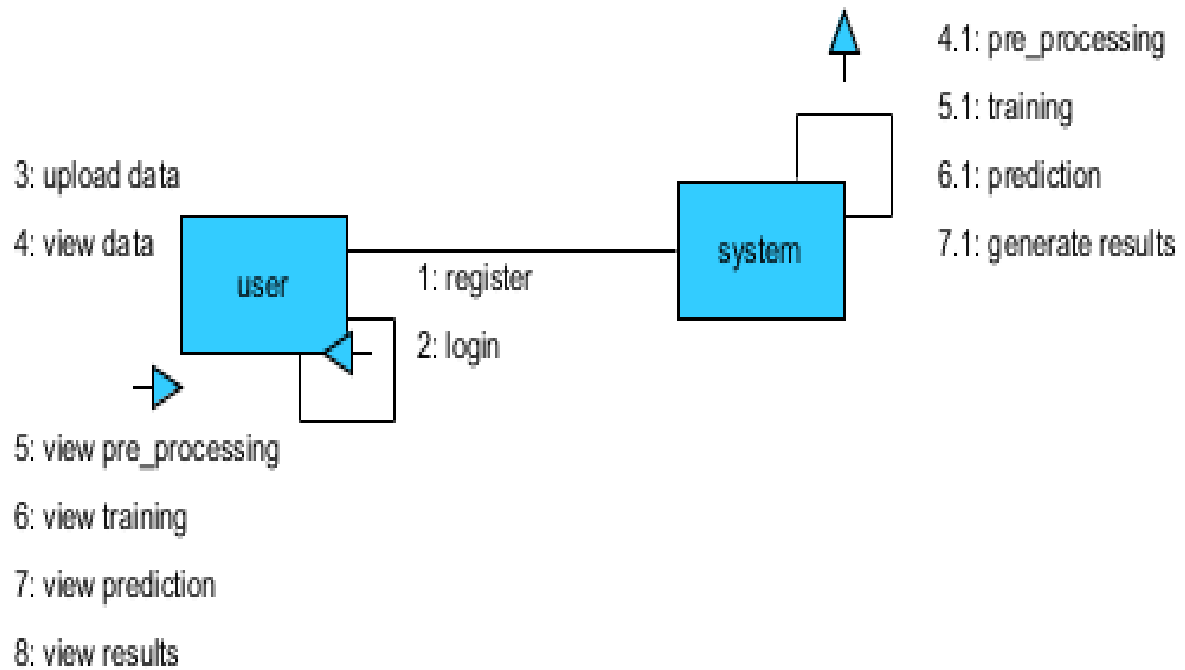


Figure 4.6: Collaboration Diagram

Figure 4.6 shows that the method call sequence is indicated by some numbering technique as shown above. The number indicates how the methods are called one after another. We have taken the same order management system to describe the collaboration diagram. The method calls are similar to that of a sequence diagram. But the difference is that the sequence diagram does not describe the object organization whereas the collaboration diagram shows the object organization.

4.3 Algorithm & Pseudo Code

4.3.1 Random Forest Algorithm

1. Load the necessary libraries such as numpy, pandas, Tensorflow, Keras, sklearn, and matplotlib.
2. Collect historical data on cyber hacking incidents, including network logs, system event logs, firewall logs, intrusion detection system (IDS) alerts, etc.
3. Calculate technical indicators such as moving averages and exponential moving averages using the pandas rolling and ewm functions.
4. Determine the scope of the project, including the types of breaches to predict and the data available.
5. Ensure the data is representative of the problem and contains features that could be useful for prediction.
6. Perform dimensionality reduction techniques if needed (e.g., PCA) to reduce computational complexity.
7. Choose appropriate machine learning algorithms for the problem (e.g., classification algorithms such as logistic regression, decision trees, random forests, or more advanced methods like neural networks).
8. Experiment with different algorithms and evaluate their performance using appropriate metrics (e.g., accuracy, precision, recall, F1-score).
9. Tune hyperparameters using techniques like grid search, random search, or Bayesian optimization to improve model performance.
10. Evaluate the trained models on the validation dataset to assess their generalization performance.
11. Deploy the trained model into production environment, making sure it can handle real-time or batch prediction requests.

4.3.2 Pseudo Code

```
1 import random_forest_library
2
3 function cyberAttackUsingRandomForest(target):
4     // Phase 1: Reconnaissance
5     gatheredInfo = reconnaissance(target)
6
7     // Phase 2: Exploitation
8     vulnerabilities = assessVulnerabilities(gatheredInfo)
9     exploit(vulnerabilities)
10
11    // Phase 3: Exfiltration
12    sensitiveData = exfiltrateData()
13
14    // Phase 4: Evasion/Covering Tracks
15    coverTracks()
16    destroyEvidence()
17
18    return sensitiveData
19
20 function reconnaissance(target):
21    // Gather information about the target using various sources
22    gatheredInfo = gatherInformation(target)
23    // Apply Random Forest algorithm to analyze and prioritize gathered information
24    analyzedInfo = random_forest_analyze(gatheredInfo)
25
26    return analyzedInfo
27
28 function assessVulnerabilities(gatheredInfo):
29    // Use Random Forest algorithm to identify potential vulnerabilities
30    vulnerabilities = random_forest_identifyVulnerabilities(gatheredInfo)
31    // Prioritize vulnerabilities based on potential impact and exploitability
32    prioritizedVulnerabilities = prioritize(vulnerabilities)
33
34    return prioritizedVulnerabilities
35
36 function exploit(vulnerabilities):
37    // Exploit identified vulnerabilities to gain unauthorized access
38    for each vulnerability in vulnerabilities:
39        if exploitSuccessful(vulnerability):
40            establishAccess()
41
42 function exfiltrateData():
43    // Extract sensitive data from the compromised system or network
44    data = extractData()
45    // Use Random Forest algorithm to obfuscate data exfiltration
46    obfuscateData(data)
47
48    return data
```

```

49
50 function coverTracks():
51     // Employ Random Forest algorithm to generate fake or misleading activity
52     generateFakeActivity()
53     // Modify logs and timestamps to hide traces of the attack
54     modifyLogs()
55
56 function destroyEvidence():
57     // Remove any remaining evidence of the attack
58     eraseLogs()
59     // Use Random Forest algorithm to predict potential forensic analysis techniques and counter
60     // them
61     counterForensicAnalysis()
62     def reconnaissance(target):
63         # Gather information about the target using various sources
64         gathered_info = gather_information(target)
65         # Apply Random Forest algorithm to analyze and prioritize gathered information
66         analyzed_info = random_forest_library.analyze(gathered_info)
67
68         return analyzed_info
69
70 def assess_vulnerabilities(gathered_info):
71     # Use Random Forest algorithm to identify potential vulnerabilities
72     vulnerabilities = random_forest_library.identify_vulnerabilities(gathered_info)
73     # Prioritize vulnerabilities based on potential impact and exploitability
74     prioritized_vulnerabilities = prioritize(vulnerabilities)
75
76     return prioritized_vulnerabilities
77
78 def exploit(vulnerabilities):
79     # Exploit identified vulnerabilities to gain unauthorized access
80     for vulnerability in vulnerabilities:
81         if exploit_successful(vulnerability):
82             establish_access()
83
84 def exfiltrate_data():
85     # Extract sensitive data from the compromised system or network
86     data = extract_data()
87     # Use Random Forest algorithm to obfuscate data exfiltration
88     obfuscate_data(data)
89
90     return data
91
92 def cover_tracks():
93     # Employ Random Forest algorithm to generate fake or misleading activity
94     generate_fake_activity()
95     # Modify logs and timestamps to hide traces of the attack
96     modify_logs()

```

4.4 Module Description:

4.4.1 Data Preprocessing

This module provides a structured approach to processing the data required for training a machine learning model to predict cyber hacking breaches. Adjustments may be needed based on the specific characteristics of your dataset and project requirements. Splits the preprocessed data into training and testing sets using train tests-split from scikit-learn.

CSV file dataset-cyberhacking.dataset

4.4.2 Feature Selection

This module provides a simple feature selection method using the chi-squared test. You can adjust the scoring function and feature selection method based on the characteristics of your data and specific project requirements. Uses the SelectKBest method from scikit-learn with the chi-squared test as the scoring function to select the top num-features features.

4.4.3 Model Training

In this segment, the trained Trains different machine learning models using the preprocessed data. Includes tasks such as hyperparameter tuning and model evaluation using cross-validation. Choose a set of candidate machine learning algorithms suitable for the problem, such as logistic regression, decision trees, random forests, support vector machines (SVM), or neural networks.

4.4.4 Prediction

In this segment, the trained model(s) to make predictions on the preprocessed new data. Obtain predicted probabilities or class labels depending on the problem (e.g., binary classification of hacking breaches). Format the predictions into a suitable format for further analysis or presentation. Depending on the application, this may involve storing predictions in a database, generating reports, or displaying them in a user interface.

4.4.5 Monitoring and Maintenance

Monitors the deployed model's performance in the production environment. Handles model retraining and updates to keep the model up-to-date with new data and evolving threats. Trigger model retraining based on predefined criteria such as a decrease in model performance or significant data drift. Retrain the model periodically with new data to keep it up-to-date and effective. Use techniques like incremental learning to update the model without retraining from scratch.

4.4.6 Visualization

This component generates visualizations to explore the data, understand feature importance, and present model evaluation results. Helps in communicating insights and findings to stakeholders effectively. Use techniques such as bar charts, heatmaps, or permutation feature importance to rank and display the importance of each feature in predicting cyber hacking breaches. Use techniques such as bar charts, heatmaps, or permutation feature importance to rank and display the importance of each feature in predicting cyber hacking breaches.

Chapter 5

IMPLEMENTATION AND TESTING

5.1 Input and Output

5.1.1 Input Design

In an information system, input is the raw data that is processed to produce output. During the input design, the developers must consider the input devices such as PC, MICR, OMR, etc. Therefore, the quality of system input determines the quality of system output. Well-designed input forms and screens.

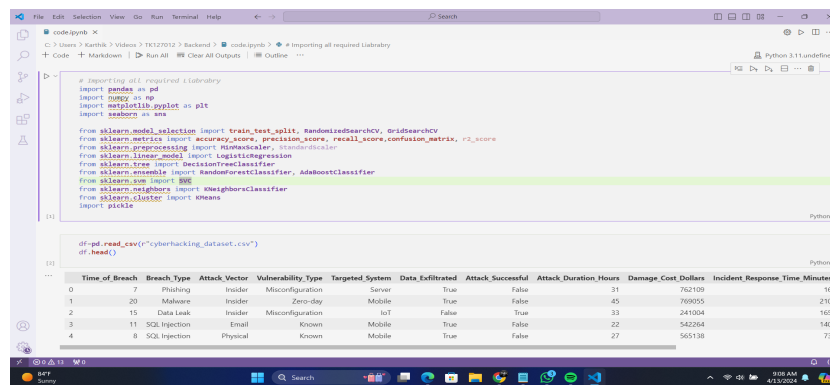


Figure 5.1: Dataset Input Diagram

Figure 5.1 shows that a comprehensive dataset comprising various features related to past breaches is essential. This dataset may include attributes such as the type of attack (e.g., malware, phishing), the target system or network, the methods used for infiltration and exfiltration, the duration of the breach, and the severity of the impact. Additionally, contextual information such as the industry sector, geographical location, and time of occurrence can provide valuable insights into the dynamics of cyber attacks. Gathering such data may involve collating information from incident reports, security logs, threat intelligence feeds, and other sources. Careful preprocessing and labeling of the dataset are crucial to ensure the quality and relevance of the data for training machine learning models.

5.1.2 Output Design

By incorporating these elements into the output design, you can create an effective and informative interface for presenting prediction results and insights from the cyber hacking breaches prediction . Adjustments and refinements may be necessary based on the specific requirements and preferences of the stakeholders. Display the prediction results prominently on the dashboard interface. Present the predicted probability or class label for each observation in a tabular format or as a list.

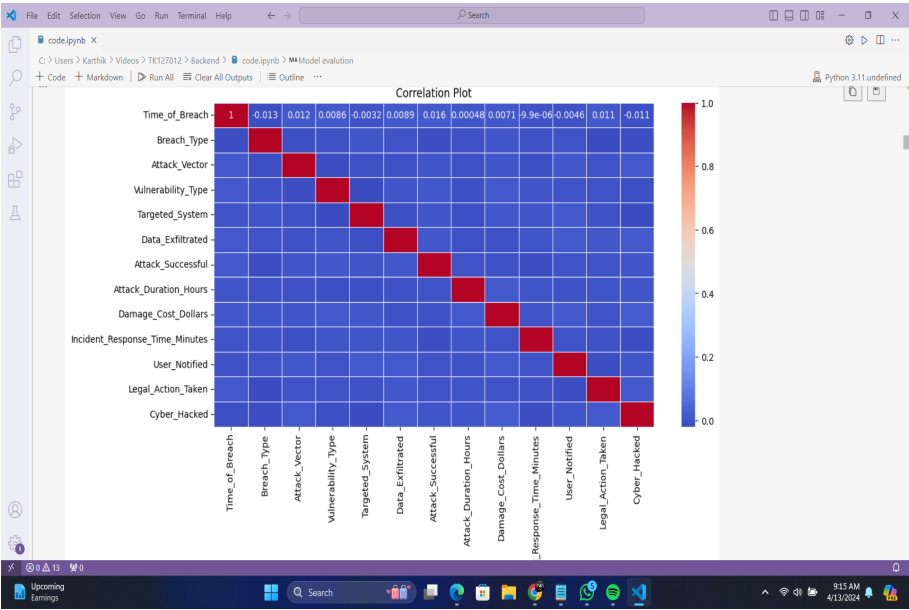


Figure 5.2: Graph Representing Types of Breaches in Output Diagram

Figure 5.2 shows that various types of attacks are observed, each with distinct characteristics and implications. Malware-based attacks involve the deployment of malicious software to compromise systems or networks, potentially leading to data theft or system disruption. Phishing attacks leverage deceptive tactics, such as fraudulent emails or websites, to trick individuals into revealing sensitive information or installing malware. Other common breach types include DDoS (Distributed Denial of Service) attacks, where malicious actors overwhelm targeted systems with excessive traffic, and SQL injection attacks, exploiting vulnerabilities in web applications to gain unauthorized access to databases. Machine learning techniques play a crucial role in detecting and mitigating these breaches by analyzing patterns and anomalies in data traffic, enhancing cybersecurity defenses against evolving threats.

5.2 Testing

The purpose of testing is to discover errors. Testing is the process of trying to discover every conceivable fault or weakness in a work product. It provides a way to check the functionality of components, sub-assemblies, assemblies and/or a finished product. It is the process of exercising software with the intent of ensuring that the Software system meets its requirements and user expectations and does not fail in an unacceptable manner. There are various types of test. Each test type addresses a specific testing requirement. Testing plays a crucial role in ensuring the reliability, accuracy, and effectiveness of a machine learning-based cyber hacking breaches prediction.

5.3 Types of Testing

5.3.1 Unit Testing

Unit testing involves the design of test cases that validate that the internal program logic is functioning properly, and that program inputs produce valid outputs. All decision branches and internal code flow should be validated. It is the testing of individual software units of the application. It is done after the completion of an individual unit before integration. This is a structural testing, that relies on knowledge of its construction and is invasive.

Unit tests perform basic tests at component level and test a specific business process, application, and/or system configuration. Unit tests ensure that each unique path of a business process performs accurately to the documented specifications and contains clearly defined inputs and expected results.

- **Define the input data:** Unit testing requires input data that can be used to test the functionality of each unit. In this case, the input data would include data about the breaches occurred, the date, and other relevant features.

- **Define the output data:** The output data would be the predicted values of leaked data. This output data would be generated using random forest classifier.

Input

```
1 # Step 1: Import necessary libraries
2 import pandas as pd # For data manipulation
3 from sklearn.model_selection import train_test_split # For splitting data into train and test sets
4 from sklearn.ensemble import RandomForestClassifier # For Random Forest classifier
5 from sklearn.metrics import accuracy_score # For evaluating model accuracy
6
7 # Step 2: Load the dataset
8 dataset = pd.read_csv("cyber_hacking_dataset.csv")
9
10 # Step 3: Prepare the data
11 # Split the dataset into features (X) and target variable (y)
12 X = dataset.drop(columns=["breach_type"]) # Features (excluding the breach type)
13 y = dataset["breach_type"] # Target variable (breach type)
14
15 # Step 4: Split the data into training and testing sets
16 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
17
18 # Step 5: Initialize the Random Forest classifier
19 rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42) # Using 100 decision
    trees
20
21 # Step 6: Train the Random Forest classifier
22 rf_classifier.fit(X_train, y_train)
23
24 # Step 7: Make predictions on the test set
25 predictions = rf_classifier.predict(X_test)
26
27 # Step 8: Evaluate the model accuracy
28 accuracy = accuracy_score(y_test, predictions)
29 print("Accuracy:", accuracy)
```

Test result

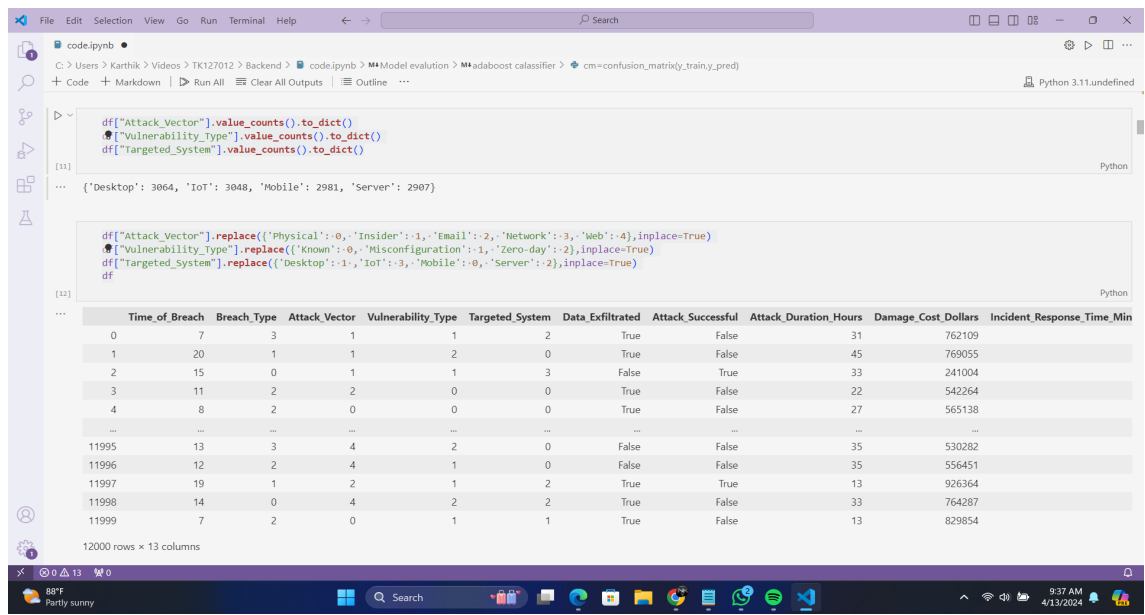


Figure 5.3: Input for Unit Testing

Figure 5.3 shows that the data preprocessing steps, such as handling missing values, encoding categorical variables, and scaling numerical features, are executed correctly. The processed data meets expected standards and is suitable for model training. Measure the computational efficiency and resource utilization of the model during training and inference. Model Training Unit T

5.3.2 Black Box Testing

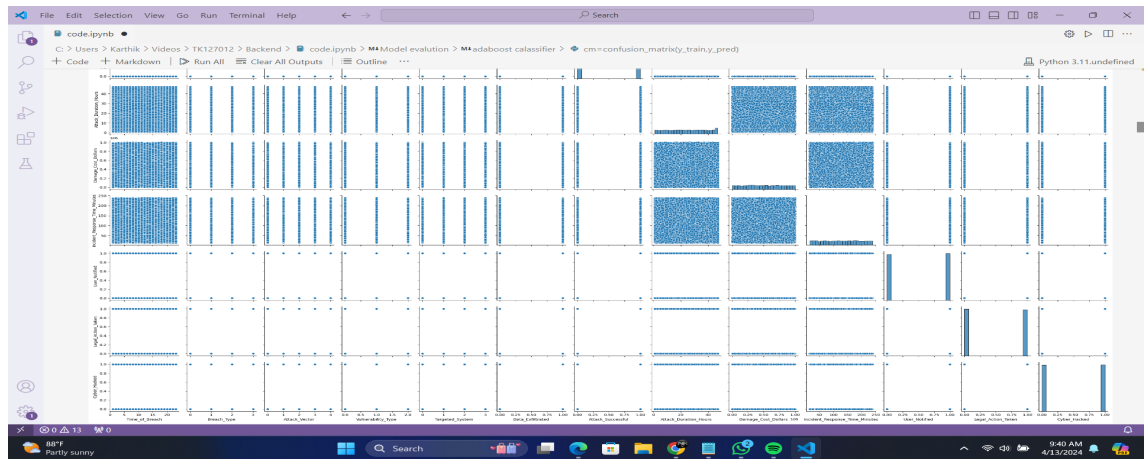


Figure 5.4: Result of Black Box Testing

Figure 5.4 illustrates the significance of black-box testing in evaluating the performance testing the software without any knowledge of the inner workings, structure or language of the module being tested. Black box tests, as most other kinds of tests, must be written from a definitive source document, such as specification or requirements document, such as specification or requirements document. It is a testing in which the software under test is treated, as a black box .you cannot “see” into it. The test provides inputs and responds to outputs without considering how the software works

5.3.3 Test Result

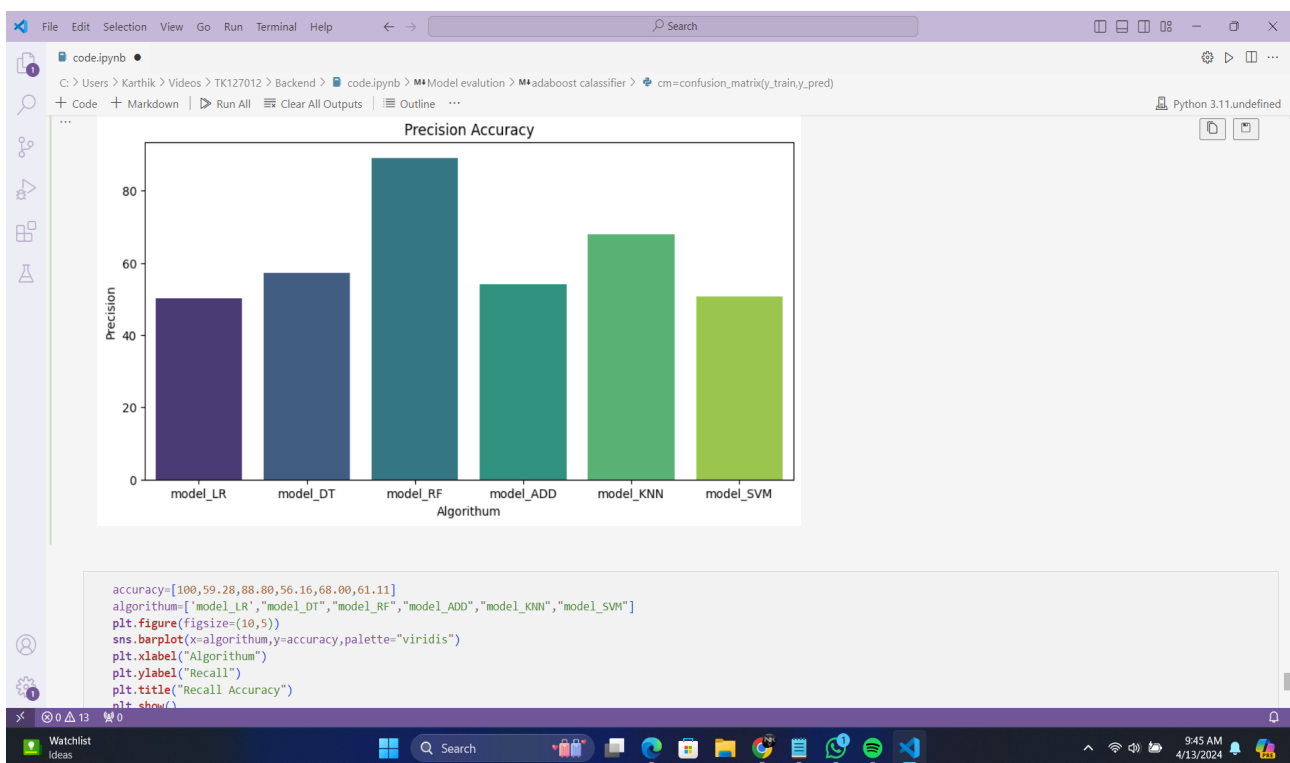


Figure 5.5: Test Image

Figure 5.5 summary of the model’s performance on the test set, including the overall accuracy. Compare the predicted labels with the actual labels in the testing dataset to calculate accuracy. $\text{Accuracy} = (\text{Number of correct predictions}) / (\text{Total number of predictions})$ Analyze the accuracy score to understand how well the model performs. A higher accuracy score indicates that the model is making more correct predictions, while a lower accuracy score suggests that the model may need improvement.

Chapter 6

RESULTS AND DISCUSSIONS

6.1 Efficiency of the Proposed System

The proposed system leverages machine learning algorithms, such as Random Forest, to analyze patterns and anomalies in data traffic, enabling the detection of potential cyber hacking breaches. The complexity of the machine learning model influences both computational requirements and predictive performance. Balancing model complexity with computational constraints is crucial for achieving a trade-off between accuracy and efficiency. Techniques such as model pruning, regularization, and algorithm optimization can help streamline model complexity without sacrificing predictive power.

The deployment architecture of the system influences its scalability, reliability, and operational efficiency. Cloud-based infrastructure, containerization, and microservices architecture enable flexible deployment and efficient resource allocation, supporting dynamic workloads and demand fluctuations. Additionally, leveraging distributed computing frameworks and parallel processing techniques can accelerate model training and inference, enhancing overall system efficiency. Performance metrics, system logs, and feedback loops facilitate real-time monitoring of model performance and resource utilization, enabling proactive identification of bottlenecks and inefficiencies. Adaptive algorithms, auto-scaling mechanisms, and periodic model retraining ensure that the system remains responsive and adaptive to evolving cyber threats and operational requirements.

Efficiency Metrics of Random Forest:

1. **Accuracy:** 0.9148888888888889
2. **Precision Score:** 0.926962457337884
3. **Recall Score:** 0.9015268864793096

6.2 Comparison of Existing and Proposed System

The existing system for detecting cyber hacking breaches may rely on traditional rule-based approaches, signature-based detection methods, or basic machine learning algorithms. While these systems may offer some level of protection, they often struggle to keep pace with the rapidly evolving nature of cyber threats. As the size and diversity of datasets increase, these systems may experience performance degradation and resource constraints, hindering their ability to provide timely and accurate threat detection.

The proposed system leverages advanced machine learning algorithms, such as Random Forest, to analyze vast amounts of data and adapt dynamically to evolving cyber threats. By learning from historical data and detecting subtle patterns and anomalies, the system can identify new and emerging attack vectors with greater accuracy and efficiency. Machine learning algorithms, trained on labeled datasets and validated against ground truth labels, can significantly reduce false positive rates compared to traditional rule-based systems. By leveraging advanced feature engineering, anomaly detection, and ensemble learning techniques, the proposed system can improve the accuracy and reliability of threat detection while minimizing false positives. System offers enhanced predictive capabilities, enabling the early detection of cyber threats and proactive mitigation of risks. By leveraging advanced algorithms such as Random Forest, deep learning, and anomaly detection, the system can identify subtle deviations from normal behavior and detect previously unseen attack patterns, enhancing cybersecurity posture and resilience.

Evaluation Metrics	Existing System(DT)	Proposed System (RF)
Accuracy	0.5607777777777778	0.9148888888888889
Precision Score	0.5568273092369478	0.926962457337884
Recall Score	0.6136313343660101	0.9015268864793096

Table 6.1: Comparing Efficiency Metrics of Existing System and Proposed System

6.3 Sample Code

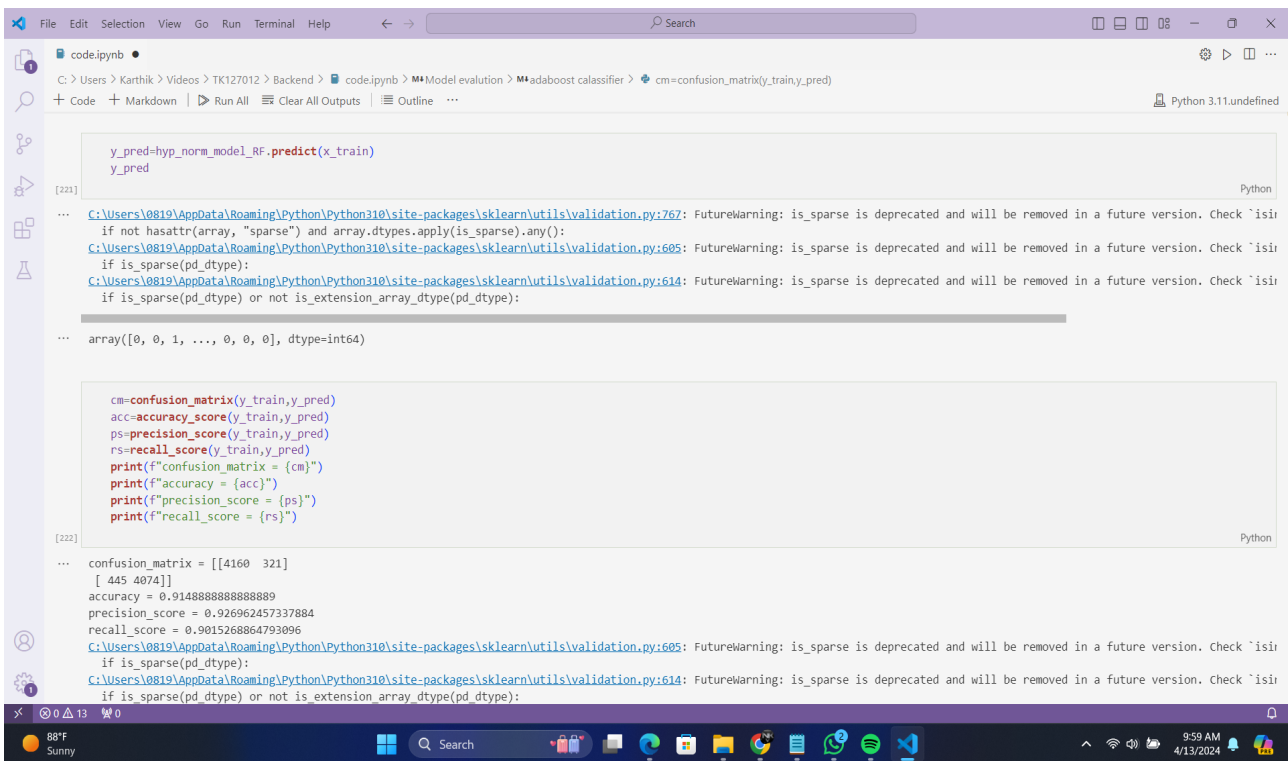
```
1 # Importing all required Liabrabry
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 from sklearn.model_selection import train_test_split, RandomizedSearchCV, GridSearchCV
8 from sklearn.metrics import accuracy_score, precision_score, recall_score, confusion_matrix, r2_score
9 from sklearn.preprocessing import MinMaxScaler, StandardScaler
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
13 from sklearn.svm import SVC
14 from sklearn.neighbors import KNeighborsClassifier
15 from sklearn.cluster import KMeans
16 import pickle
17
18 df=pd.read_csv(r"cyberhacking_dataset.csv")
19 df.head()
20
21 Time_of_Breach          0
22 Breach_Type             0
23 Attack_Vector          0
24 Vulnerability_Type     0
25 Targeted_System        0
26 Data_Exfiltrated       0
27 Attack_Successful      0
28 Attack_Duration_Hours  0
29 Damage_Cost_Dollars    0
30 Incident_Response_Time_Minutes  0
31 User_Notified          0
32 Legal_Action_Taken     0
33 Cyber_Hacked           0
34 dtype: int64
35 Index(['Time_of_Breach', 'Breach_Type', 'Attack_Vector', 'Vulnerability_Type',
36        'Targeted_System', 'Data_Exfiltrated', 'Attack_Successful',
37        'Attack_Duration_Hours', 'Damage_Cost_Dollars',
38        'Incident_Response_Time_Minutes', 'User_Notified', 'Legal_Action_Taken',
39        'Cyber_Hacked'],
40        dtype='object')
41 df["Breach_Type"].value_counts().to_dict()
42 df["Breach_Type"].replace({'Data Leak': 0, 'Malware': 1, 'SQL Injection': 2, 'Phishing': 3},
43                             inplace=True)
44 df
45 df["Attack_Vector"].value_counts().to_dict()
46 df["Vulnerability_Type"].value_counts().to_dict()
47 df["Targeted_System"].value_counts().to_dict()
```

```

47 dfdf["Attack_Vector"].replace({'Physical': 0, 'Insider': 1, 'Email': 2, 'Network': 3, 'Web': 4},
    inplace=True)
48 df["Vulnerability_Type"].replace({'Known': 0, 'Misconfiguration': 1, 'Zero-day': 2},inplace=True)
49 df["Ta
50 df["Data_Exfiltrated"].value_counts().to_dict()
51 df["Data_Exfiltrated"].replace({True: 1, False: 0},inplace=True)
52 df["Attack_Successful"].replace({True: 1, False: 0},inplace=True)
53 df["User_Notified"].replace({True: 1, False: 0},inplace=True)
54 df["Legal_Action_Taken"].replace({True: 1, False: 0},inplace=True)
55 df["Cyber_Hacked"].replace({True: 1, False: 0},inplace=True)
56 <class 'pandas.core.frame.DataFrame'>
57 RangeIndex: 12000 entries, 0 to 11999
58 Data columns (total 13 columns):
59 #    Column                                Non-Null Count  Dtype
60 ---  -
61 0    Time_of_Breach                        12000 non-null  int64
62 1    Breach_Type                          12000 non-null  int64
63 2    Attack_Vector                        12000 non-null  int64
64 3    Vulnerability_Type                  12000 non-null  int64
65 4    Targeted_System                    12000 non-null  int64
66 5    Data_Exfiltrated                    12000 non-null  int64
67 6    Attack_Successful                   12000 non-null  int64
68 7    Attack_Duration_Hours                12000 non-null  int64
69 8    Damage_Cost_Dollars                 12000 non-null  int64
70 9    Incident_Response_Time_Minutes      12000 non-null  int64
71 10   User_Notified                       12000 non-null  int64
72 11   Legal_Action_Taken                  12000 non-null  int64
73 12   Cyber_Hacked                       12000 non-null  int64
74 dtypes: int64(13)
75
76 corr=df.corr()
77 plt.figure(figsize=(12, 6))
78 sns.heatmap(corr, annot=True, cmap='coolwarm', linewidths=0.5)
79 plt.title('Correlation Plot')
80 plt.show()
81
82 plt.figure(figsize=(10,10))
83 sns.pairplot(df)
84 plt.show()
85
86 x=df.drop("Cyber_Hacked",axis=1)
87 y=df["Cyber_Hacked"]
88 11995    0
89 11996    1
90 11997    1
91 11998    0
92 11999    1
93 Name: Cyber_Hacked, Length: 12000, dtype: int64
94 x_train, x_test, y_train, y_test=train_test_split(x,y, test_size=0.25, random_state=1, stratify=y)
95 x_train.shape, x_test.shape, y_train.shape, y_test.shape

```

Output



```
code.ipynb
C:\Users\Karthik > Videos > TK127012 > Backend > code.ipynb > M4 Model evaluation > M4 adaboost classifier > cm=confusion_matrix(y_train,y_pred)
+ Code + Markdown + Run All Clear All Outputs Outline ... Python 3.11.1 undefined

y_pred=hyp_norm_model_RF.predict(x_train)
y_pred

[222]
... C:\Users\0819\AppData\Roaming\Python\Python310\site-packages\sklearn\utils\validation.py:767: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isr
if not hasattr(array, "sparse") and array.dtypes.apply(is_sparse).any():
C:\Users\0819\AppData\Roaming\Python\Python310\site-packages\sklearn\utils\validation.py:605: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isr
if is_sparse(pd_dtype):
C:\Users\0819\AppData\Roaming\Python\Python310\site-packages\sklearn\utils\validation.py:614: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isr
if is_sparse(pd_dtype) or not is_extension_array_dtype(pd_dtype):

... array([0, 0, 1, ..., 0, 0, 0], dtype=int64)

cm=confusion_matrix(y_train,y_pred)
acc=accuracy_score(y_train,y_pred)
ps=precision_score(y_train,y_pred)
rs=recall_score(y_train,y_pred)
print(f"confusion_matrix = {cm}")
print(f"accuracy = {acc}")
print(f"precision_score = {ps}")
print(f"recall_score = {rs}")

[222]
... confusion_matrix = [[4160 321]
[ 445 4074]]
accuracy = 0.9148888888888889
precision_score = 0.926962457337884
recall_score = 0.9015268864793096
C:\Users\0819\AppData\Roaming\Python\Python310\site-packages\sklearn\utils\validation.py:605: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isr
if is_sparse(pd_dtype):
C:\Users\0819\AppData\Roaming\Python\Python310\site-packages\sklearn\utils\validation.py:614: FutureWarning: is_sparse is deprecated and will be removed in a future version. Check `isr
if is_sparse(pd_dtype) or not is_extension_array_dtype(pd_dtype):
```

Figure 6.1: Random Forest output Metrics

Figure 6.1 shows that the model's performance on the test set, including the overall accuracy. A confusion matrix that shows how many images were correctly classified and misclassified for each class in the CIFAR-10 dataset. The diagonal elements of the matrix represent the number of correctly classified images, while the off-diagonal elements represent the number of misclassified images. A classification report that provides metrics such as precision, recall, and F1-score for each class in the CIFAR 10 dataset. These metrics give a more detailed understanding of the model's performance in each class. A grid of images from the test set, along with their predicted labels and probabilities. The predicted labels are shown in bold, while the true labels are shown in normal font. The probabilities represent the confidence of the model in its predictions. A set of images that the model misclassified, along with their predicted labels and true labels. These images can help to identify areas where the model needs to improve.

Chapter 7

CONCLUSION AND FUTURE ENHANCEMENTS

7.1 Conclusion

The conclusion of the proposed approach for cyber hacking breaches prediction lies in its innovative utilization of machine learning algorithms to proactively identify and forecast potential cyber threats. Unlike traditional methods, this framework leverages advanced data analytics to discern subtle patterns and anomalies within vast datasets, enabling early detection of potential hacking breaches. The integration of machine learning models empowers the system to adapt and evolve with emerging cyber threats, enhancing its predictive capabilities. This novel approach not only contributes to the field of cybersecurity but also establishes a proactive and dynamic paradigm for anticipating and mitigating cyber risks, thus fortifying the resilience of digital ecosystems. The application of machine learning techniques for detecting cyber hacking breaches represents a significant advancement in cybersecurity capabilities. By leveraging the power of advanced algorithms such as Random Forest, deep learning, and anomaly detection, organizations can enhance their ability to detect, prevent, and mitigate cyber threats effectively.

The theoretical analysis highlights the key advantages of the proposed system, including its adaptability to evolving threats, scalability to handle large-scale datasets, and ability to reduce false positive rates while enhancing predictive capabilities. By transitioning from static rule-based systems to dynamic machine learning-based solutions, organizations can stay ahead of cyber adversaries and strengthen their cybersecurity posture.

7.2 Future Enhancements

Future enhancements could involve integrating more advanced machine learning models, such as deep learning algorithms, to further improve the accuracy and robustness of cyber hacking breach predictions. Additionally, incorporating real-time data streams and leveraging anomaly detection techniques could enhance the timeliness and effectiveness of threat detection. Furthermore, exploring ensemble learning methods to combine the strengths of multiple models could yield even more reliable predictions. Moreover, developing automated response systems that can dynamically adjust security measures based on predictive insights would bolster proactive defense strategies..

- Integration of Advanced Algorithms:** Explore the integration of more advanced machine learning algorithms such as deep learning, ensemble methods, and gradient boosting algorithms. These algorithms may capture more complex patterns and relationships in the data, potentially leading to improved prediction accuracy.

- Model Interpretability:** Enhance the interpretability of the machine learning models to provide insights into how predictions are made. Techniques such as SHAP (SHapley Additive exPlanations) values, LIME (Local Interpretable Model agnostic Explanations), and model-agnostic interpretation methods can help explain the model's predictions to stakeholders.

- Real-time Monitoring:** Implement real-time monitoring capabilities to detect cyber hacking breaches as they occur or shortly thereafter. This involves integrating the prediction model into existing security systems and continuously analyzing incoming data streams for suspicious activity.

Chapter 8

INDUSTRY DETAILS

8.1 Industry name

ENCORA PVT.LTD

8.1.1 Duration of Internship (From Date - To Date)

17/01/2024 - 05/2024

8.1.2 Duration of Internship in months

6 months

8.1.3 Industry Address

Pallavaram-Thoraipakkam Road,Chennai.Tamil Nadu -600097

8.2 Internship Offer Letter



Dear ,

Nalluri Karthik

Reg ID- VTU_15337

Greetings from **Encora Innovation Pvt,Ltd**

Encora is Driven by a motto Digital Engineering for innovating at scale and enabling enterprise Modernization helps to provide companies and enterprises with complete IT Solutions to help improve their operations. Our expertise in IT Solutions helps us provide end-to-end support for IT Solutions, which includes planning, deploying, and maintaining the systems. Our ability to address customers' needs and concerns with diverse expectations and help them with new, refreshing ideas help us set ourselves apart. With the best blend of energy and expertise.the proffered partner for digitally native and enterprise clients

In reference to your application we would like to congratulate you on being selected for an internship for **"Associate Trainee"** with Encora Innovation Pvt Ltd. Your internship is scheduled to start effectively from **17-01-2024** for a Period of **120** days

The project details and technical platform will be shared with you on or before commencement of training.

Again, congratulations and we look forward to working with you.

You should report for training at the following address:

Reporting Office Address:

Pallavaram -Thoraipakkam Road, Chennai,Tamil Nadu- 600097

Contact Person: AnuhyaReddy

Contact Number- 9885488485



Dear ,

Kota Manish

Reg ID- VTU_17731

Greetings from **Encora Innovation Pvt,Ltd**

Encora is Driven by a motto Digital Engineering for innovating at scale and enabling enterprise Modernization helps to provide companies and enterprises with complete IT Solutions to help improve their operations. Our expertise in IT Solutions helps us provide end-to-end support for IT Solutions, which includes planning, deploying, and maintaining the systems. Our ability to address customers' needs and concerns with diverse expectations and help them with new, refreshing ideas help us set ourselves apart. With the best blend of energy and expertise.the proffered partner for digitally native and enterprise clients

In reference to your application we would like to congratulate you on being selected for an internship for **"Associate Trainee"** with Encora Innovation Pvt Ltd. Your internship is scheduled to start effectively from **17-01-2024** for a Period of **120** days

The project details and technical platform will be shared with you on or before commencement of training.

Again, congratulations and we look forward to working with you.

You should report for training at the following address:

Reporting Office Address:

Pallavaram -Thoraipakkam Road, Chennai,Tamil Nadu- 600097

Contact Person: AnuhyaReddy

Contact Number- 9885488485



Project Commencement Form

Name of the Industry: ENCORA

Address: Pallavaram -Thoraipakkam Road, Chennai,Tamil Nadu- 600097

Team Details:

S.No	ID No	Student Name	Degree & Branch
1.	VTU15337	NALLURI KARTHIK	B.TECH/CSE
2.	VTU17731	KOTA MANISH	

Date of reporting for project work: 10/02/2024

Name of the Industry Supervisor : Kuppam vasanthi

Department : IT

Designation : Internship Manager

Contact Number :9885488485

Email ID : Kuppamvasanthi@gmail.com

Name of the Internal Supervisor : Anuhya Reddy

Contact No. :9885488485

Email ID : encoraanuhya@gmail.com

Tentative Project Title / Project domain: Cyber hacking breaches using machine learning

Chapter 9

PLAGIARISM REPORT

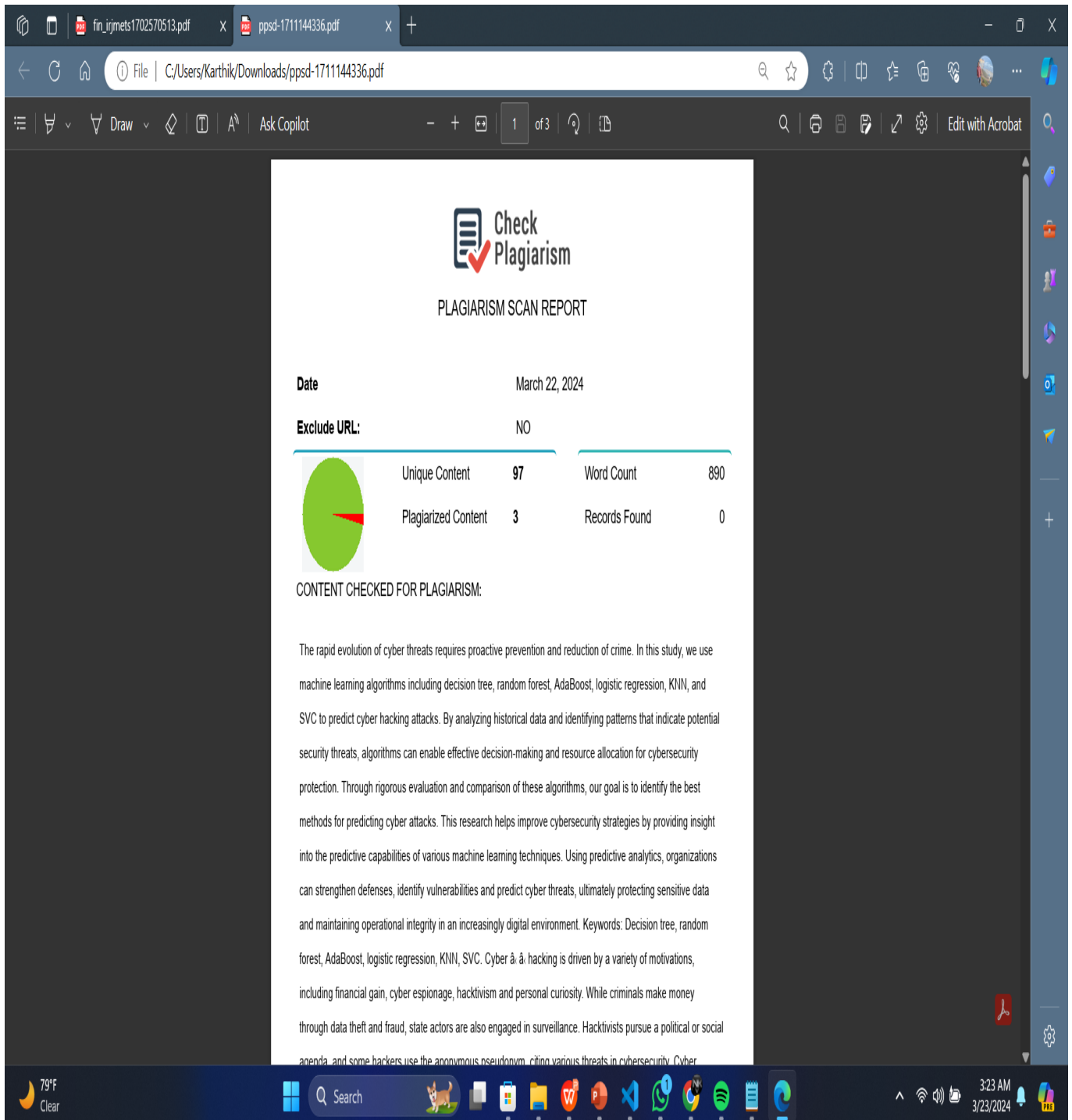


Figure 9.1: Plagiarism Report

Chapter 10

SOURCE CODE & POSTER PRESENTATION

10.1 Source Code

```
1 from flask import Flask, render_template, request, url_for, flash, redirect, session
2 import pandas as pd
3 import numpy as np
4 from sklearn.ensemble import ExtraTreesClassifier
5 from sklearn.svm import SVC
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.ensemble import AdaBoostClassifier
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.model_selection import train_test_split
10 from sklearn import preprocessing
11 from sklearn.metrics import accuracy_score
12 from sklearn.tree import DecisionTreeClassifier
13 import mysql.connector
14 db=mysql.connector.connect(user="root",password="",port='3306',database='cyber_attack')
15 cur=db.cursor()
16
17 app=Flask(__name__)
18 app.secret_key="CBJcb786874wrf78chdchsdcv"
19
20 @app.route('/')
21 def index():
22     return render_template('index.html')
23
24 @app.route('/about')
25 def about():
26     return render_template('about.html')
27
28
29 @app.route('/login',methods=['POST','GET'])
30 def login():
31     if request.method=='POST':
32         useremail=request.form['useremail']
33         session['useremail']=useremail
34         userpassword=request.form['userpassword']
35         sql="select * from user where Email='%s' and Password='%s'"%(useremail,userpassword)
```



```

36     cur.execute(sql)
37     data=cur.fetchall()
38     db.commit()
39     if data ==[]:
40         msg="user Credentials Are not valid"
41         return render_template("login.html",name=msg)
42     else:
43         return render_template("load.html",myname=data[0][1])
44     return render_template('login.html')
45
46 @app.route('/registration',methods=["POST","GET"])
47 def registration():
48     if request.method=="POST":
49         username=request.form['username']
50         useremail = request.form['useremail']
51         userpassword = request.form['userpassword']
52         conpassword = request.form['conpassword']
53         Age = request.form['Age']
54         address = request.form['address']
55         contact = request.form['contact']
56         if userpassword == conpassword:
57             sql="select * from user where Email='%s' and Password='%s'"%(useremail ,userpassword)
58             cur.execute(sql)
59             data=cur.fetchall()
60             db.commit()
61             print(data)
62             if data ==[]:
63
64                 sql = "insert into user(Name,Email,Password,Age,Address,Contact) values(%s,%s,%s,%s,%s,%s)"
65                 val=(username ,useremail ,userpassword ,Age, address , contact)
66                 cur.execute(sql , val)
67                 db.commit()
68                 flash("Registered successfully","success")
69                 return render_template("login.html")
70             else:
71                 flash("Details are invalid","warning")
72                 return render_template("registration.html")
73         else:
74             msg = "Password doesn't match"
75             return render_template("registration.html",msg=msg)
76     return render_template('registration.html')
77
78 @app.route('/load',methods=["GET","POST"])
79 def load():
80     global df, dataset
81     if request.method == "POST":
82         data = request.files['data']
83         df = pd.read_csv(data)
84         dataset = df.head(100)

```

```

85     msg = 'Data Loaded Successfully'
86     return render_template('load.html', msg=msg)
87 return render_template('load.html')
88
89
90 @app.route('/view')
91 def view():
92     print(dataset)
93     print(dataset.head(2))
94     print(dataset.columns)
95     return render_template('view.html', columns=dataset.columns.values, rows=dataset.values.tolist()
96                             )
97
98 @app.route('/preprocess', methods=['POST', 'GET'])
99 def preprocess():
100     global x, y, x_train, x_test, y_train, y_test, hvectorizer, df
101     if request.method == "POST":
102         size = int(request.form['split'])
103         size = size / 100
104         dataset["Breach_Type"].replace({'Data Leak': 0, 'Malware': 1, 'SQL Injection': 2, 'Phishing':
105                                         : 3}, inplace=True)
106         dataset["Attack_Vector"].replace({'Physical': 0, 'Insider': 1, 'Email': 2, 'Network': 3, '
107                                         Web': 4}, inplace=True)
108         dataset["Vulnerability_Type"].replace({'Known': 0, 'Misconfiguration': 1, 'Zero-day': 2},
109                                                inplace=True)
110         dataset["Targeted_System"].replace({'Desktop': 1, 'IoT': 3, 'Mobile': 0, 'Server': 2},
111                                              inplace=True)
112         dataset["Data_Exfiltrated"].replace({True: 1, False: 0}, inplace=True)
113         dataset["Attack_Successful"].replace({True: 1, False: 0}, inplace=True)
114         dataset["User_Notified"].replace({True: 1, False: 0}, inplace=True)
115         dataset["Legal_Action_Taken"].replace({True: 1, False: 0}, inplace=True)
116         dataset["Cyber_Hacked"].replace({True: 1, False: 0}, inplace=True)
117
118     # Assigning the value of x and y
119     x=dataset.drop("Cyber_Hacked",axis=1)
120     y=dataset["Cyber_Hacked"]
121     x_train, x_test, y_train, y_test=train_test_split(x,y, test_size=size, random_state=1, stratify=y)
122
123     # describes info about train and test set
124     print("Number transactions X_train dataset: ", x_train.shape)
125     print("Number transactions y_train dataset: ", y_train.shape)
126     print("Number transactions X_test dataset: ", x_test.shape)
127     print("Number transactions y_test dataset: ", y_test.shape)
128
129     print(x_train, x_test)
130
131     return render_template('preprocess.html', msg='Data Preprocessed and It Splits Successfully'
132                             )

```

```

129     return render_template('preprocess.html')
130
131 @app.route('/model', methods=['POST', 'GET'])
132 def model():
133     if request.method == "POST":
134         global model
135         print('cccccccccccccccccccccccccccccccccccccccccccccccccccccccc')
136         s = int(request.form['algo'])
137         if s == 0:
138             return render_template('model.html', msg='Please Choose an Algorithm to Train')
139         elif s == 1:
140             print('aaaaaaaaaaaaaaaaabbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbb')
141             from sklearn.ensemble import AdaBoostClassifier
142             ad = AdaBoostClassifier()
143             ad.fit(x_train, y_train)
144             y_pred = ad.predict(x_train)
145             ac_ad = accuracy_score(y_train, y_pred)
146             ac_ad = ac_ad * 100
147             print('aaaaaaaaaaaaaaaaaaaaa')
148             msg = 'The accuracy obtained by AdaBoost Classifier is ' + str(ac_ad) + str('%')
149             return render_template('model.html', msg=msg)
150         elif s == 2:
151             classifier = DecisionTreeClassifier(max_leaf_nodes=39, random_state=0)
152             classifier.fit(x_train, y_train)
153             y_pred = classifier.predict(x_train)
154
155             ac_dt = accuracy_score(y_train, y_pred)
156             ac_dt = ac_dt * 100
157             msg = 'The accuracy obtained by Decision Tree Classifier is ' + str(ac_dt) + str('%')
158             return render_template('model.html', msg=msg)
159         elif s == 3:
160             from sklearn.ensemble import RandomForestClassifier
161             rf=RandomForestClassifier(n_estimators = 50,
162                                     min_samples_split = 3,
163                                     min_samples_leaf = 2,
164                                     max_features = 'log2',
165                                     max_depth = 10,
166                                     bootstrap = True)
167
168             rf.fit(x_train, y_train)
169             rf=rf.fit(x_train, y_train)
170             y_pred = rf.predict(x_train)
171
172             ac_rf = accuracy_score(y_train, y_pred)
173             ac_rf = ac_rf * 100
174             msg = 'The accuracy obtained by random Forest Classifier is ' + str(ac_rf) + str('%')
175             return render_template('model.html', msg=msg)
176         elif s == 4:
177             from sklearn.neighbors import KNeighborsClassifier
178             knn = KNeighborsClassifier(n_neighbors=12)
179             knn.fit(x_train, y_train)

```

```

179         y_pred = knn.predict(x_train)
180
181         ac_knn = accuracy_score(y_train, y_pred)
182         ac_knn = ac_knn * 100
183         msg = 'The accuracy obtained by K-Nearest Neighbour is ' + str(ac_knn) + str('%')
184         return render_template('model.html', msg=msg)
185     elif s == 5:
186         svc = SVC()
187         svc.fit(x_train, y_train)
188         y_pred = svc.predict(x_train)
189
190         ac_svc = accuracy_score(y_train, y_pred)
191         ac_svc = ac_svc * 100
192         msg = 'The accuracy obtained by support vector Classifier is ' + str(ac_svc) + str('%')
193         return render_template('model.html', msg=msg)
194     elif s == 6:
195         from sklearn.linear_model import LogisticRegression
196
197         lr = LogisticRegression()
198         lr.fit(x_train, y_train)
199         y_pred = lr.predict(x_train)
200
201         ac_lr = accuracy_score(y_train, y_pred)
202         ac_lr = ac_lr * 100
203         msg = 'The accuracy obtained by Logistic Regression is ' + str(ac_lr) + str('%')
204         return render_template('model.html', msg=msg)
205     return render_template('model.html')
206
207 @app.route('/prediction', methods=['GET', 'POST'])
208 def prediction():
209     if request.method == "POST":
210         # f1=int(request.form['city'])
211         f1 = float(request.form['Time_of_Breach'])
212         f2 = float(request.form['Breach_Type'])
213         f3 = float(request.form['Attack_Vector'])
214         f4 = float(request.form['Vulnerability_Type'])
215         f5 = float(request.form['Targeted_System'])
216         f6 = float(request.form['Data_Exfiltrated'])
217         f7 = float(request.form['Attack_Successful'])
218         f8 = float(request.form['Attack_Duration_Hours'])
219         f9 = float(request.form['Damage_Cost_Dollars'])
220         f10 = float(request.form['Incident_Response_Time_Minutes'])
221         f11 = float(request.form['User_Notified'])
222         f12 = float(request.form['Legal_Action_Taken'])
223
224
225         print(f2)
226         print(type(f2))
227
228         li = [f1, f2, f3, f4, f5, f6, f7, f8, f9, f10, f11, f12]


```

```

229     print(li)
230
231     # model.fit(X_transformed, y_train)
232
233     # print(f2)
234     # import pickle
235     # filename = 'randomforest.sav'
236     # model = pickle.load(open(filename, 'rb'))
237     # result = model.predict([li])
238     from sklearn.ensemble import RandomForestClassifier
239     rf=RandomForestClassifier(n_estimators = 50,
240                               min_samples_split = 3,
241                               min_samples_leaf = 2,
242                               max_features = 'log2',
243                               max_depth = 10,
244                               bootstrap = True)
245
246     rf.fit(x_train, y_train)
247     rf=rf.fit(x_train, y_train)
248     result = rf.predict([li])
249     print(result)
250     print('result is ', result)
251     # (Normal = 0, Cyber_Hacked = 1 )
252     if result == 0:
253         msg = 'There is No-Cyber Hacked'
254         return render_template('prediction.html', msg=msg)
255     else:
256         msg = 'There is Cyber Hacked'
257         return render_template('prediction.html', msg=msg)
258
259     return render_template('prediction.html')
260
261
262
263
264
265
266
267
268 if __name__ == '__main__':
269     app.run(debug=True)

```

10.2 Poster Presentation



Vel Tech
Rangarajan Dr. Sagunthala
Vellore Institute of Technology
Vellore, Tamil Nadu 620 015

CYBER HACKING BREACHES PREDICTION USING MACHINE LEARNING
Department of Computer Science and Engineering
School of Computing
1156CS701-MAJOR PROJECT
INTERNSHIP THROUGH DIND
ENCORA PVT.LTD
WINTER SEMESTER 2023-2024

Batch: (2020-2024)

ABSTRACT

Cyber-physical systems (CPS) have made significant progress in many dynamic applications due to the integration between physical processes, computational resources, and communication capabilities. However, cyber-attacks are a major threat to these systems. Unlike faults that occur by accidents cyber-physical systems, cyber-attacks occur intelligently and stealthily. Some of these attacks which are called deception attacks, inject false data from sensors or controllers, and also by compromising with some cyber components, corrupt data, or enter misinformation into the system. If the system is unaware of the existence of these attacks, it won't be able to detect them, and performance may be disrupted or disabled altogether. Therefore, it is necessary to adapt algorithms to identify these types of attacks in these systems. It should be noted that the data generated in these systems is produced in very large number, with so much variety, and high speed, so it is important to use machine learning algorithms to facilitate the analysis and evaluation of data and to identify hidden patterns.

TEAM MEMBER DETAILS

<Student 1. vtu 15337/Nalluri Karthik>
<Student 2. vtu 17731/Kota Manish >
<Student 1. Phone no-9353354620>
<Student 2. Phone no-9618012585>
<Student 1. -vtu15337@veltech.edu.in>
<Student 2. -vtu17731@veltech.edu.in>

INTRODUCTION

The advent of the digital age has ushered in unprecedented opportunities for connectivity and innovation. However, it has also given rise to a growing menace – cyber hacking breaches. In recent years, cyberattacks have become increasingly sophisticated and devastating, posing a significant threat to individuals, businesses, and even nations. This project aims to comprehensively investigate and analyze cyber hacking breaches that occurred in the past year. We will delve into the methods, motivations, and impacts of these breaches to gain a deeper understanding of the evolving landscape of cyber threats. By examining a range of high-profile cases, we intend to identify common vulnerabilities and attack vectors.

Our objectives include mapping the tactics employed by hackers, assessing the effectiveness of security measures, and evaluating the financial and reputational costs incurred by victims. Additionally, we will explore the ethical, legal, and regulatory aspects surrounding cyberattacks and data breaches.

METHODOLOGIES

Take the Dataset: The system accepts and processes the dataset provided by the user. This dataset forms the foundation for building the predictive model.

Preprocessing: Before training a predictive model, the system preprocesses the dataset. This includes handling missing data, data cleaning, and feature extraction. Preprocessing ensures that the data is in a suitable format for modeling.

Training: The system uses machine learning techniques and Python modules to train a model based on the preprocessed dataset. The model learns patterns and relationships within the data, allowing it to make predictions.

Generate Results: Once the model is trained, the system can generate results based on user input values. These results typically indicate whether the input data corresponds to a specific condition, event, or prediction, such as detecting predict the type of Disease.

RESULTS

Random forest classifier

confusion_matrix = [[4160 321]
[445 4074]]

accuracy = 0.9148888888888889
precision_score = 0.926962457337884
recall_score = 0.9015268864793096

Table 1. Label in Zipit Callbot.

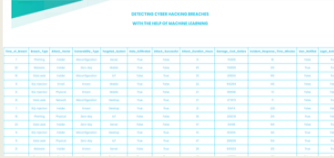


Figure 1. Label in Zipit Callbot.




Chart 1. Label in Zipit Callbot.


STANDARDS AND POLICIES

General Data Protection Regulation (GDPR):
GDPR establishes rules for data protection and privacy for individuals within the European Union (EU) and the European Economic Area (EEA). It imposes obligations on organizations regarding the collection, processing, and protection of personal data.

NIST Cybersecurity Framework:
Developed by the National Institute of Standards and Technology (NIST), the NIST Cybersecurity Framework provides a set of guidelines, standards, and best practices for improving cybersecurity risk management.

Ethical Guidelines for AI and Machine Learning:
Adhere to ethical guidelines and principles for the development and deployment of artificial intelligence (AI) and machine learning systems. This includes considerations such as fairness, transparency, accountability, and the avoidance of bias in algorithmic decision-making.

Figure 2. Label in Zipit Callbot.



CONCLUSIONS

The conclusion of the proposed approach for cyber hacking breaches prediction lies in its innovative utilization of machine learning algorithms to proactively identify and forecast potential cyber threats. Unlike traditional methods, this framework leverages advanced data analytics to discern subtle patterns and anomalies within vast datasets, enabling early detection of potential hacking breaches.

ACKNOWLEDGEMENT

1. Project Supervisor Name- DR.P.J.BESLIN PAJILA,ME.,Ph.D.,
2. Project supervisor Contact No--918056458326
3. Project supervisor Mail ID-drbeslinpajila@veltech.edu.in

Figure 10.1: Poster Presentation

References

- [1] Kure, H.I., Islam, S., Ghazanfar, M. et al. Asset criticality and risk prediction for an effective cybersecurity risk management of cyber-physical system. *Neu 47 ral Comput Applic* 34, 493–514 (2022). PP. 1-6, 2021.
- [2] Kwon, Cheolhyeon, Weiyi Liu, and Inseok Hwang. "Security analysis for cyber physical systems against stealthy deception attacks." In 2013 American control conference, IEEE (2013): 3344-3349,2013
- [3] Mandal, S., Saha, B., Nag, R. (2020). Exploiting Aspect-Classified Sentiments for CyberCrime Analysis and Hack Prediction. In: Kar, N., Saha, A., Deb, S. (eds) Trends in Computational Intelligence, Security and Internet of Things. IC CISIoT 2020. Communications in Computer and Information Science, vol 1358.PP. 1542-1552, 2020.
- [4] M. Eling and W. Schnell, "What can we realize cyber risk and cyber risk insurance?" *J. Risk Finance*, vol. 17, no. 5, pp. 474–491, 2016.
- [5] Pajic, Miroslav, James Weimer, Nicola Bezzo, Oleg Sokolsky, George J. Pappas, and Insup Lee. "Designandimplementationof attack-resilient cyberphysical systems: With a focus on attack-resilient state estimators." *IEEE Control Systems Magazine* 37, no. 2 (2017): 66-81.
- [6] Poyraz, O.I., Canan, M., McShane, M. et al. Cyber assets at risk: monetary impact of U.S. personally identifiable information mega data breaches. *Geneva Pap Risk Insur Issues Pract* 45, 616–638 (2020).
- [7] R. R.Subramanian, R.Avula, P. S. SuryaandB.Pranay, "ModelingandPredict ing Cyber Hacking Breaches," 2021 5th International Conference on Intelligent Computing and Control Systems (ICICCS), 2021, pp. 288-293
- [8] Sun, Hongtao, Chen Peng, Taicheng Yang, Hao Zhang, and Wangli He. "Re silient control of networked control systems with stochastic denial of service attacks." *Neurocomputing* 270 (2017): 170-177.
- [9] Sheng, Long, Ya-Jun Pan, and Xiang Gong. "Consensus formation control for a class of networked multiple mobile robot systems." *Journal of Control Science and Engineering* 2012

- [10] Zeng, Went, and Mo-Yuen Chow. “Resilient distributed control in the presence of misbehaving agents in networked control systems.” *IEEE transactions on cybernetics* 44, no. 11 (2014): 2038-2049.