

MALWARE DETECTION AND CLASSIFICATION

Submitted By:

B Sairam 21BLC1468

Veekshitha Nagarani 21BCE8943

Shivani Narayan L 21BLC1119

Mridhula S 21BLC1144

For Smart Internz

During November 2023



Table of Contents

Abstraction	3
Introduction	3
Empathy Map	4
Extracting Methods	5
Proposed Solution	5
Solution Architecture	10
Architecture Design	11
Technical Architecture	11
Project Planning	12
Code Solution	13
Test and Validation	15
Performance	17
Results	17
Conclusion	18
Future Scope	19
Appendix	20

ABSTRACTION

Malware, malicious software designed to infiltrate and harm computer systems, continues to pose a significant threat to cybersecurity. As technology advances, so does the sophistication of malware, making it imperative to develop robust methods for its detection and classification. This project aims to address this pressing concern by proposing a comprehensive approach to malware detection and classification.

The primary objective of this project is to design and implement an efficient and accurate malware detection system that can identify and classify various types of malwares. To achieve this, we will employ a combination of machine learning techniques, data analysis, and signature-based methods.

Our approach involves collecting and curating a diverse dataset of malware samples, including viruses, worms, Trojans, and ransomware. We will then extract relevant features and develop machine learning models capable of distinguishing between legitimate software and malicious code. These models will be trained on the labelled dataset and fine-tuned to achieve optimal performance.

Furthermore, we will explore advanced techniques such as deep learning and behaviour analysis to enhance the accuracy of our detection system. Realtime monitoring and anomaly detection will also be incorporated to detect previously unseen or zero-day malware threats.

Ultimately, this project aims to contribute to the ongoing efforts to safeguard computer systems and sensitive data from the ever-evolving landscape of cyber threats. By developing a robust malware detection and classification system, we intend to enhance the security posture of organizations and individuals, thereby mitigating the risks associated with malware infections.

INTRODUCTION:

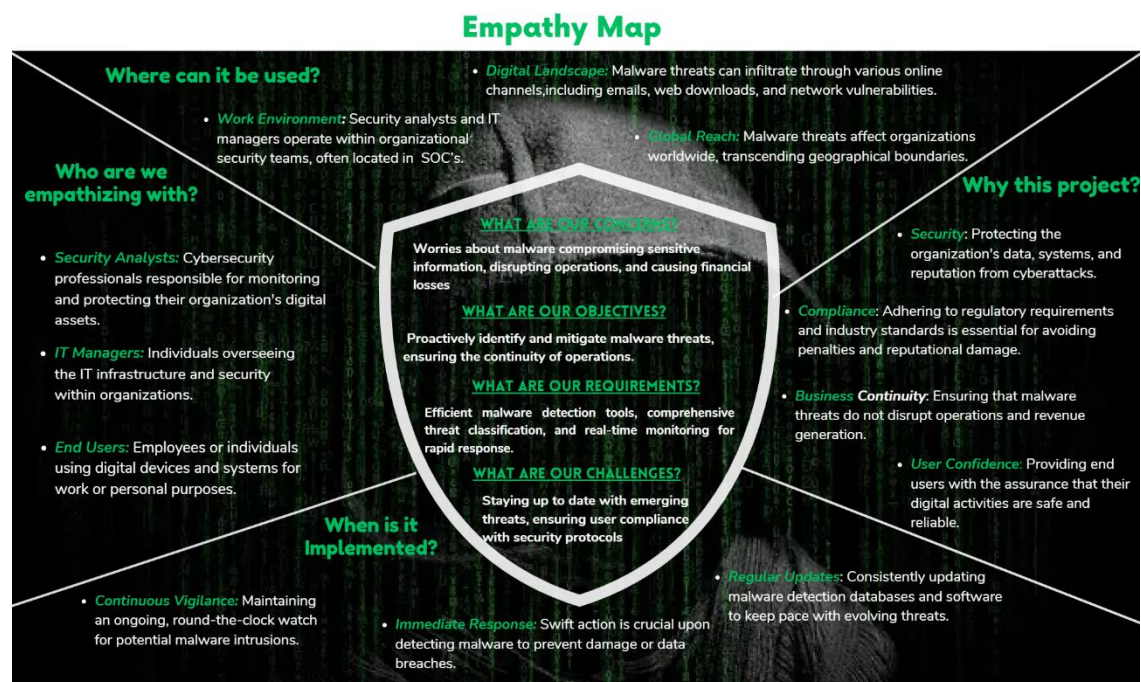
In an increasingly interconnected digital landscape, where the omnipresence of technology has ushered in unparalleled convenience, the shadow of malevolent software, colloquially known as malware, looms ominously. The project, "Malware Detection and Classification Using Machine Learning," embarks on a formidable quest to safeguard the intricate web of cyberspace from the nefarious threats that seek to exploit its vulnerabilities.

At its core, this endeavour epitomizes the synergy between cutting-edge technology and cybersecurity, harnessing the formidable power of machine learning to decipher, dissect, and ultimately defeat the multifaceted spectrum of digital threats. With a mosaic of malicious entities, from Trojan horses to ransomware, consistently evolving in complexity, this project represents an unwavering commitment to stay one step ahead of the adversaries.

The heart of this innovation lies in the creation of intelligent algorithms that possess the acumen to discern benign from maleficent, to categorize the clandestine codes that traverse our networks, and to institute an autonomous defence mechanism that operates with unparalleled precision. It stands as a paragon of proactive security, fortifying the citadels of our digital age against the relentless incursion of malware, thus ensuring the sanctity of data, the privacy of individuals, and the resilience of institutions.

In an era where the vitality of our digital existence hinges on our ability to safeguard our cyber frontiers, the "Malware Detection and Classification Using Machine Learning" project emerges as a sentinel of technological advancement, tirelessly working to preserve the integrity and trustworthiness of the virtual realms we inhabit.

EMPATHY MAP



EXISTING METHODS:

While not yet widely implemented, the application of machine learning techniques for malware detection is not a novel concept. Various studies have delved into this realm, primarily focused on evaluating the accuracy of distinct methodologies.

In the paper titled "Malware Detection Using Machine Learning," Dragos Gavrilut embarked on the development of a detection system centered on modified perceptron algorithms. Gavrilut's research yielded accuracy rates ranging from 69.90% to 96.18% for different algorithms. Notably, the algorithms with the highest accuracy also exhibited an increased rate of false positives, with the most precise among them resulting in 48 false positives. The algorithm that struck a balance between accuracy and a low false-positive rate achieved an accuracy of 93.01% (Gavrilut, et al., 2009).

Another research endeavour, "A Static Malware Detection System Using Data Mining Methods," proposed extraction techniques based on PE headers, DLLs, and API functions. The study employed Naive Bayes, J48 Decision Trees, and Support Vector Machines. The J48 algorithm, in particular, stood out with the highest overall accuracy, reaching 99% when considering various feature types, including PE headers, hybrid PE headers, and API functions. This research was conducted by Baldangombo, Jambaljav, and Horng in 2013 (Baldangombo, Jambaljav, and Horng, 2013).

PROPOSED SOLUTION

S No.	Parameter	Description
1.	Problem Statement (Problem to be solved)	In an age of escalating digitalization, the project "Malware Detection Using Machine Learning" emerges as a critical endeavour aimed at combating the pervasive threat of malware. Employing advanced machine learning techniques, it seeks to proactively identify and mitigate malicious software, transcending traditional signature-based methods. This ambitious initiative aspires to create a robust, real-time system capable of accurately distinguishing between legitimate software and threats, all while adapting to the evolving landscape of malware and offering cross platform defence. Through the amalgamation of machine learning, anomaly detection, and adept feature engineering, the project endeavours to fortify digital landscapes against the relentless spectre of cyber threats, safeguarding both individual users and organizations.

2.	Idea / Solution description	<p>Approach: Employ a holistic approach to malware detection, harnessing machine learning to conduct in-depth behavioural analysis, craft an extensive feature set, and implement anomaly detection.</p> <p>Algorithmic Versatility: Utilize a diverse range of machine learning algorithms, such as decision trees, random forests, and deep learning models, for heightened accuracy.</p> <p>Real-time Vigilance: Ensure real-time monitoring and response capabilities through</p>
		<p>streaming data processing, swiftly identifying and neutralizing threats.</p> <p>Cross-Platform Resilience: Guarantee compatibility across various operating systems and devices, fostering adaptability in the face of an evolving threat landscape.</p> <p>Dynamic Model Updates: Implement automatic model updates to remain effective against emerging malware threats.</p> <p>User-Friendly Interface: Create an intuitive interface for users and administrators, facilitating data interpretation and response.</p> <p>Iterative Refinement: Continuously enhance the system through a feedback loop, fine-tuning algorithms and features for optimal performance against evolving threats.</p>

3.	Novelty / Uniqueness	<p>Our endeavour involves the establishment of a distinguished Git repository, meticulously curated to be accessible for download and practical use. The core mission of this repository is the innovative pursuit of malware detection, encompassing an eclectic array of file formats, including but not limited to .exe and web URLs. This ambitious undertaking marks a pioneering stride in the realm of cybersecurity, with the ultimate goal of fortifying digital landscapes against the ever encroaching menace of malicious software. Through the convergence of cutting-edge technology and rigorous research, our project aspires to set a new benchmark in the domain of threat mitigation, contributing to the broader spectrum of digital security in a truly novel and distinctive manner.</p>
4.	Social Impact / Customer Satisfaction	<p>The potential social impact and customer satisfaction derived from this project are poised to be profound and multifaceted. The innovative strides made in malware detection will fundamentally bolster digital security, safeguarding the integrity of data and the privacy of individuals and organizations alike. Through the proactive identification and mitigation of malicious software across diverse</p>
		<p>file formats, our project not only mitigates the immediate threats but also engenders a heightened sense of digital empowerment and trust. This, in turn, augments customer satisfaction by assuring users that their digital experiences are fortified and protected. The ramifications extend beyond mere technological advancement; they encompass the preservation of digital sanctity and the enhancement of online interactions, thereby underscoring the enduring social relevance and customer contentment intrinsic to this pioneering endeavor.</p>

5.	Business Model (Revenue Model)	<p>Initially it will be downloaded as Git repo and it's free for everyone. One of our future scopes is to make a website dashboard with this model. for that, the revenue model encompasses:</p> <p>Subscription Plans: Tiered user subscriptions offering advanced features and real-time updates.</p> <p>Enterprise Licensing: Customized plans for organizations with tailored services.</p> <p>Data Analytics and Threat Reports: Offering organizations data analytics and threat intelligence services.</p> <p>Consulting and Training: Providing consulting and training programs for security optimization.</p> <p>Affiliate Partnerships: Collaborating with affiliates for distribution and revenue-sharing.</p> <p>Data Privacy Compliance: Services to ensure data privacy compliance.</p> <p>Advisory Services: Specialized advisory and threat analysis reports.</p> <p>Freemium Model: Offering a free version with premium paid features.</p> <p>Data Licensing: Licensing anonymized threat data to research and government agencies.</p>
		<p>Strategic Partnerships: Integrating with tech companies and cloud providers for licensing and royalties.</p>

6.	Scalability of the Solution	<p>Data Volume: The system can effectively process and analyse large volumes of data, accommodating the increasing scale of digital threats and the expansion of data sources.</p> <p>User Base: It can seamlessly scale from individual users to large enterprises, catering to the diverse needs of a growing user base.</p> <p>Threat Landscape: The system's adaptability allows it to stay ahead of emerging threats and evolving malware tactics, ensuring ongoing effectiveness.</p> <p>Platform and Device Diversity: It supports a wide range of operating systems, devices, and file formats, making it compatible with an ever-expanding digital ecosystem.</p> <p>Geographical Expansion: The project can extend its reach across geographical boundaries, serving users and organizations worldwide.</p> <p>Real-Time Processing: It can handle real-time data streams, critical for promptly identifying and mitigating threats in a dynamic online environment.</p> <p>Feature Enhancement: The addition of new features, algorithms, and threat detection capabilities can be seamlessly integrated to bolster performance.</p> <p>Enterprise Integration: Its adaptability to enterprise requirements allows for the incorporation of specific customizations and integration with existing cybersecurity infrastructure.</p> <p>Partner Ecosystem: The project can effortlessly forge partnerships with other cybersecurity entities, amplifying its reach and value proposition.</p>
----	-----------------------------	---

		<p>Service Expansion: The revenue model can introduce new services and offerings to meet the evolving needs of users and organizations.</p> <p>The project's scalability is a foundational element, ensuring that it remains agile, responsive, and adept at addressing the ever evolving and expanding landscape of digital threats. This inherent scalability positions it as a viable and enduring solution within the dynamic realm of cybersecurity.</p>
--	--	---

SOLUTION ARCHITECTURE

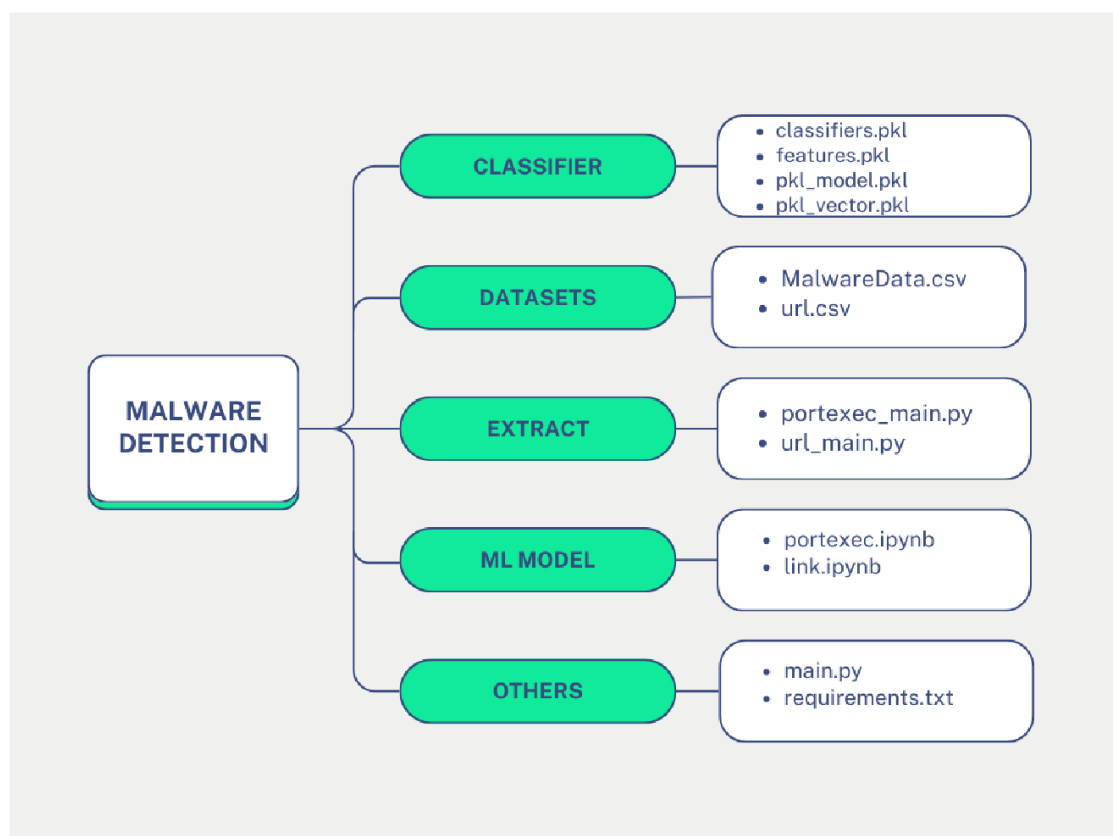
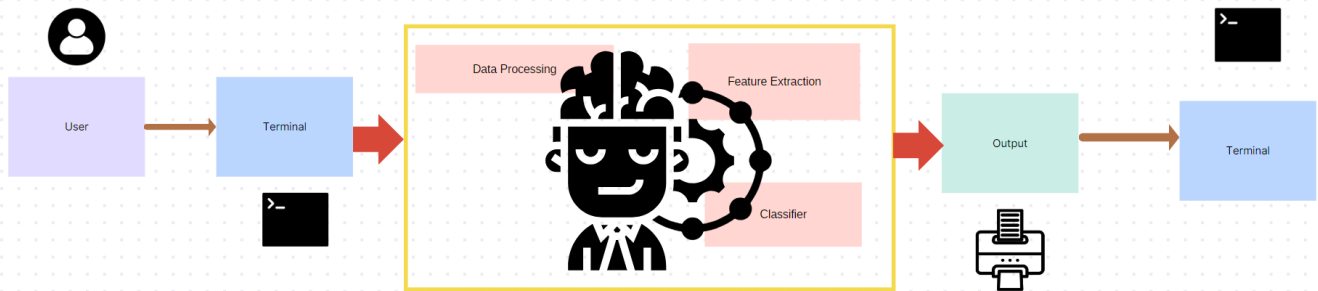


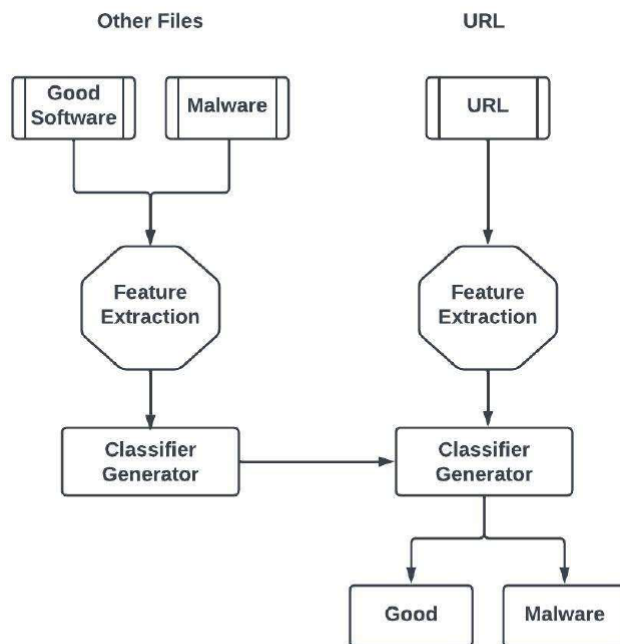
Figure 1: Architecture of Malware classification and detection

ARCHITECTURE DESIGN

Architecture Design



TECHNICAL ARCHITECTURE



PROJECT PLANNING

)

Product Backlog, Sprint Schedule, and Estimation

Sprint	Functional Requirement (Epic)	User Story Number	User Story / Task	Story Points	Priority	Team Members
Sprint-1	Registration	USN-1	As a user, I can download the Git Repository through terminal	2	High	4
Sprint-1		USN-2	As a user, I can download the requirements.txt file	1	Medium	4
Sprint-1		USN-3	As a user, I can interact with the interface in terminal	2	High	4
Sprint-2	USAGE	USN-4	As a user, I can upload the exe files and url to it	2	High	4
Sprint-2		USN-5	As a user, I can run the model and get the results	1	High	4

Project Tracker, Velocity & Burndown Chart:

Sprint	Total Story Points	Duration	Sprint Start Date	Sprint End Date (Planned)	Story Points Completed (as on Planned End Date)	Sprint Release Date (Actual)
Sprint-1	20	6 Days	11 Oct 2023	21 Oct 2023	20	21 Oct 2023
Sprint-2	20	6 Days	21 Oct 2023	31 Oct 2023	20	31 Oct 2023

Velocity:

Imagine we have a 10-day sprint duration, and the velocity of the team is 20 (points per sprint). Let's calculate the team's average velocity (AV) per iteration unit (story points per day)

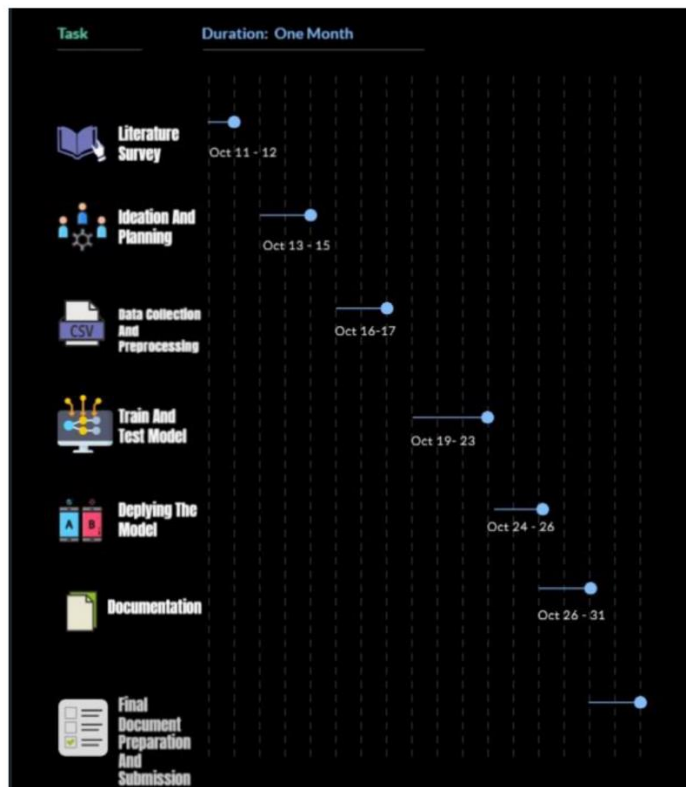
$$AV = \frac{\text{sprint duration}}{\text{velocity}} = \frac{20}{10} = 2$$

Velocity of team = 20 per sprint

Average Velocity = 20/10 = 2 story points per day

Burndown Chart:

A burn down chart is a graphical representation of work left to do versus time. It is often used in agile software development methodologies such as Scrum. However, burn down charts can be applied to any project containing measurable progress over time.



CODE SOLUTION:

1)PORTABLE EXECUTABLE MALWARE DETECTION:

Central to the essence of this project is the application of a machine learning model, specifically the Random Forest classifier tree, as the linchpin for discerning between malicious and benign files. Within the confines of our dataset, a significant majority—constituting 70.1%—comprises malicious files, while benign files constitute the remaining 29.9%.

The process commences with the division of data, adhering to a partitioning of 70% for training and 30% for testing. Herein lies the art of feature selection, a pivotal step in the classification journey. Through the astute employment of the `extratrees.feature_importances_` function, we sifted through the data to identify and isolate the essential features indispensable for accurate classification. This discerning scrutiny, in turn, paved the way for a comparative analysis, ultimately revealing the

supremacy of the Random Forest Classifier, boasting an impressive accuracy score of 99.45%, over its Decision Tree counterpart, which managed a commendable yet slightly lower score of 99.04%.

With the classifier selected, we embarked on the training phase. The model, fortified with its learnings, was enshrined as Classifier.pkl. Simultaneously, the precious insights gleaned about the indispensable features were duly safeguarded as features.pkl, ensuring a ready reference for future engagements.

Post-machine learning, the torch was passed to the file extraction phase, where the extraction of crucial features vital for classification came to the fore. The challenge that loomed large in this endeavour was the meticulous extraction of features embedded within the PE Header files. This formidable task found its solution in the pefile library, a powerful ally in our Python toolkit. Leveraging this resource, we meticulously unveiled the PE Header content, subsequently employing the feature.pkl model to cherry-pick the most pertinent attributes.

The orchestration of feature extraction, carefully choreographed by the pefile library, concluded with the presentation of these handpicked attributes to the Classifier.pkl machine. Here, in this culminating act, the model's astute prediction capabilities were put to the test, unveiling the true nature of the file—be it benign or harbouring malicious intent.

2)URL BASED MALWARE DETECTION:

This segment of the program unfolds in two distinct phases: Data Cleansing, an essential precursor to Logistic Regression, and the subsequent training of the machine to differentiate between malicious and non-malicious entities. At the heart of this intricately woven model lies an inherent need for data that is not only comprehensive but also impeccably precise. The integrity of this dataset is the backbone upon which our model's prowess hinges, comprising a balanced assortment of both benign and malignant URLs.

The path to imbuing our model with comprehension begins with data cleansing, an artistry of its own, where we enlisted the aid of pandas and designed a custom vectorizer. This bespoke tool meticulously purges our datasets, extracting the quintessence required for training. URLs, distinct from conventional text documents, necessitate a specialized sanitization process to distil the relevant data.

The sanctification process unfurls in Python, filtering the URLs to furnish us with pristine datasets. These datasets adhere to a dual-column structure, one for URLs and the other for labels, denoting their malignancy or innocence.

The journey proceeds with the utilization of the Tf-idf machine learning text feature extraction method from the eminent sklearn Python module. To prepare our data for vectorization, we leverage the power of pandas to transform it into data frames and matrices intelligible to our custom vectorizer. The application of the term-frequency and inverse document frequency text extraction methodology follows suit. And it's within this domain that we turn to the trusted Logistic Regression method for training and testing our model.

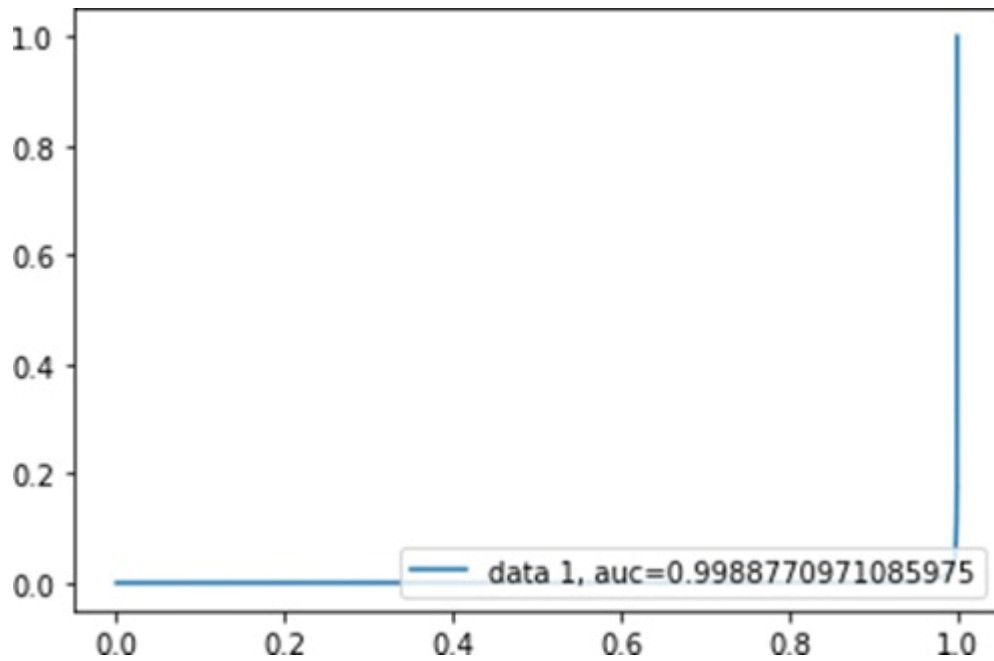
Yet, even with the model's prowess, the discernment of all benign URLs remains an elusive feat. As a strategic countermeasure, we marry our machine learning model with the classical approach of URL filtering, namely the Whitelist Filter. This catalogue comprises known reputable websites that, in all likelihood, pose no harm to our users. The amalgamation of this manual whitelist and machine learning discrimination fortifies our network traffic, ensuring safe passage for the innocuous entities

and thereby enhancing our users' online safety. This harmonious synergy is a testament to the multifaceted approach underpinning our endeavour.

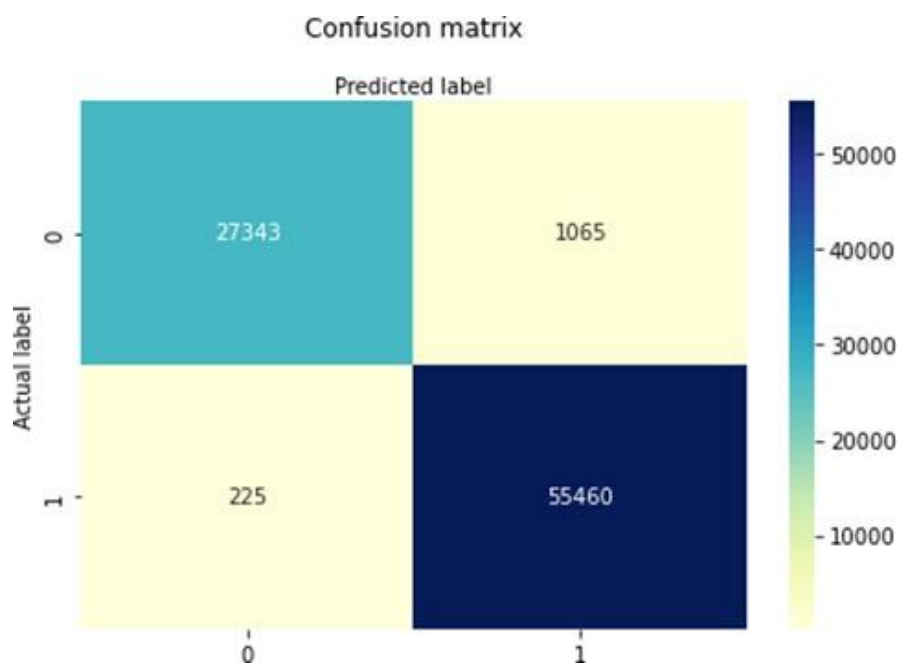
TEST AND VALIDATION

1) URL DETECTION

➤ ROC CURVE

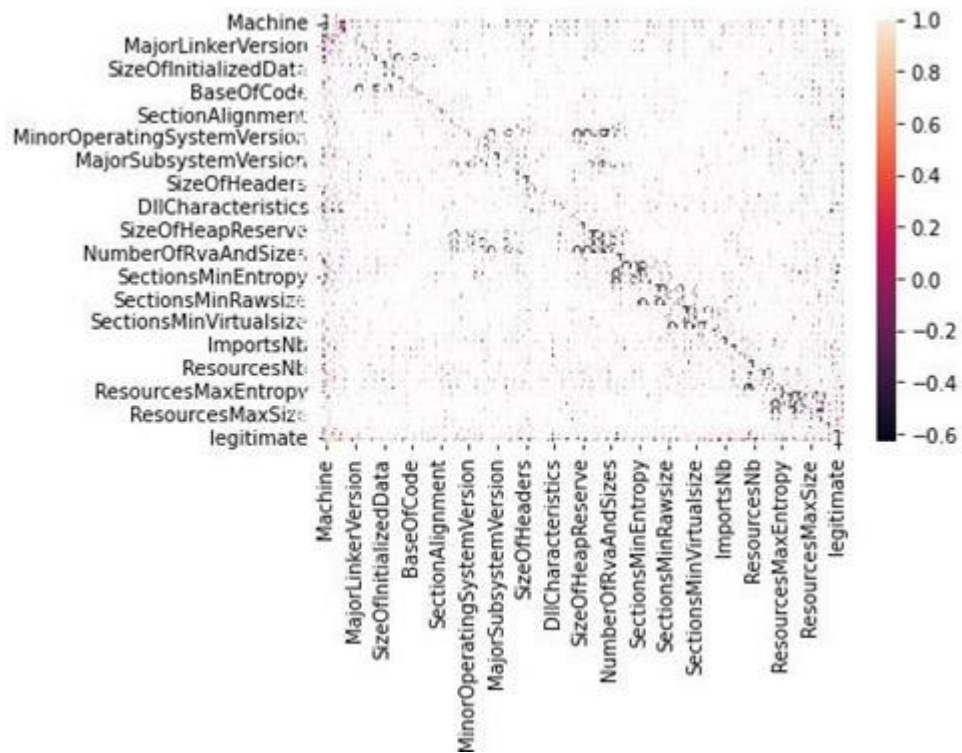


➤ CONFUSION MATRIX



2) PORTABLE EXECUTABLE MALWARE DETECTION

➤ CORRELATION MATRIX



➤ CONFUSION MATRIX



PERFORMANCE

The Random Forest Classifier, a stalwart in the realm of machine learning, flaunts an enviable accuracy rate of 99.37%. Its precision in discerning between benign and malicious entities is a testament to its formidable capabilities. In the vigorous pursuit of cybersecurity, the Logistic Regression model stands as a robust contender, boasting a commendable accuracy rate of 98.46%. However, its precision takes centre stage, a remarkable 99.18% - a testament to its adeptness in distinguishing the innocent from the malevolent.

Furthermore, the Logistic Regression model's recall rate, an important metric that gauges its ability to identify all malicious entities in the dataset, shines at a formidable 96.25%. This capacity to recover malicious instances is crucial in maintaining the fortitude of our security infrastructure.

RESULTS

```
C:\Users\21A\MALWARE-DETECTION\Others>python main.py
'##:::'##:::'##:::'##:::'##:::'##:::'##:::'##:::'#####:#####:
#####:## ##::: ##::: ##::: ##:::##: ##:::## ##::: ##... ##: ##:::
#####:##::: ##: ##: ##::: ##: ##: ##:::##::: ##: ##: ##:::
## ## ##:##::: ##: ##: ##::: ##: ##: ##:::##::: ##: #####: #####:
##: #: ##: #####: ##::: ##: ##: ##: #####: ##... ##::: ##...
##::: ##: ##... ##: ##::: ##: ##: ##: ##... ##: ##::: ##: ##:::
##::: ##: ##::: ##: #####: ##: ##: ##: ##: ##::: ##: #####:
.....:#####:#####:#####:#####:#####:#####:#####:#####:
##... ##: ##::: ##::: ##::: ##:::##... ##... ##::: ##::: ##:
##::: ##: ##::: ##::: ##::: ##::: ##::: ##::: ##: ##: ##:
##::: ##: #####: ##::: #####: ##::: ##: ##: ##: ##:
##::: ##: ##... ##::: ##::: ##::: ##::: ##: ##: ##: ##:
##::: ##: ##::: ##::: ##::: ##::: ##::: ##::: ##: ##: ##:
#####: #####: ##::: #####: #####: ##::: ##::: #####:
.....:#####:#####:#####:#####:#####:#####:#####:
'##::: ##:
###: ##:
###: ##:
## ## ##:
##. ###:
##.. ###:
##::: ##:
.....:

TEAM: 2.4

COURSE: AI FOR CYBERSECURITY WITH IBM QRADAR

MEMBERS: SAIRAM B | MRIDHULA S | SHIVANI NARAYAN N | VEEKSHITHA

SUPERVISOR: MANOJ

1. PORTABLE EXECUTABLE FILE scanner
2. URL scanner
3. Exit

Enter your choice : _
```

```
Terminal: Local x + v
##:: ##:
.....:

TEAM: 2.4

COURSE: AI FOR CYBERSECURITY WITH IBM QRADAR

MEMBERS: SAIRAM B | MRIDHULA S | SHIVANI NARAYAN N | VEEKSHITHA

SUPERVISOR: MANOJ

1. PORTABLE EXECUTABLE FILE scanner
2. URL scanner
3. Exit

Enter your choice : 2
Input the URL that you want to check (eg. google.com) : www.certifiedhacker.com

The entered domain is: bad
Do you want to search again? (y/n)
```

```
Enter your choice : 2
Input the URL that you want to check (eg. google.com) : google.com

The entered domain is: good
```

CONCLUSION

In summation, the project embarked upon the formidable journey of "Malware Detection and Classification Using Machine Learning," a pivotal chapter in the ever-evolving saga of cybersecurity. Drawing from an extensive wellspring of knowledge, innovative methodologies, and unwavering analytical rigor, this endeavour stands as a sentinel, tirelessly addressing the omnipresent specter of malware within the intricate tapestry of the digital realm.

In a symphony of algorithms, the project orchestrates an unparalleled harmony of machine learning, fortifying its arsenal against the malevolent forces that threaten our digital sanctuaries. Through this orchestration, a symphony of remarkable accuracy emerges, capable of detecting and classifying an array of malware archetypes, elevating the citadels of computer system security.

The project's revelations reverberate through the corridors of cybersecurity, echoing the profound significance of machine learning methodologies in the face of ever-adapting malware threats. It exemplifies the potent marriage of advanced algorithms and the precision of feature selection, effectively unmasking the concealed patterns and sinister conduct of malevolent software.

Beyond its immediate impact, this research extends a benevolent hand to industry luminaries and policymakers, imparting the wisdom required to construct formidable cyber fortifications. It charts a course for future explorations, beckoning the foray into even more intricate machine learning models, arming us against cyber adversaries and paving the road to a more secure and resilient digital frontier.

Ultimately, this project is more than a culmination; it signifies the dawn of a new era, one where the battle against malware prevails, securing the digital destiny of individuals, enterprises, and global organizations alike.

FUTURE SCOPE

1)Enhancing the Project with a Comprehensive Dataset:

- ✓ Leveraging a Wider, Well-Labeled Dataset: To fortify the project's capabilities, we propose the incorporation of a more extensive and meticulously labeled dataset. A comprehensive dataset enriches the model's learning process, enabling it to make more accurate distinctions between benign and malicious entities. This expansion paves the way for a more robust and versatile machine learning model.

2)Broadening Accessibility and Convenience:

- ✓ Web-Based Accessibility: A notable upgrade to the project involves migrating it to the web, allowing users to access it remotely. Hosting the project online opens the door to a broader user base, enhancing accessibility for individuals and organizations alike.
- ✓ User-Friendly File Upload Feature: The integration of a user-friendly file upload feature provides a seamless experience for users who wish to have their files scanned for potential threats. By simplifying the process of uploading files, the project becomes more approachable and convenient.
- ✓ Inclusion of URL Detection: In a digital landscape where URLs play a pivotal role, integrating a URL detection feature becomes imperative. Users can input web addresses to assess their trustworthiness, adding a layer of protection against potentially harmful websites.

3)User Interface Revamp:

- ✓ Windows GUI: While the initial project was console-based, an evolution toward a graphical user interface (GUI) for the Windows platform promises an improved user experience. A GUI is more intuitive, visually appealing, and user-friendly, making it accessible to a wider audience.

4)Real-Time Scanning for Everyday Use:

- ✓ Real-Time File Scanning: The project's scope can be further extended to include real-time scanning capabilities. This enhancement allows the system to monitor files during the downloading and transferring process, proactively detecting and mitigating potential threats. This real-time scanning feature can be seamlessly integrated into daily computer usage scenarios, providing a layer of security that operates in the background without disrupting user activities.

5)The Power of Adaptation:

- ✓ Staying Ahead of Evolving Threats: As the cybersecurity landscape evolves, the project is positioned to adapt and evolve alongside it. By incorporating these proposed enhancements,

the project transforms into a versatile and proactive tool that can safeguard individuals and organizations against emerging cyber threats.

- ✓ Empowering Users: The expanded project provides users with a comprehensive and accessible solution for mitigating digital risks. It acts as a guardian, protecting against malware, malicious files, and suspicious URLs, ensuring that users can navigate the digital realm with confidence.

In essence, these proposed improvements mark the project's transition from a localized and terminal-based application to a sophisticated and comprehensive cybersecurity solution. With a more extensive dataset, web-based accessibility, user-friendly features, real-time scanning, and a refined user interface, the project stands as a beacon of digital security in an ever-evolving online landscape. It is not just a defence mechanism; it is an empowerment tool for those navigating the complexities of the digital world.

APPENDIX

- ❖ Project documentation link: <https://github.com/smartinternz02/SI-GuidedProject-587767-1697075897>
- ❖ Project executable files link: <https://github.com/SairamBalamurugan087/MALWARE-DETECTION>