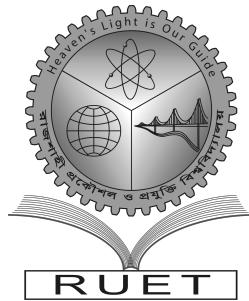


Heaven's Light is Our Guide



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

Rajshahi University of Engineering & Technology, Bangladesh

Vision Based Malware Classification Framework Based on Neural Network

Author

Shah Ahmed Saad Rupai

Roll No. 1703069

Department of Computer Science & Engineering
Rajshahi University of Engineering & Technology

Supervised by

Prof. Dr. Md. Ali Hossain

Professor

Department of Computer Science & Engineering
Rajshahi University of Engineering & Technology

ACKNOWLEDGEMENT

I would first like to express my gratitude to Allah Ta'ala for providing us with the chance and motivation to finish our thesis work.

We want to sincerely thank, appreciate, and respect our supervisor, Prof. Dr. Md. Ali Hossain, Professor of the Department of Computer Science and Engineering, Rajshahi University of Engineering and Technology, Rajshahi. He has continuously encouraged us, offered us advise, helped us, and cooperated sympathetically with us whenever he felt like it throughout the year in addition to providing us with the technical instructions, guidance, and documentation we needed to finish the task. The most effective tool for helping us attain our goal was his ongoing assistance. He was there for us at any time of day if we ever got stuck in a complex problem or circumstance. Without his genuine concern, this work would not have taken the current form that it does.

Additionally, I would like to express my gratitude to all of the instructors at Rajshahi University of Engineering and Technology's Computer Science and Engineering departments for their time, effort, and insightful advice.

Last but not least, I want to express my gratitude to my parents, friends, and well-wishers for their ongoing inspiration and numerous helpful suggestions throughout this endeavor.

August 14, 2023

RUET, Rajshahi Shah Ahmed Saad
Rupai

Shah Ahmed Saad Rupai

Heaven's Light is Our Guide



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING

Rajshahi University of Engineering & Technology, Bangladesh

CERTIFICATE

*This is to certify that this thesis report entitled “**Vision Based Malware Classification Framework Based on Neural Network**” submitted by **Shah Ahmed Saad Rupai, Roll:1703069** in partial fulfillment of the requirement for the award of the degree of Bachelor of Science in Department of Computer Science & Engineering of Rajshahi University of Engineering & Technology, Bangladesh is a record of the candidates’ own work carried out by them under my supervision. This thesis has not been submitted for the award of any other degree.*

Supervisor

External Examiner

Prof. Dr. Md. Ali Hossain

Professor

Department of Computer Science &
Engineering
Rajshahi University of Engineering &
Technology
Rajshahi-6204

Department of Computer Science &
Engineering
Rajshahi University of Engineering &
Technology
Rajshahi-6204

ABSTRACT

In the field of cybersecurity, developing strong and efficient methodologies for malware detection categorization is a crucial challenge due to the constant evolution of malware and the rising cyber threats. The rapid creation of innovative and polymorphic malware variants is a challenge for conventional approaches, which frequently rely on signatures and behaviors. The unique vision-based malware classification methodology presented in this thesis uses computer vision to evaluate and classify malware samples. Traditional methods for classifying malware frequently rely on manually created features and rule-based techniques, which can have limitations in their capacity to recognize the complicated patterns and variations seen in complex malware pictures. These techniques might have a hard time keeping up with the constantly changing and new malware varieties, which would decrease their accuracy and raise the number of false positives. Furthermore, the manual process of feature engineering might not be scalable enough to handle the enormous amount of data produced by contemporary malware attacks and might be time-consuming. Contrarily, neural network-based methods have the potential to automatically identify pertinent features from unprocessed data, resulting in malware classification models that are more reliable and adaptive and can distinguish between benign and harmful samples more effectively. This framework's main objective is to treat malware binaries differently by treating them as visual things. This innovative viewpoint enables the implementation of well-established image processing methods for thorough feature extraction and classification.

The proposed models showed classification accuracy of 97.57% (Convolutional Neural Network), 97.49% (Sparse Principal Component Analysis & Multi Layer Perceptron) and 97.13% (Principal Component Analysis & Multi Layer Perceptron). Concolutional Neural Network showed best accuracy.

Key words: Malware Classification, Computer Vision, Vision-Based Framework, Image Processing

CONTENTS

	Pages
ACKNOWLEDGEMENT	ii
CERTIFICATE	iii
ABSTRACT	iv
CHAPTER 1 Introduction	1
1.1 Introduction	1
1.2 Problem Identification	3
1.3 Challenges	4
1.4 Motivation	5
1.5 Goals	6
1.6 Research Question	7
1.7 Research Objectives	7
1.8 Thesis Organization	8
1.9 Conclusion	9
CHAPTER 2 Background Study and Literature Review	10
2.1 Introduction	10
2.2 Machine Learning	11
2.2.1 Supervised Learning	12
2.2.2 Neural Network	13
2.3 Literature Review : Traditional Approaches & Impact of Neural Network on Vision Based Malware Classification	14
2.4 Conclusion	17
CHAPTER 3 Methodology	18
3.1 Introduction	18

3.2	Overview of the Proposed System	19
3.3	Data Collection & Description	19
3.4	Data Analysis	20
3.4.1	Data Set Specification	21
3.5	Data Preprocessing	22
3.6	Feature Extraction	23
3.7	Proposed Models	24
3.7.1	Supervised Machine Learning	25
3.7.2	Convolutional Neural Network (CNN)	26
3.7.3	Multi-Layer Perceptron (MLP)	27
3.8	Conclusion	28
CHAPTER 4 Vision Based Malware Classification Using Neural Network		29
4.1	Introduction	29
4.2	Implementation Tools	29
4.3	Detailed Data Analysis	30
4.3.1	Data Set Information	30
4.4	Feature Extraction: Principal Component Analysis (PCA)	32
4.5	Feature Extraction: Sparse Principal Component Analysis (SPCA)	33
4.6	Machine Learning Models Implementation	35
4.6.1	Convolutional Neural Network (CNN)	35
4.6.2	Multi Layer Perceptron (MLP)	36
4.7	Conclusion	36
CHAPTER 5 Result and Performance Analysis		37
5.1	Introduction	37
5.2	Data Analysis and Preprocessing	37
5.2.1	Feature Extraction	38
5.3	Performance Evaluation Metrics	39
5.4	Results of the classification models:	41
5.4.1	Convolutional Neural Network (CNN)	41
5.4.1.1	Figure of confusion matrix :	42
5.4.1.2	Precision of each class :	42

5.4.2	Principal Component Analysis (PCA) & Multi Layer Perceptron (MLP)	43
5.4.2.1	Figure of confusion matrix :	43
5.4.2.2	Figure of ROC curve:	44
5.4.2.3	Precision of each class :	44
5.4.3	Sparse Principal Component Analysis (SPCA) & Multi Layer Perceptron (MLP)	45
5.4.3.1	Figure of confusion matrix :	45
5.4.3.2	Figure of ROC curve:	46
5.4.3.3	Precision of each class :	46
5.5	Conclusion	47
CHAPTER 6 Conclusion and Future Works		48
6.1	Introduction	48
6.2	Thesis Summary	48
6.3	Limitations	50
6.4	Future Works	51
6.5	Conclusion	51
REFERENCES		52

LIST OF TABLES

Sl	Table Name	Pages
5.1	Performance Analysis of my work	41

LIST OF FIGURES

Sl	Figure Name	Pages
1.1	Types of Malware	2
2.1	Areas of machine learning fields	12
2.2	Neural Network	14
3.1	Work flow Proposed System	19
3.2	Malware Image	20
3.3	Process of Creating Malware Image From Malware Binaries	21
3.4	25 Different Families of Malware Images	22
3.5	Machine Learning model's workflow	25
3.6	Supervised Machine Learning model's workflow	26
3.7	Convolutional Neural Network (CNN) Model	27
3.8	Multi-Layer Perceptron Model	28
4.1	Folders of 25 Families From Dataset	31
4.2	Malware Images From Single Folder	31
4.3	PCA Explained Visually	33
4.4	PCA-Sparse PCA Explained Visually	34
5.1	Visualization of Data Using First Two Principal Components	38
5.2	Confusion Matrix	39
5.3	Confusion matrix of CNN Model	42
5.4	Precision Of Each Class	42
5.5	Confusion matrix of PCA+MLP Model	43
5.6	ROC Curve of PCA+MLP Model	44
5.7	Precision Of Each Class	44
5.8	Confusion matrix of SPCA+MLP Model	45
5.9	ROC Curve of SPCA+MLP Model	46
5.10	Precision Of Each Class	46

Chapter 1

Introduction

1.1 Introduction

The cybersecurity environment is a key area of conflict in the digital age, as technology affects every aspect of contemporary life. Malware stands out from the wide range of digital dangers as a chameleon-like foe that is constantly evolving and adapting to take advantage of weak spots in the complex web of interconnected systems. This thesis goes deeply into the complex world of malware, examining its multidimensional makeup and the wide-ranging effects it has on people, businesses, and the global digital economy.

The term "malware," a portmanteau of "malicious software," refers to a variety of programs created with malicious purpose. Its forms are as varied as they are sneaky: from viruses and worms that multiply themselves over networks to trojans that hide themselves behind ostensibly harmless software to ransomware that keeps crucial data hostage in exchange for money. These hostile organisms thrive in the remote areas of the internet, finding weaknesses with the surgical accuracy of a scalpel.

Malware's impacts go far beyond the boundaries of ones and zeros. A malware assault that is effective can create a path of destruction and undermine the tenets of trust, data integrity, and privacy. Theft of personal information, financial fraud, and identity theft target specific people. No organization, regardless of size, is immune to the threat of cyberattacks. Governments and key infrastructure that were formerly thought to be resistant are now at risk from cyberattacks that cross international borders.

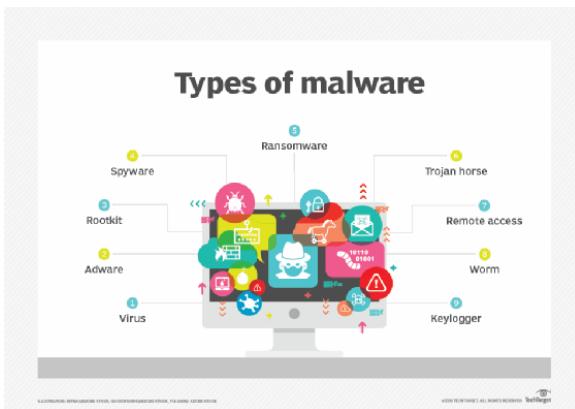


Figure 1.1: Types of Malware [1]

The methods and strategies used by those who create malware change along with the digital environment. One method used to get over conventional protections is polymorphism. Other methods include obfuscation and social engineering. Cybersecurity experts that work to remain ahead of the curve face a significant challenge due to the malware industry's continuous pace of innovation.

The dynamic and constantly changing nature of current software presents difficulties for traditional techniques to malware classification. These techniques mostly rely on predetermined criteria and signatures, which are prone to becoming obsolete as attackers use new strategies to avoid detection. High false negative rates might result from the failure to respond fast to new threats, allowing malicious software to go undetected. Traditional methods might also have trouble addressing the wide variety of file formats and obfuscation strategies employed by virus developers to avoid detection. As a result, these techniques could not be scalable enough or precise enough to effectively counter the landscape of malware attacks, which is becoming more and more sophisticated.

In the ever-evolving landscape of cybersecurity, innovative approaches are essential to tackle the relentless growth of malware variants. One such pioneering approach is vision-based malware classification, which harnesses the power of computer vision to analyze the visual characteristics of malware samples. By transforming binary code into visual representations, this methodology unlocks new dimensions for analysis. Through the lenses of image processing and pattern recognition, this technique seeks to unveil unique visual patterns that differentiate

benign software from malicious code. This introduction sets the stage for a deep dive into the realm of vision-based malware classification, where the fusion of visual data and cybersecurity promises to reshape our strategies against digital threats[2].

1.2 Problem Identification

The struggle against malware is becoming more and more intense in the constantly changing cybersecurity scene. In all of its manifestations, malicious software has demonstrated extraordinary flexibility, consistently outperforming conventional detection techniques. The main issue is that malware programmers are increasingly able to alter the code, making signature-based and behavior-based detection systems ineffective. Furthermore, the sheer volume of malware samples is growing exponentially, necessitating novel methods that can quickly and reliably classify and categorize these threats[3].

The question of how to extend malware detection beyond normal bounds arises as a compelling problem to address these difficulties. The larger issue of keeping up with quickly evolving malware strains while accommodating the expanding scope of the digital threat ecosystem is captured by the question. A fascinating approach that aims to alter how we view and evaluate malware is vision-based malware detection. We seek to extract novel insights that can evade conventional approaches by examining malware binaries as visual patterns. The objective is to provide a fresh and efficient method that can not only recognize known malware but also potentially identify new threats with greater accuracy.

Given how dynamic and constantly changing current malware is, traditional approaches to malware classification sometimes struggle to handle it properly. These techniques mostly rely on predetermined criteria and signatures, which can quickly become out of date as attackers use new strategies to avoid detection. High false negative rates brought on by a failure to respond rapidly to new threats can let dangerous software slip through the cracks. Furthermore, due to the variety of file formats and obfuscation strategies employed by malware developers to avoid detection, standard methodologies could have trouble addressing them. Because of this, these techniques could not be scalable or precise enough to effectively counter the variety of malware

attacks that are becoming more complex. Additionally, manual feature engineering is frequently a major component of older techniques, which can be time-consuming and may not fully capture the range of virus behaviors. Malware classification mistakes or missed detections may result from a lack of a comprehensive grasp of malware dynamics and behaviors. Additionally, conventional approaches frequently concentrate on a particular element of malware while ignoring possible synergies between many traits. This constraint may limit their ability to reliably distinguish between legitimate and harmful software across a variety of aspects, which may reduce their effectiveness in combating contemporary, polymorphic malware threats.

This issue deals with the urgent need for cybersecurity resilience, making it more than merely an intellectual challenge. The effects of malware attacks range from individual privacy violations to global disruptions as society becomes more computerized and linked. Our security tactics could be revolutionized by finding a solution to the vision-based malware detection issue, which would offer an essential protection against the unrelenting cleverness of malware writers. To create the foundation for a safer and more secure digital future, computer vision, machine learning, and domain knowledge must come together.

An integrated and interdisciplinary approach is necessary to understand the complexity of vision-based malware detection. This goal involves the creation of novel algorithms, frameworks, and approaches, drawing on knowledge in computer vision, machine learning, software analysis, and cybersecurity. The end goal is to create a paradigm shift that transforms our security mechanisms and makes it possible to proactively identify and reduce malware threats.

1.3 Challenges

Multiple obstacles must be carefully considered in the realm of vision-based malware categorization. First off, developing reliable and generalizable models is essential due to the diversity and dynamic nature of malware strains. The conversion of malware binaries into images is a tough challenge. Another difficulty is ensuring scalability and efficiency while dealing with enormous datasets and intricate architectures. Additionally, addressing the interpretability of deep learning models as well as coping with unbalanced class distributions are crucial concerns.

Furthermore, the accuracy of deep learning models employed for malware classification may be compromised by adversarial attacks. Last but not least, to keep ahead of emerging threats, the ongoing evolution of malware tactics necessitates adaptive and dynamic models. For vision-based malware classification systems to be successfully implemented, certain obstacles must be overcome.

1.4 Motivation

The urgent need to close the rising gap between the effectiveness of conventional detection techniques and the growing arsenal of malware is what spurred researchers to explore the field of vision-based malware detection. The inadequacies of conventional methods have been made clear by the emergence of complex malware variants and their capacity to adapt and circumvent conventional protections. The subject of how to defend our digital environments against these dangers grows urgent and more pertinent as malware becomes more polymorphic and evasive.

Vision-based malware detection has the ability to completely change how we perceive malware behavior, which is what makes it so alluring. We overcome the constraints of syntax manipulation and delve into the distinctive characteristics that arise in the visual representation of malware binaries by shifting the focus from code-based signatures to visual patterns. This approach offers a rare chance to identify complex variants and evasive tactics used by malware writers—patterns that could otherwise elude conventional detection approaches.

The goal of adopting a preventative strategy for cyber security is included in the broad motivation. In order to defend against attacks that are fundamentally hard to forecast, vision-based malware detection offers a revolutionary method that holds the promise of anticipatorily identifying future threats. The capacity to recognize malware via the prism of visual patterns ushers in a new era of detection—one that is in line with the inventive and dynamic nature of malware and provides defenders with cutting-edge tools to combat these always changing cyber threats.

The motivation also encompasses the ability to adapt. The field of creating malware is always developing along with the digital environment. A cyber security paradigm that can adapt in

real-time is required due to the usage of obfuscation techniques, polymorphism, and innovative attack vectors. By basing its principles on visual patterns, vision-based malware detection shows promise in its ability to adapt to new malware variants and families, resulting in a strong and resilient protection system.

In the end, the drive to investigate malware detection which is vision-based comes from a shared desire to protect our interconnected digital world. It draws on the collective wisdom of cyber security practitioners, researchers, and specialists who work to outsmart bad actors' cunning. This topic's study is in line with the overarching objective of strengthening our cyber defenses through creativity, cooperation, and a never-ending quest for knowledge in the rapidly developing field of cyber security.

The primary objective of this thesis is to train our model with the dataset by addressing all of the complexities that arise in the process. This will make the model more resilient to challenging circumstances and enable it to produce satisfactory results regardless of the circumstances.

1.5 Goals

The pursuit of vision-based malware categorization is supported by a wide range of objectives, many of which have the potential to completely alter the cyber security landscape. These objectives include the construction of a cutting-edge classification paradigm that deciphers malware through visual patterns by fusing computer vision and image analysis. By improving evasion detection, addressing polymorphic malware, enabling early threat identification, visualizing malware behavior, encouraging interdisciplinary collaboration, ensuring the practical implementation of solutions, and adding to the larger body of research, this approach seeks to go beyond the limitations of conventional techniques. By adopting these goals, we set out on a revolutionary path to strengthen our defenses against the constantly evolving dangers in the digital sphere and promote a more secure and resilient cyber environment.

1.6 Research Question

This study aims to address the following research questions:

- (a) What are the limits of the conventional heuristic- and signature-based malware classification methods?
- (b) How can important elements from malware photos be extracted using computer vision techniques to increase classification accuracy?
- (c) What conclusions may be drawn from contrasting the effectiveness of the suggested vision-based neural network approach with conventional techniques using real-world malware datasets?

1.7 Research Objectives

The research projects in the area of computer vision-based malware classification are directed by a number of specific goals that together work to advance our understanding of malware detection and classification. These goals include both smaller technical goals and larger strategic ones :

1. Develop visual feature extraction techniques.
2. Creation of stable classification models.
3. Performance evaluation of detection.

By addressing these research goals, the study hopes to advance our understanding of emerging malware classification techniques and reveal how vision-based techniques combined with neural networks may be used to advance cybersecurity. Through these research objectives, we embark on a journey to unlock the potential of vision-based malware classification, envisaging a future where the fusion of computer vision and cybersecurity transforms our ability to combat evolving digital threats.

1.8 Thesis Organization

Following is how the remaining chapters are organized:

Chapter 2 : Background Study And Literature Review

In this chapter, the contributions, methodologies, and limits of the works related to the categorization of malware based on various methods are discussed. This article discusses the background of vision-based malware classification, the importance of classification models, and related topics.

Chapter 3 : Methodology

This section gives an overview of the technique that was used to write the thesis as well as the suggested method designs.

Chapter 4 : Implementation

The application of malware categorization, along with data preprocessing and machine learning, are all covered in this chapter. The topic of neural network models is also covered. Additionally, this chapter provides examples of the data set description, system architecture, and implementation.

Chapter 5 : Result and Performance Analysis

In this chapter, the outcomes of our tests and a performance analysis of the model are discussed. This section also includes a description of the metrics that were employed to assess our model.

Chapter 6 : Conclusion and Future Works

The thesis project is completed in this chapter. Its flaws are explored, and a course for further research is suggested.

1.9 Conclusion

The confluence of malware categorization and computer vision in the developing field of cybersecurity presents a promising paradigm change in our protection against changing cyber threats. Given the complexity of malware and the limits of existing detection techniques, the introduction of vision-based malware classification highlights the urgency of this goal. The first step of our journey is to reveal the hidden potential of visual patterns found within malware files. This journey combines cutting-edge technology, cross-disciplinary cooperation, and steadfast dedication to a safer digital space. The following parts will analyze this ground-breaking tactic in detail, exposing methods, difficulties, and revelations that shed light on the way forward in the ever-changing field of cybersecurity.

Chapter 2

Background Study and Literature Review

2.1 Introduction

The importance of having a thorough awareness of the corpus of existing information increases in the dynamic field of cyber security where innovation and adaptation are constants. An important pillar in the investigation of vision-based malware categorization is the part on background research and literature reviews. The goal of this section is to explain the entire panorama of pertinent research, approaches, and technological developments that have prepared the way for the convergence of malware analysis with computer vision. We learn about the development of malware detection over time, the advent of computer vision in cyber security, and the complexities of visual pattern recognition by digging into the annals of literature. This preparation sets the stage for our adventure and offers a strong basis for the unique vision-based malware categorization method. The outlines of our individual contribution to this dynamic and essential junction of areas are ultimately shaped by our exploration of the synergies and gaps in this rich tapestry of study.

Malware attack statistics highlight the alarming truth of the vulnerabilities of the digital age. Cyber attacks have increased exponentially in recent years, and many different types of malware are to blame. According to reports, millions of malware samples are produced each year, illustrating the size of the threat environment. Particularly ransomware assaults have seen a rapid rise, with both people and businesses being targets of data extortion. The attack surface has also increased as a result of the growth of mobile devices and the Internet of Things (IoT), opening up new channels for malware distribution. These figures demonstrate the widespread

effects of cyberattacks, including monetary losses, data breaches, and business interruptions across all sectors. Understanding these numbers emphasizes the essential need for cutting-edge methods, like vision-based malware classification, to bolster our defenses against these relentless digital enemies. Malware is evolving in sophistication and diversity, and this trend will only continue.[4].

The process of looking for algorithms that can infer general hypotheses from examples provided by an external source is known as supervised machine learning. These general hypotheses can then be applied to predictions about brand-new examples. In the context of various Supervised Machine Learning (ML) classification methodologies, compares and contrasts a variety of supervised learning approaches and decides the most efficient classification algorithm based on the data set, the quantity of instances, and the variables (features).[5].

2.2 Machine Learning

A key component of modern artificial intelligence, machine learning represents a paradigm change in computer problem-solving. Machine learning goes beyond conventional rule-based programming since it is based on the idea that computers can automatically learn from experience and get better over time. It includes a wide range of methods that let algorithms to evaluate data, spot patterns, and come to wise judgments without having to be explicitly programmed. Supervised learning, in which algorithms learn from labeled examples to generalize and predict results for unobserved data, is at the core of it. Data's inherent structures are uncovered by unsupervised learning, exposing hidden linkages and groups. Reinforcement learning guides agents through dynamic environments as they experiment and learn to maximize rewards.

It is clear that machine learning has a significant impact across industries. By evaluating medical images and projecting patient outcomes, it helps in the diagnosis of diseases in the healthcare industry. It is used by financial firms to spot fraudulent transactions and forecast market movements. Autonomous vehicles are changing transportation, a feat made possible by machine learning's capacity to decipher challenging real-time data. Systems for entertainment and recommendations make use of this capability to curate content and customize user experiences.

But there are obstacles in the way of the technology's potential. The significance of responsible implementation is underlined by ethical considerations, data privacy, algorithm bias, and the requirement for interpretable models. As machine learning's capabilities develop, multi-disciplinary cooperation between computer scientists, statisticians, subject matter experts, and ethicists is essential to maximizing its advantages while minimizing hazards.

In conclusion, machine learning has transformed how we engage with data, producing insights, automation, and innovation that were previously thought to be unachievable. It has the potential to disrupt industries, spur scientific advancements, and alter the very foundation of our modern life as it moves forward.

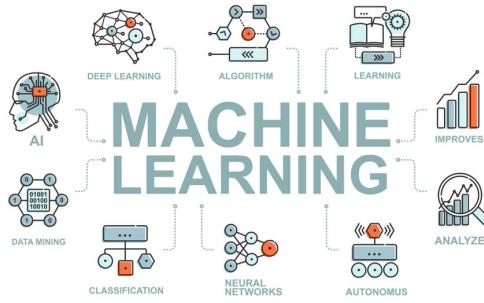


Figure 2.1: Areas of machine learning fields [6]

2.2.1 Supervised Learning

Through the route of supervised learning, a spectacular trip in the field of machine learning emerges. By helping computers to learn, adapt, and make defensible conclusions from data, this fundamental technique has cemented its place as a key factor in the development of predictive analytics. Supervised learning reveals the potential to endow machines with the astonishing capacity to recognize and generalize patterns, turning raw data into useful insights. Applications range from analyzing the sentiment of text to diagnosing medical diseases.[7]. In its most basic form, supervised learning entails the instruction of algorithms using labeled datasets. In essence, these datasets educate the algorithm the fundamental relationship between the two by

giving it pairs of inputs and associated intended outcomes. Imagine a sizable collection of pictures, each one painstakingly identified as containing a dog or a cat. The computer carefully examines these photographs, recognizes the visual clues that distinguish feline from canine species, and eventually masters the ability to distinguish between cats and dogs in obscure images. The versatility of supervised learning across domains is what makes it so attractive. It aids in the diagnosis of diseases using medical data in healthcare while forecasting stock values and identifying fraudulent transactions in finance. It transforms client interactions by translating languages and interpreting natural language to power chatbots. The field of autonomous vehicles thrives on supervised learning, where algorithms pick up the finer points of precise road navigation from human drivers. Through the use of supervised learning, machines are able to comprehend the complexities of our environment and serve as a light of predicting prowess. Its uses are only limited by our creativity, transforming industries, advancing research, and redefining the idea of making informed decisions. As technology advances hand in hand with supervised learning, its path becomes a monument to our effort to develop computers that are not only smarter but also profoundly wiser.

2.2.2 Neural Network

Computer simulations of brain neurons are called neural networks. Layers are made up of linked nodes or neurons. Each neuron takes in inputs, weighs them, and then activates the output that results from that process. Neural networks are used in artificial intelligence and machine learning to discover patterns and relationships in data. MLPs are models of neural networks. Input, hidden, and output layers are present. Each input layer neuron's feature value is propagated by the network. In the hidden layers, where the majority of calculations take place, weighted connections and activation functions modify the input. The output layer classifies or predicts things. MLPs are used to map input data to output labels during supervised learning. Backpropagation is used to optimize the weights and biases of neuron connections, and these parameters are changed to reduce the discrepancy between expected and actual outputs when training an MLP. MLPs are helpful for image classification, natural language processing, and regression because they can learn complex non-linear data relationships. Even though they have contributed to the development of more sophisticated neural network topologies like CNNs and RNNs, MLPs remain crucial in artificial neural networks.

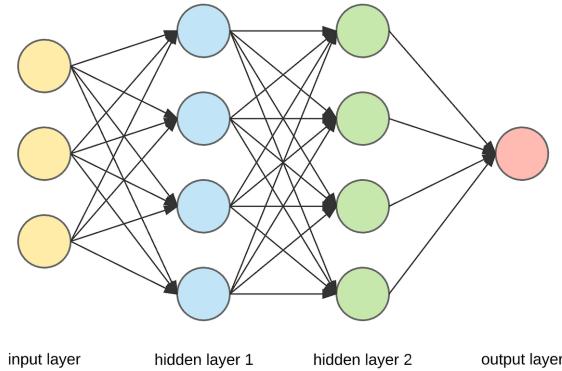


Figure 2.2: Neural Network [8]

2.3 Literature Review : Traditional Approaches & Impact of Neural Network on Vision Based Malware Classification

In order to effectively defend against the constantly changing world of digital threats, malware detection, a cornerstone of cyber security, employs a wide range of methodologies [9]. Heuristic analysis finds suspicious behaviors that can point to the existence of malware, in contrast to signature-based detection, which uses established patterns to identify known malware strains. Dynamic anomaly detection methods find departures from typical system behavior and potential attacks. While behavior-based detection keeps track of software activities to spot malicious activity, sandboxing executes files in closed contexts to watch their behavior. Artificial intelligence assists in recognizing complex and developing threats, while machine learning identifies patterns from massive datasets to find anomalies. On the other hand, intrusion detection systems monitor network traffic to look for illegal access and the spread of malware. Together, these strategies create a multifaceted defense that exemplifies the dynamic character of the continuing conflict between cyber security experts and malicious actors.

Signature-based malware detection and classification is characterized by its simplicity and efficiency [10]. It involves matching software or code against a predefined set of known signatures, enabling swift identification of familiar malware strains. This approach is highly effective for well-established threats, providing rapid responses to known attacks. Signature-based

methods also yield low false positive rates, as the signatures are specific to known malicious patterns. This strategy, nevertheless, has several serious drawbacks. Due to their difficulty in detecting novel or changing malware strains, signature-based techniques are useless against zero-day assaults and polymorphic malware that often modifies its source code. There is a latency between the appearance of a new threat and its detection since signatures need to be updated often. The use of evasion strategies like code obfuscation, in which attackers alter malware to avoid signature detection, might likewise be used against this strategy. If the malware's signature is not in the database, the reliance on predefined signatures leaves signature-based systems open to false negatives.

In comparison to signature-based techniques, heuristic-based malware detection [11] and categorization offer a more flexible approach. These methods are able to recognize malware strains that were previously unidentified by examining program behavior and features. This adaptability is particularly useful against polymorphic malware and zero-day assaults, which signature-based techniques have trouble identifying. Heuristic approaches are improved in their capacity to recognize changing threats by their ability to detect obfuscated code or variants of recognized malware. Heuristic-based approaches have their own set of difficulties, though. Due to their reliance on spotting patterns that mimic malevolent conduct, they risk producing false positives. Heuristics may overlook novel attack vectors or modifications that don't fit preset patterns because they are built on assumptions and rules. Additionally, as attackers develop new tactics to get around heuristic rules, these techniques must always be improved. Expert expertise is required to create precise and thorough heuristics, and there is always a chance of false negatives if an unknown assault doesn't set off any predetermined heuristics.

In the research where the authors classified and detected malwares based on malware behavior has its pros and cons [12]. The field of cyber security faces both clear benefits and difficulties with behavior-based virus identification. Its greatest asset is adaptability, which it uses to efficiently identify both known and unidentified dangers by examining software activity. This makes it possible to defend against emerging assaults, such as those using zero-day vulnerabilities. As long as behavior patterns are constant, polymorphic malware, which alters its code to avoid detection, can be successfully combated. But putting behavior-based systems into practice necessitates sophisticated algorithms, which may make them complex. Real-time behavior

analysis may need a lot of resources, which could affect system performance. Privacy concerns occur owing to data monitoring and the need to balance accuracy and potential false negatives. As a result, behavior-based detection offers dynamic defense but also calls for careful consideration of technical, resource, and privacy considerations.

Another publication [13] summarizes a ground-breaking framework that transforms the cybersecurity landscape. This study addresses a crucial issue: the identification and subsequent detection of malware. It does so by utilizing deep neural networks in conjunction with image-based analysis. The work provides a paradigm change in traditional detection approaches by creatively fusing visual features acquired from malware files, opening up new possibilities for strengthened cybersecurity. The authors provide a brilliant strategy with a careful attention to deep neural networks that is distinguished by its precision and robustness. This work opens the door to rethinking defensive methods necessary to counter the ongoing growth of digital threats in addition to providing a new perspective for malware identification. We set out on an exploration of the interaction between cutting-edge technology, interdisciplinary cooperation, and a steadfast dedication to building a more secure online ecosystem as we dig deeper into the subject. The parts that follow dive into the intricate layers of this ground-breaking framework, illuminating the approaches, obstacles, and realizations that direct cybersecurity's course in a changing environment.

The identification and categorization of malware using deep learning offers a cutting-edge method with a number of significant benefits [14]. These techniques excel at automatically deriving complicated patterns and characteristics from big datasets, enabling them to recognize intricate and altering malware strains. They have exceptional adaptability and can recognize both well-known and novel threats, making them efficient against polymorphic malware and zero-day attacks. Deep learning models have the ability to analyse a variety of data kinds, including text, pictures, and network traffic, enabling a comprehensive analysis of potential threats. Additionally, their accuracy is improved by their capacity to recognize minute correlations in the data. But there are other difficulties with deep learning-based approaches. They are constrained in domains where there is a dearth of readily available training data because to their need for significant amounts of labeled training data. Model interpretability can be a problem since deep neural networks frequently operate as "black boxes," making it challenging

to comprehend why particular classifications are produced. Inadequate management of the data can also lead to overfitting, when models perform well on training data but badly on new data. Scalability concerns may arise because to the complexity of these models, which necessitates heavy computational load for both training and inference.

2.4 Conclusion

A thorough examination of malware detection and classification methods [9] in the field of cybersecurity via the prism of background research and literature analysis uncovers a varied landscape of strategies. A reliable method for detecting known threats, signature-based detection is based on predetermined patterns. By closely examining behavior and characteristics, heuristic-based approaches are adaptable and capable of spotting new and developing malware. A cutting-edge approach called deep learning makes use of detailed patterns discovered through thorough data analysis to provide agility and adaptability in the face of complicated challenges. Each strategy has advantages and disadvantages of its own [9]. While signature-based techniques excel at quickly recognizing known threats, they falter when dealing with malware that is unknown or that is fast changing. Although heuristic-based solutions are adaptable, handling false positives and staying up to date with changing attack vectors can be difficult. Deep learning proves to be a formidable force, demonstrating adaptability, comprehensive analysis, and subtle pattern recognition—even if it takes a lot of data, computer power, and knowledge.

Chapter 3

Methodology

3.1 Introduction

The methodology chapter of a thesis acts as the core of the research project by explaining the methodical strategy taken to answer the research questions and accomplish the stated goals. The steps, methods, and instruments used in data collection, processing, analysis, and interpretation are carefully described in this chapter. The methodology chapter provides readers with a detailed grasp of the rigor and validity of the research by outlining the methods taken to produce results. This chapter creates a link between the intellectual underpinnings of the study and its actual application through a clear and organized explanation. It not only clarifies the direction of the chosen research but also makes it easier to reproduce the results, which promotes the development of new knowledge in the area. As a result, the methodology chapter guides readers through the complexities of the research process, provides insights into how the methodology aligns with the study objectives, and ensures the objectivity of the results.

We will go into great detail on the methodology used in the thesis in this chapter. In the field of machine learning, computers learn from data without explicit programming in order to produce forecasts, judgments, or classifications. In supervised machine learning, which is a subset of machine learning, models are trained on labeled data to find correlations between inputs and desired outputs, enabling them to make precise predictions on brand-new, untainted data. Here, two supervised learning models are introduced, including the multi-layer perceptron neural network method and convolutional neural network.

3.2 Overview of the Proposed System

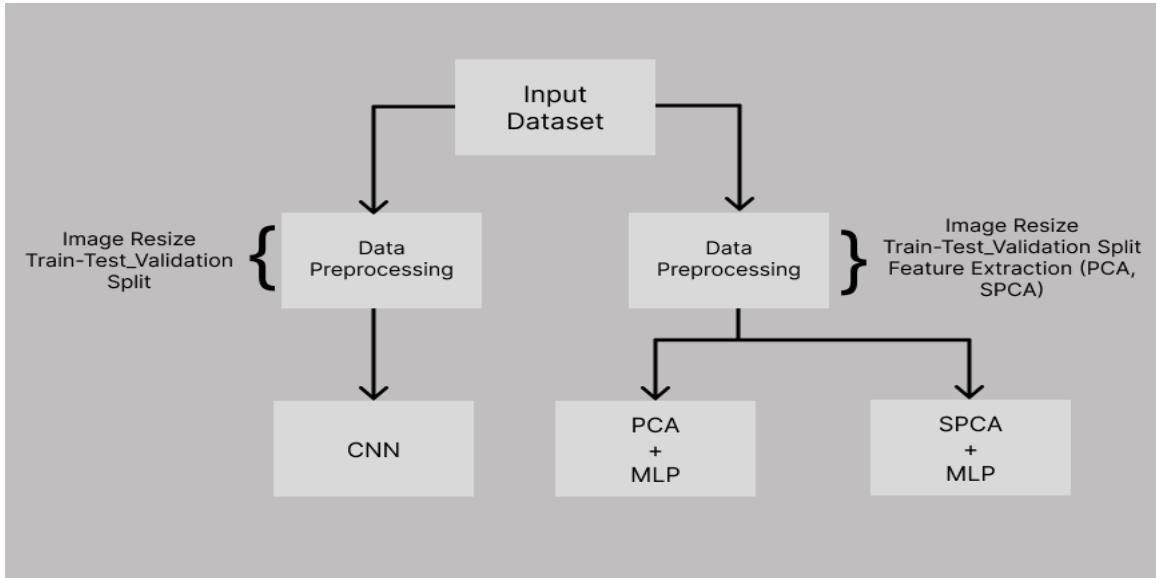


Figure 3.1: Work flow of Proposed System

3.3 Data Collection & Description

The Malimg dataset [15], a key resource in the field of image-based malware classification, is the subject of this section of the study, which focuses on its thorough compilation and characterization. This part examines the painstaking procedure for compiling and preparing the dataset [16], illuminating its special characteristics, organizational scheme, and content.

The methodical procedure for gathering malware pictures from various sources is described in this subsection. It explains how the dataset was carefully selected to include a diverse range of malware families, with the goal of capturing the visual characteristics that were unique to each. The procedures used to guarantee the dataset's diversity and representativeness are described in depth, underlining the need of a well-rounded collection. This section elaborates on how the

dataset naturally represents malware through visuals. The images' dimensions and grayscale makeup shed light on the opportunities and difficulties that come with visual malware investigation. The dataset's new approach to malware classification is underlined by the grayscale

palette, which was intended to concentrate on structural patterns rather than color information.

The properties of the dataset are painstakingly defined in this important section. Transparent information is provided regarding the quantity of samples, how they were distributed among malware families, and any imbalances. The discussion of the dataset’s supporting metadata, including filenames and labels, emphasizes the dataset’s function in facilitating supervised learning experiments. The dataset contains 9339 malware images of 25 different malware families.

3.4 Data Analysis

The 9339 photos were created using malware samples from Microsoft’s Big Dataset [17]. Malware binaries are frequently generated using binary code, which is a pattern of alternating 1s and 0s. This binary information is frequently expressed as integers as the first step. For every integer value, a grayscale pixel value can be assigned. Malware binaries have been used to make malware images in this manner.

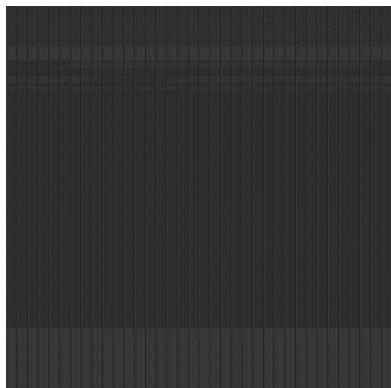


Figure 3.2: Malware Image

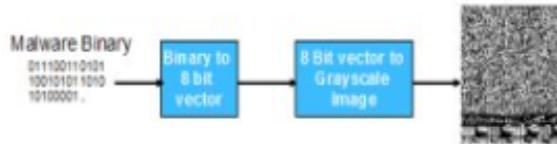


Figure 3.3: Process of Creating Malware Image From Malware Binaries

3.4.1 Data Set Specification

The following discussion includes a full overview of the dataset, how it is used in the study, and attribute information:

(a) Source:

From Dropbox:

Collaborator : Nataraj, Lakshmanan and Karthikeyan, Sreejith and Jacob, Gregoire and Manjunath, Bangalore S.

Author : Nataraj

The data was obtained from :

<https://www.dropbox.com/s/ep8qjakfwh1rzk4/malimgdataset.zip?dl=0>

(b) Relevant Papers:

Nataraj, L., Karthikeyan, S., Jacob, G., Manjunath, B.S.: Malware images: visualization and automatic classification. In: Proceedings of the 8th International Symposium on Visualization for Cyber Security, VizSec '11, pp. 4:1–4:7. ACM, New York, NY, USA (2011). . [15]

<https://dl.acm.org/doi/abs/10.1145/2016904.2016908>

(c) Feature Information: The dataset is predominantly utilized in classification tasks. The dataset is comprised of 9339 records of malware images which represent 25 malware families.

No.	Family	Family Name	No. of Variants
01	Dialer	Adialer.C	122
02	Backdoor	Agent.FYI	116
03	Worm	Allaple.A	2949
04	Worm	Allaple.L	1591
05	Trojan	Alueron.gen!J	198
06	Worm:AutoIT	Autorun.K	106
07	Trojan	C2Lop.P	146
08	Trojan	C2Lop.gen!G	200
09	Dialer	Dialplatform.B	177
10	Trojan Downloader	Dontovo.A	162
11	Rogue	Fakerean	381
12	Dialer	Instantaccess	431
13	PWS	Lolyda-AA 1	213
14	PWS	Lolyda-AA 2	184
15	PWS	Lolyda-AA 3	123
16	PWS	Lolyda.AT	159
17	Trojan	Malex.gen!J	136
18	Trojan Downloader	Obfuscator.AD	142
19	Backdoor	Rbot!gen	158
20	Trojan	Skintrim.N	80
21	Trojan Downloader	Swizzor.gen!E	128
22	Trojan Downloader	Swizzor.gen!I	132
23	Worm	VB.AT	408
24	Trojan Downloader	Wintrim.BX	97
25	Worm	Yuner.A	800

Figure 3.4: 25 Different Families of Malware Images

3.5 Data Preprocessing

The data pre-processing stage of the machine learning life cycle is one of the most crucial stages since it simplifies data analysis, which in turn improves the algorithms' precision and efficiency [18].

The Malimg dataset, also known as the Nataraj Malware Image Dataset, is a collection of grayscale pictures made from malware binaries. Each image in the dataset is a visual representation of a malware binary's integer content that has been converted into pixel values. The dataset offers a varied collection of photos that correlate to various malware families or categories in an effort to aid research on image-based malware classification. Machine learning models can be built and tested for the goal of detecting and classifying malware based on its visual properties thanks to these photos' labeling with their associated malware categories. Researchers and practitioners can create and test image-based classification methods using the Malimg dataset as a helpful resource.

The generation of batches of picture data is necessary for preprocessing. Then, it is necessary to generate batches of data from a directory that has subdirectories for various classes. Splitting the dataset into train and test portions and normalizing the pixel values are also necessary. It is recommended to create one-hot encode label arrays, where each row corresponds to a one-hot encoded vector showing the presence of a specific class.

3.6 Feature Extraction

When it comes to improving model performance, reducing dimensionality, and spotting significant patterns in data, feature extraction is a crucial stage in machine learning and data analysis. Extrapolating pertinent features from raw picture data is crucial for the classification of malware using images. The importance of feature extraction is examined in this part, along with important methods that improve the effectiveness of the classification framework, such as Principal Component Analysis (PCA) and Sparse Principal Component Analysis (SPCA).

Raw photographs are high-dimensional data representations that can be computationally taxing and suffer from the dimensionality curse in the context of classifying malware based on images. In order to help machine learning models better identify underlying patterns, feature extraction attempts to reduce these complicated images into more concise, representative, and useful representations. Faster training, less overfitting, and improved generalization are made possible by this method, which boosts the classification framework's overall effectiveness and accuracy.

In the realm of feature extraction, PCA is a popular linear method. The original features are changed into a fresh set of orthogonal features known as main components. The ability of these elements to capture the greatest amount of data variance determines their order. The dimensionality of the data can be considerably decreased while retaining as much information as feasible by simply keeping a portion of the principal components. This method is very useful for classifying malware based on images since it reduces the impact of noise while revealing important image properties.

The drawback of the conventional approach in capturing sparse and confined features is ad-

dressed by SPCA, a PCA extension. By encouraging sparsity and ensuring that only a small subset of the retrieved features significantly contribute to the representation, SPCA improves feature extraction in the context of image analysis. This is especially important for malware classification based on images since sparse characteristics can reveal discriminatory patterns that separate malware families. The sparseness of SPCA-derived features improves classification accuracy and sheds light on the distinctive characteristics of different malware subtypes.

In conclusion, feature extraction plays a key preparatory role in the identification of malware using images. Uncovering the most informative parts of picture data, reducing dimensionality, and improving the accuracy and effectiveness of classification models are all made possible by techniques like PCA and SPCA. These techniques give the framework the ability to distinguish malware families visually while simultaneously handling the difficulties presented by high-dimensional data.

3.7 Proposed Models

In order to complete this thesis, I used several machine learning and neural network models. The field of artificial intelligence has seen machine learning emerge as a transformative force that has fundamentally altered how data is handled and used. The value of machine learning is best understood when it is applied to classification tasks, where the main goal is to assign data points to previously defined groups or classes. Machine learning methods enable computers to automatically identify patterns from data, streamlining the classification process. This enhances the categorization process' accuracy and effectiveness. Due to the complexity of the patterns they are able to discern and their capacity to adapt to continuously changing situations, these algorithms are able to identify nuanced relationships within the data that human analysts may find challenging to recognize. These techniques are essential for producing precise predictions and well-informed judgments across a wide range of application domains due to their scalability, adaptability, and capacity to generalize patterns to new, unknown data. As the industry develops, machine learning's role in classification will remain crucial, but it will also continue to push the boundaries of what is possible with data-driven automation and intelligence.

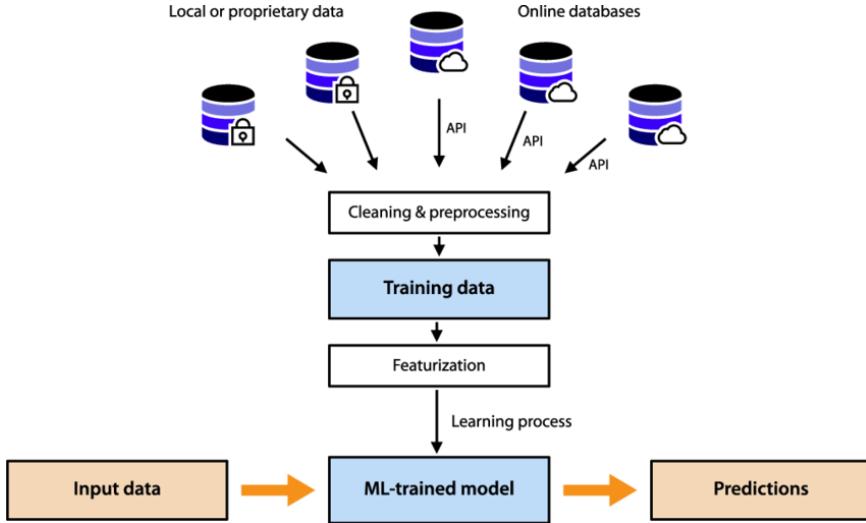


Figure 3.5: Machine Learning model's workflow [19]

3.7.1 Supervised Machine Learning

Intelligent Systems typically do supervised classification. Thus, many AI (Logic-based, Perceptron-based) and statistics (Bayesian Networks, Instance-based) methodologies have been developed. Supervised learning builds a predictive feature-based class label distribution model. The classifier is used to assign class labels to testing instances with known predictor feature values but unknown class labels. The idea of discovering patterns and relationships within labeled data is the cornerstone of the basic technique in the field of artificial intelligence known as supervised machine learning [20]. In a supervised learning situation, a dataset with input features and associated target outputs is given to the algorithm. The algorithm's goal is to discover underlying patterns and relationships in the data in order to learn a mapping from inputs to outputs. The system can now accurately anticipate or categorize previously unknown data points thanks to this learned mapping. Classification, regression, and even some tasks involving natural language processing are among the many activities that fall under the umbrella of supervised learning. The range of its applications includes everything from picture identification and recommendation systems to healthcare and banking. The efficacy of supervised learning has been proved in multiple real-world applications, establishing its significance as a cornerstone in contemporary machine learning approaches. The effectiveness of supervised learning depends on the availability of high-quality labeled data.

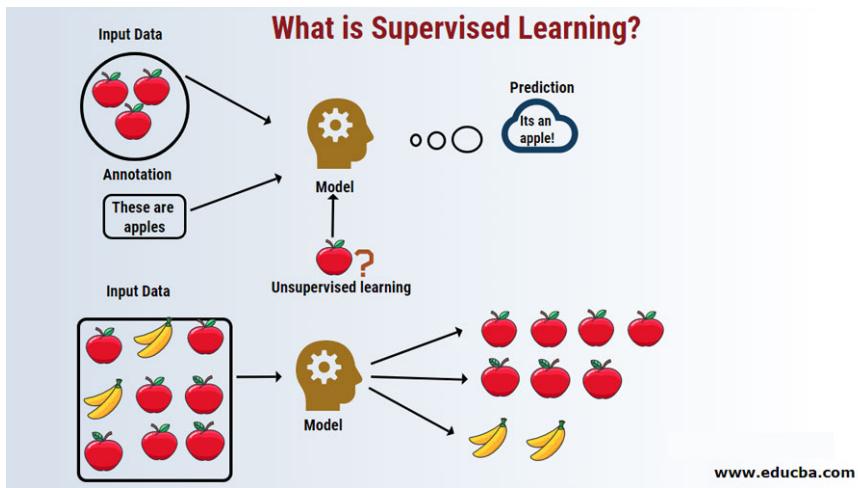


Figure 3.6: Supervised Machine Learning model's workflow [21]

3.7.2 Convolutional Neural Network (CNN)

Convolutional Neural Network (CNN) is a neural network architecture that is commonly used for classification tasks, such as predicting the presence or absence of cardiac disease based on multiple input features. The MLP is composed of layers of interconnected nodes, or neurons, organized as an input layer, one or more concealed layers, and an output layer. In the context of predicting cardiac disease, the input layer represents the classification-contributing features, such as age, cholesterol levels, blood pressure, and more. Each neuron in the input layer is associated with a distinct feature. Using weighted connections and activation functions, the hidden layers, which may contain one or multiple layers, transform the input data. During the training process, each neuronal connection acquires an associated weight. The hidden layers incorporate non-linearity into the model, enabling the MLP to discover complex relationships in the data that may not be apparent to simple linear models such as Logistic Regression. In this scenario, the output layer of the MLP consists of a single neuron representing the binary classification outcome: whether or not the patient has heart disease (1). Usually a sigmoid function, the activation function of the output neuron converts the final weighted sum of inputs into a probability score.

to benefit from the predictive potential of RF and the clarity of Logistic Regression.

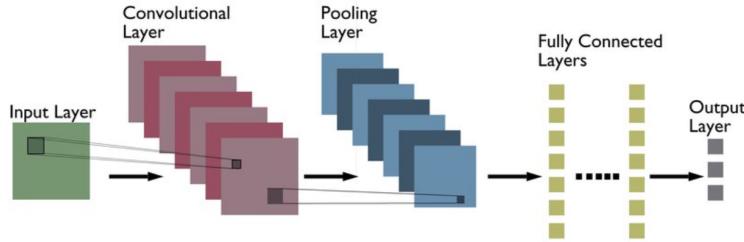


Figure 3.7: CNN Model [22]

3.7.3 Multi-Layer Perceptron (MLP)

Multilayer Perceptron (MLP) is a neural network architecture that is commonly used for classification tasks, such as predicting the presence or absence of cardiac disease based on multiple input features. The MLP is composed of layers of interconnected nodes, or neurons, organized as an input layer, one or more concealed layers, and an output layer. In the context of predicting cardiac disease, the input layer represents the classification-contributing features, such as age, cholesterol levels, blood pressure, and more. Each neuron in the input layer is associated with a distinct feature. Using weighted connections and activation functions, the hidden layers, which may contain one or multiple layers, transform the input data. During the training process, each neuronal connection acquires an associated weight. The hidden layers incorporate non-linearity into the model, enabling the MLP to discover complex relationships in the data that may not be apparent to simple linear models such as Logistic Regression. In this scenario, the output layer of the MLP consists of a single neuron representing the binary classification outcome: whether or not the patient has heart disease (1). Usually a sigmoid function, the activation function of the output neuron converts the final weighted sum of inputs into a probability score. To benefit from the predictive potential of RF and the clarity of Logistic Regression. A fundamental neural network architecture with several applications in numerous industries is the multilayer perceptron (MLP). The MLP exhibits its adaptability in tasks like classification and regression because it consists of interconnected nodes structured into layers, including an input layer, one or more hidden layers, and an output layer. Each node in the input layer represents a feature, and the hidden layers process and change the input data to reveal complex patterns through weighted connections and activation functions. The network modifies weights throughout training to enhance predictions. MLP is an effective tool for tackling complex problems that are beyond the scope of linear models due to its capacity to capture non-linear relationships in data.

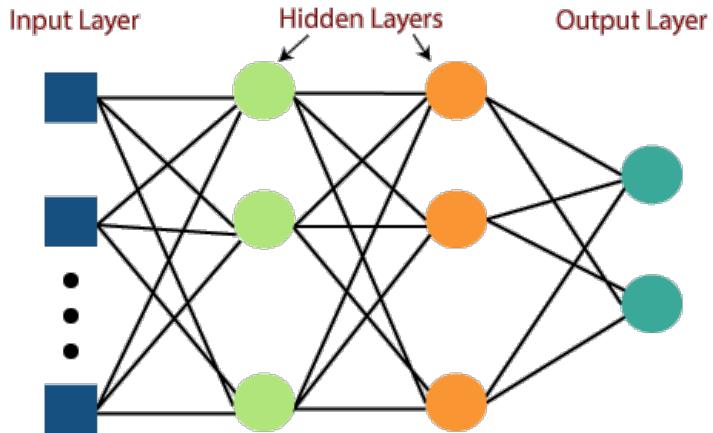


Figure 3.8: Multi-Layer Perceptron Model [23]

3.8 Conclusion

An extensive analysis and overview of the research done for this thesis are provided in this chapter. These discussions have touched on a variety of subjects, including those associated with data collection, analysis, and processing, among other things. Then we gave a brief review of the ML models that had been applied in order to categorize malware families. We looked at each model's individual accuracy as well as how it behaved on the datasets in order to determine which one was the greatest fit for this thesis work.

Chapter 4

Vision Based Malware Classification Using Neural Network

4.1 Introduction

This chapter extends the previous one by giving a detailed theoretical explanation, and it does so by following it with a thorough computational explanation of how the concepts we discussed have been implemented. There are four main steps in the implementation process.

- (a) Detailed Data Information
- (b) Data Preprocessing
- (c) Feature Extraction
- (d) Proposed Models Implementation

4.2 Implementation Tools

Python version 3.7 was the programming language that was utilized, and the Google Colab platform was the environment in which the implementation was carried out. Throughout the entirety of this thesis, a variety of different frameworks, including TensorFlow, NumPy, Pandas, Matplotlib, Seaborn, and Sci-kit Learn, were utilized.

4.3 Detailed Data Analysis

A crucial step on the path from unprocessed data to insightful conclusions is detailed data analysis. It entails systematically exploring, examining, and interpreting data to find trends, anomalies, patterns, and important data that may be concealed within the dataset. Statistical analysis, data visualization, and exploratory techniques are all used in this process, which goes beyond simple observation. Detailed data analysis seeks to expose the data's underlying structure, spot potential correlations between variables, and offer a basis for reasoned decision-making. This analytical method is crucial in converting raw data into usable knowledge, enabling people and organizations to get the most out of their data assets, whether for research, corporate planning, or problem-solving.

4.3.1 Data Set Information

Malimg Dataset [15] was employed. It [16] is made up of 9339 grayscale photos of 25 malware families, with 30% of the total data being utilized for testing and 70% being used for training. A given malware binary can be read as a vector of 8-bit unsigned integers grouped into a 2-dimensional array, and this can be viewed as a grayscale image in the [0, 255] range, with 0 denoting the background and 255 denoting white. Depending on the size of their families, the image varies in size.

A frequently used benchmark dataset for the classification of malware using images is the Malimg dataset. It comprises of a number of grayscale photos that were taken from malware files and each show visual patterns that correspond to various malware families. The dataset contains a wide range of visual traits and complexities with over 9,339 images spread across 25 different virus families. Convolutional neural networks (CNNs) are one type of machine learning model that researchers use to build and test in order to detect and categorize malware based on its visual content. This dataset is essential to the creation of efficient cyber security solutions as well as the advancement of image-based malware analysis. It's crucial to remember that while photos from the same family may appear to resemble one another, they actually include various textures that reveal different aspects of the infection.



Figure 4.1: Folders of 25 Families From Dataset

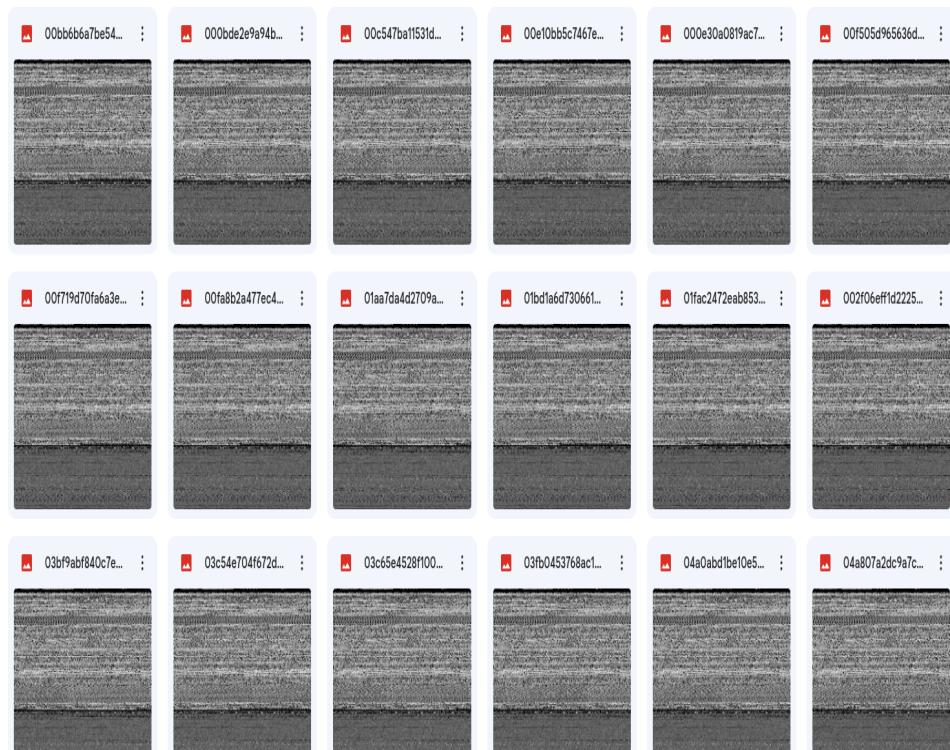


Figure 4.2: Malware Images From Single Folder

4.4 Feature Extraction: Principal Component Analysis (PCA)

Feature selection is an essential practice within the field of machine learning, with the primary objective of improving the quality and efficiency of classification models. The process entails the identification and preservation of the most informative features, while simultaneously eliminating redundant or less significant ones, within a given dataset. Through this process, the chosen features play a significant role in enhancing the model's capacity to identify patterns and generate precise predictions. This procedure not only enhances the prediction performance but also mitigates the risk of overfitting by diminishing the model's dependence on noise or irrelevant features. Moreover, the process of feature selection facilitates the enhancement of interpretability, hence rendering the outcomes of the model more comprehensible and feasible for implementation. Additionally, it optimizes computational resources, resulting in expedited training and prediction durations, a crucial aspect when dealing with extensive datasets. In general, the skillful process of selecting relevant features is a crucial step in attaining classification models that are both dependable and efficient, since they accurately capture the genuine underlying relationships within the data. By restructuring complex data along different axes, referred to as main components, Principal Component Analysis (PCA) functions as a transformative tool that makes complex data simpler. Finding the directions of maximum variance within the data is the basic tenet of PCA. To do this, it first computes the covariance matrix of the original data, and then applies eigenvalue decomposition to separate out the eigenvectors and eigenvalues. The associated eigenvalues describe the quantity of variance along those directions, and the accompanying eigenvectors represent the directions of highest variability. PCA creates a new feature space that largely encompasses the variability in the data by choosing a subset of the most important eigenvectors. This procedure effectively reduces dimensionality while keeping as much variance as feasible by projecting the data onto a lower-dimensional subspace. Data analysis and machine learning are two areas where Principal Component Analysis (PCA), a potent dimensionality reduction method, is applied. PCA tries to retain as much of the variability of the original data as possible while downscaling a high-dimensional dataset into a lower-dimensional space. By locating the principal components—new orthogonal axes that capture the majority of the variance in the data—it is able to accomplish this. The order of these elements reflects the relative importance of their respective eigenvalues. PCA enables for enhanced model performance, noise reduction, and data visualization by choosing a subset of these elements.

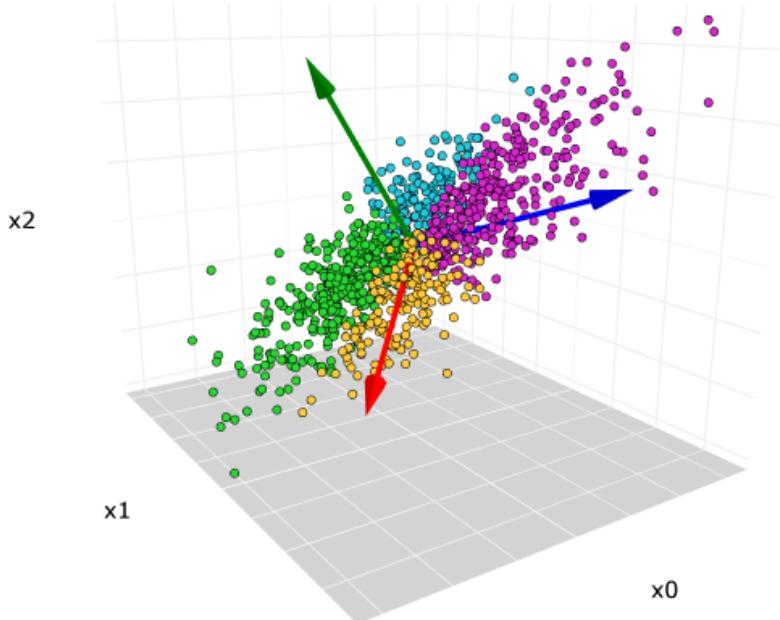


Figure 4.3: PCA Explained Visually [24]

4.5 Feature Extraction: Sparse Principal Component Analysis (SPCA)

A sophisticated variation on the principal component analysis (PCA) method designed to extract localized and sparse characteristics from high-dimensional data is called sparse principal component analysis (SPCA). While PCA looks for orthogonal axes with the greatest variance, SPCA takes a step further by encouraging sparsity in the derived components and highlighting the fact that only a small number of attributes have a major impact on the representation. This distinguishing quality has great utility in a variety of domains, including image analysis, where sparse features effectively capture distinctive traits and detailed patterns.

In order to discover sparse principal component coefficients while reducing reconstruction error, SPCA uses iterative optimization strategies. To find the most distinct and sparse features,

the approach starts by initializing a set of coefficients and then updating them iteratively.

As it supports the selection of sparse features that frequently correlate to significant picture components, SPCA is particularly effective in image-based applications at locating localized patterns, edges, and textures within images. With a focus on sparse representations, SPCA not only improves data reduction but also identifies crucial characteristics that set items or classes apart.

Before supplying data to machine learning algorithms, SPCA can be utilized as a feature extraction technique in a larger context. It has shown promise in increasing model generalization, lowering overfitting, and improving classification accuracy. By placing more emphasis on localized and sparse features, SPCA improves the readability of models and gives domain experts access to information about the precise characteristics influencing predictions.

In conclusion, sparse principal component analysis offers a novel method for feature extraction, especially in cases where localized and sparse properties are crucial. It is a powerful tool for spotting hidden patterns, strengthening model performance, and better comprehending complex information because to its iterative optimization process and emphasis on sparse representations.

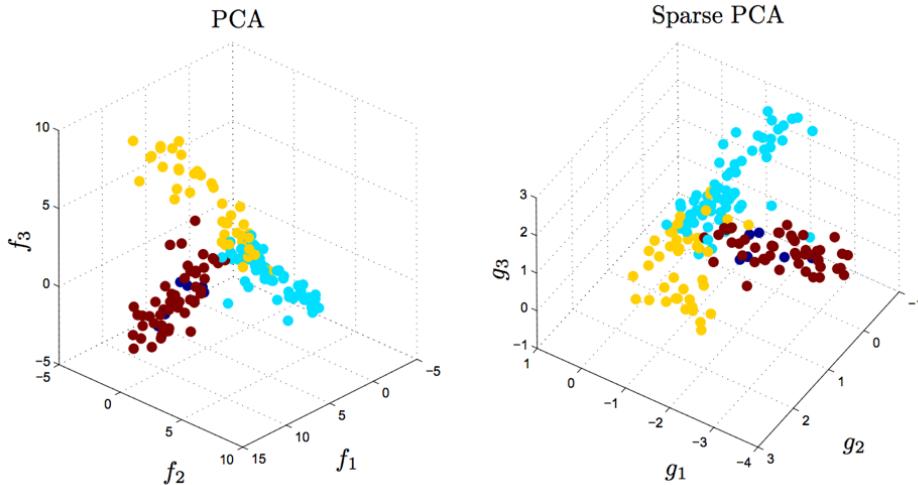


Figure 4.4: PCA-Sparse PCA Explained Visually [25]

4.6 Machine Learning Models Implementation

4.6.1 Convolutional Neural Network (CNN)

The implementation of a Convolutional Neural Network (CNN) model for image classification tasks using the Keras framework is covered in length in this section of the thesis. The provided code covers the model's architecture, configuration, and compilation, setting the groundwork for image classification based on extracted characteristics.

The initialization of a sequential architecture, suggestive of a linear layer-by-layer organization, marks the start of the model structure. A convolutional layer, identified by the Conv2D function, is the first layer added. By using a 3x3 kernel, this layer seeks to extract 128 feature maps. In order to inject non-linearity into the model and enable it to capture complicated relationships within the data, the ReLU activation function, designated by "relu," is used.

A MaxPooling layer is added using the MaxPooling2D function after the convolutional layer. The model's efficiency is increased by this layer's downsampling using a 2x2 pool size, which reduces spatial dimensions and extracts dominating features.

A second convolutional layer is incorporated after the first convolutional-pooling block, using a 3x3 kernel and ReLU activation once more to extract more features from the data. A second MaxPooling layer follows, significantly compressing the data while keeping crucial patterns.

A Flatten layer is added using the Flatten function to change from convolutional layers to fully connected layers. In order to be input into later fully connected layers, this process reshapes the feature maps into a linear array.

There are two completely connected layers added, starting with a Dense layer made up of 128 neurons. The application of the ReLU activation function enables the capturing of complex non-linear relationships in the data. There are 25 neurons in the final output layer, indicated by the Dense function, which is equal to the number of classification categories present in the dataset. Multi-class classification is made possible because to the sigmoid activation function, which helps to create class probabilities.

4.6.2 Multi Layer Perceptron (MLP)

This section provides specifics on how to implement a Multilayer Perceptron (MLP) model for picture categorization. The neural network model is constructed and trained using a module. Two hidden layers, each having 64 and 32 neurons, make up the MLP's defined design. Rectified Linear Unit (ReLU), which is renowned for managing data non-linearities effectively, was selected as the activation function. The algorithm for adaptive learning rate optimization has been selected as the solver, and its name is "adam". A seed number of 33 was chosen at random to assure reproducibility.

The MLP model was implemented in two ways. First with PCA and then with SPCA. The Multilayer Perceptron (MLP) model was implemented using two different strategies. Principal Component Analysis (PCA), a dimensionality reduction method, was used to create the model in the first strategy. By locating the most crucial orthogonal components that accurately captured the variance of the data, PCA made it possible to convert the input characteristics into a lower-dimensional space. The MLP model was then fed these modified components. Another dimensionality reduction technique, Sparse Principal Component Analysis (SPCA), was used to create the model in the second strategy. SPCA produced a more condensed and insightful representation by including sparsity in addition to reducing the dimensionality of the input data.

4.7 Conclusion

This chapter demonstrates how our suggested system is implemented methodically. We've shown you how we gathered and examined the data, as well as the flaws our exploratory data analysis revealed and the fixes we think we can come up with. The generative model, the factors that were utilized to train the model, and some fictitious data samples have also been covered.

Chapter 5

Result and Performance Analysis

5.1 Introduction

The methodical application of our recommended system is analogized in this chapter. We have described the procedures we followed for data collection and analysis, as well as the problems our exploratory data analysis revealed and the fixes we recommend. We have also shown how we gathered and analyzed the data. Furthermore, we went over the generative model, the parameters that were used to train the model, and several samples of synthetic data.

5.2 Data Analysis and Preprocessing

This part delves further into a thorough examination of the various data preparation and analysis techniques used throughout the course of this thesis. We obtain insights into the intrinsic qualities of the dataset by carefully examining several types of data analysis methodologies, revealing crucial patterns and trends that serve as the foundation for our following decisions. Additionally, a thorough analysis of the data preparation techniques used enables us to polish and shape the dataset, ensuring it is optimum for successful machine learning model training and evaluation. This section acts as the cornerstone for our efforts to use cutting-edge approaches for vision-based malware detection and classification.

5.2.1 Feature Extraction

We applied Principal Component Analysis(PCA) to improve the effectiveness and interpretability of our feature space. We wanted to keep the 30 primary components from the original dataset that were the most important, so we fixed the number of components to 30. PCA maintained a sizable amount of the variance in the high-dimensional data while orthogonally transforming it into a lower-dimensional space. This made it easier to express the data succinctly while maximizing processing resources without compromising important data.

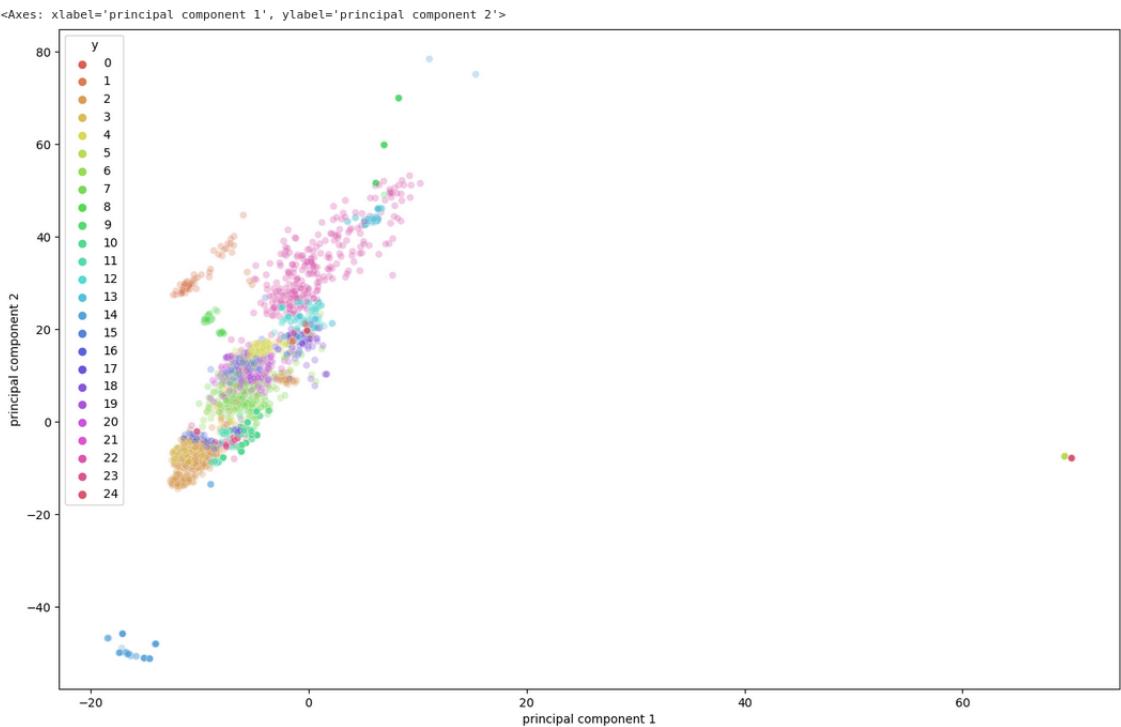


Figure 5.1: Visualization of Data Using First Two Principal Components

We used Sparse Principal Component Analysis (SPCA) to try to improve the feature representation. We instructed the algorithm to remove 30 significant sparse components from the input data using the parameter. We successfully balanced the sparsity-promoting and L2 regularization terms in the optimization process with an alpha value of 0.1 and a ridge_alpha of 0.01. The SPCA model included the distinctive features of our flattened training dataset using an iterative methodology with a maximum of 1000 iterations.

5.3 Performance Evaluation Metrics

Metrics for evaluating performance are indispensable for determining the efficacy of models, algorithms, or systems in numerous disciplines, such as machine learning, data analysis, and decision-making. These metrics provide quantitative measures for evaluating the quality, precision, and efficacy of a model's predictions or classifications.

Confusion Matrix:

The utilization of a confusion matrix is prevalent in classification tasks as a means to visually represent the efficacy of a machine learning model. It accomplishes this by presenting the quantities of True Positive, True Negative, False Positive, and False Negative predictions. This analysis offers valuable insights regarding the precision and types of errors made by the model when classifying various categories. The following is an analysis of the constituent elements of a confusion matrix:

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN

Figure 5.2: Confusion Matrix [26]

True Positive (TP): The number of instances that were correctly predicted as positive (as members of the positive class).

True Negative (TN): The number of instances that were correctly predicted as negative (belonging to the negative class).

False Positive (FP): Incorrectly predicted as positive when they truly belong to the negative class (Type I error).

False Negative (FN): The number of instances incorrectly predicted as negative when they are truly positive (Type II error).

Recall(Sensitivity or True Positive Rate):

Recall is the number of correct predictions of true positives compared to the total number of true positives in the collection. A high recall means that the model does a good job of finding true positives, but it could also mean that there are more false positives.

$$Recall = \frac{TruePositive}{TruePositive + FalseNegative} \quad (5.1)$$

Precision:

Precision is the number of correct positive predictions out of all the positive predictions made by the model. A high precision means that the model is likely to be correct if it predicts a positive instance. But this could mean losing out on some positive instances.

$$Precision = \frac{TruePositive}{TruePositive + FalsePositive} \quad (5.2)$$

F1 Score:

When classes are distributed unevenly, with one class having much more instances than the other, the F1 score is especially helpful. In these situations, accuracy alone may not be a valid indicator of a model's effectiveness because the model may be able to predict the majority class with high accuracy.

$$F1Score = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (5.3)$$

The key performance indicators for evaluating machine learning are recall, precision, and F1 score. Recall, commonly referred to as the true positive rate, gauges how well the model is able to recognize all pertinent events. The accuracy of the model's positive predictions is indicated by precision, which is the percentage of accurately predicted positive cases among all anticipated positives. The F1 score combines recall and accuracy, offering a fair evaluation of a model's performance by taking into account both false positives and false negatives. When taken as a whole, these metrics provide insightful information about how well a model performs in various categorization task contexts.

5.4 Results of the classification models:

The effectiveness of models, algorithms, or systems must be evaluated using performance metrics in a variety of fields, including machine learning, data analysis, and decision-making. These metrics offer numerical measurements for assessing the accuracy, quality, and effectiveness of a model's predictions or classifications.

We used three methods on our malimg dataset [16]. At first, we used CNN. In terms of second method, we first used PCA(30 components) to extract important features which have most impact on variance and then used MLP as classifier. For the third method, we extracted features using SPCA(30 components) and again used MLP for classification. The results with accuracy, recall, precision and f1 score are below :

Table 5.1: Performance Analysis of my work

Models	Accuracy	Recall	Precision	F1 Score
CNN	97.57%	95.68%	95.68%	95.68%
PCA+MLP	97.13%	94.69%	94.23%	94.00%
SPCA+MLP	97.49%	95.44%	94.82%	94.51%

5.4.1 Convolutional Neural Network (CNN)

The CNN model processes visual data using a multi-layered architecture. It begins with two convolutional layers that each have 128 feature mappings and a 3x3 kernel size. The Rectified Linear Unit (ReLU) activation function is used to add nonlinearity to the system. The feature maps are essentially shrunk by the subsequent max pooling layers, which have a 2x2 pooling size. The 2D feature maps are then converted into a 1D vector by a flatten layer and fed into a layer that is completely linked and made up of 128 neurons that have been activated by ReLU. Predictions for the multi-label classification task are produced by the final output layer, which has 25 neurons and a sigmoid activation function. The Categorical Cross-Entropy loss function is used to build the model, which is then optimized using the Adam optimizer and assessed using the accuracy measure.

5.4.1.1 Figure of confusion matrix :

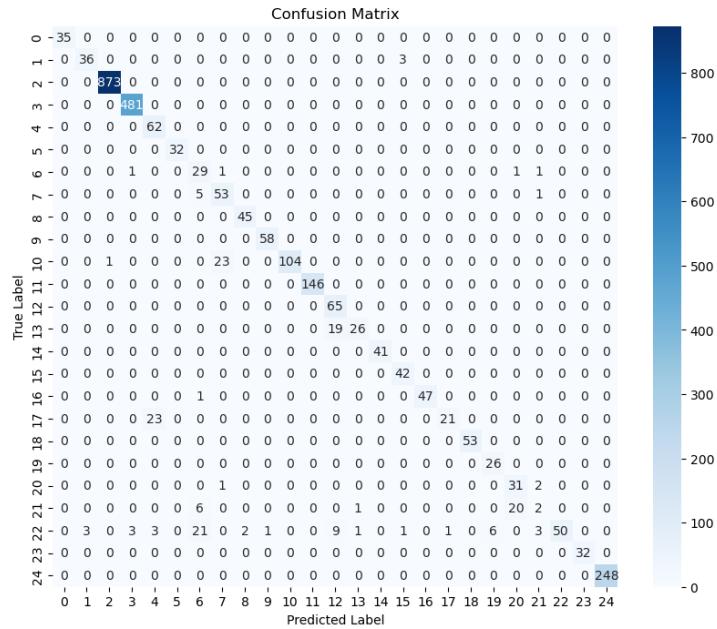


Figure 5.3: Confusion matrix of CNN Model

5.4.1.2 Precision of each class :

precision	
0	1.00
1	0.92
2	1.00
3	0.99
4	0.70
5	1.00
6	0.47
7	0.68
8	0.96
9	0.98
10	1.00
11	1.00
12	0.70
13	0.93
14	1.00
15	0.91
16	1.00
17	0.95
18	1.00
19	0.81
20	0.60
21	0.22
22	1.00
23	1.00
24	1.00

Figure 5.4: Precision Of Each Class

5.4.2 Principal Component Analysis (PCA) & Multi Layer Perceptron (MLP)

Principal Component Analysis (PCA) helps to capture important information while reducing noise by reducing the dimensionality of data. The Multi-Layer Perceptron (MLP), a neural network architecture renowned for its capacity to represent complex relationships, is then given the reduced information. This combination takes advantage of PCA's data reduction and MLP's strong learning capabilities, offering a potent approach for effective classification tasks like malware detection. This section details the mechanics of how we used PCA as a feature extractor to create a Multilayer Perceptron (MLP) model for image categorization. We used 30 principal components as they had the most impact on variance. Using a module, the neural network model is built and trained. The MLP is defined by two hidden layers, each with 64 and 32 neurons. As the activation function, the Rectified Linear Unit (ReLU), which is recognized for successfully controlling data non-linearities, was used. The solver chosen is called "adam" and it is a method for adaptive learning rate optimization. 33 was selected at random as the seed number to ensure reproducibility.

5.4.2.1 Figure of confusion matrix :

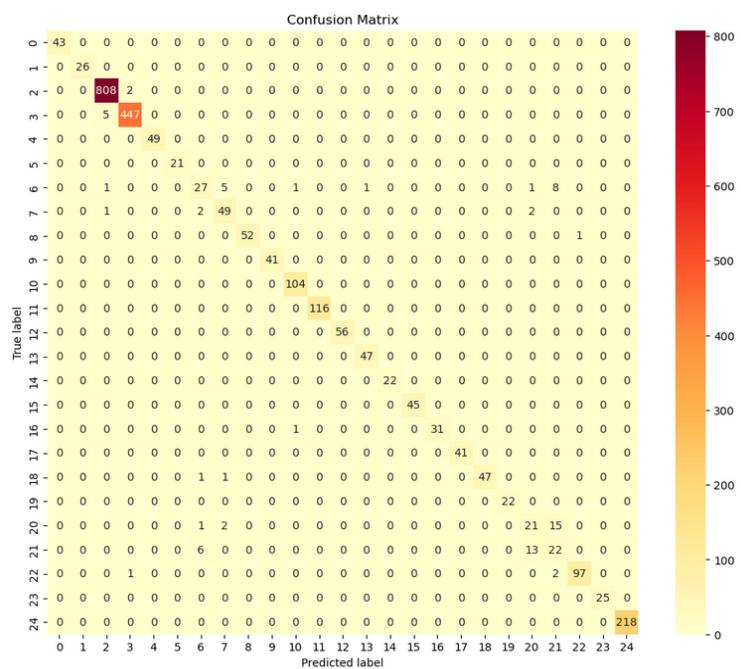


Figure 5.5: Confusion matrix of PCA+MLP Model

5.4.2.2 Figure of ROC curve:

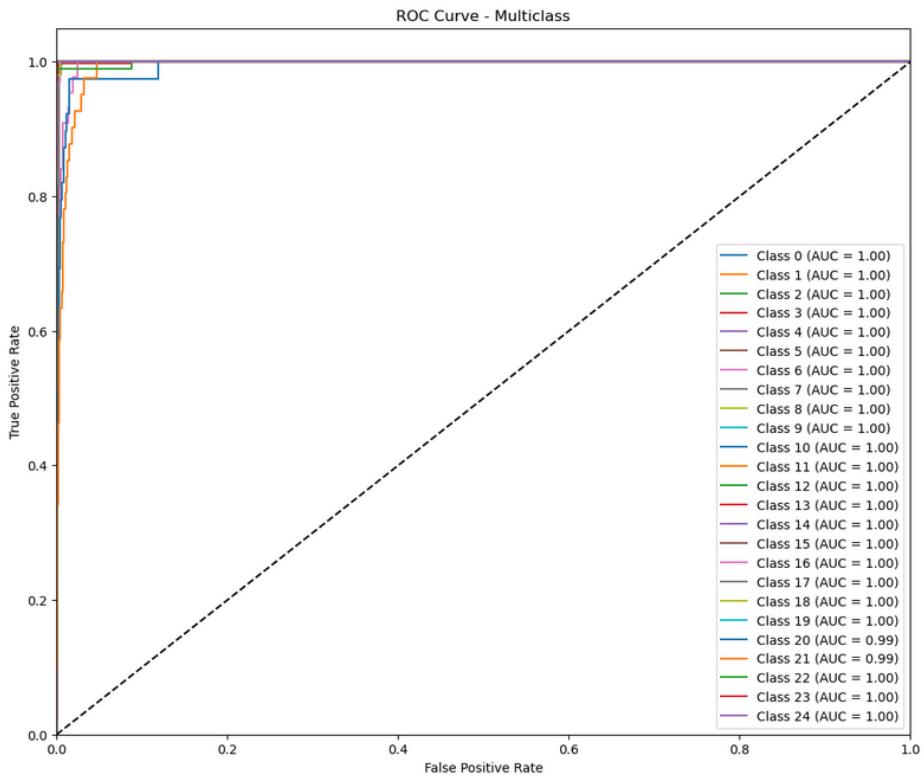


Figure 5.6: ROC Curve of PCA+MLP Model

5.4.2.3 Precision of each class :

```
Precision for class 0: 1.0
Precision for class 1: 1.0
Precision for class 2: 0.9914110429447853
Precision for class 3: 0.9933333333333333
Precision for class 4: 1.0
Precision for class 5: 1.0
Precision for class 6: 0.7297297297297297
Precision for class 7: 0.8596491228070176
Precision for class 8: 1.0
Precision for class 9: 1.0
Precision for class 10: 0.9811320754716981
Precision for class 11: 1.0
Precision for class 12: 1.0
Precision for class 13: 0.9791666666666666
Precision for class 14: 1.0
Precision for class 15: 1.0
Precision for class 16: 1.0
Precision for class 17: 1.0
Precision for class 18: 1.0
Precision for class 19: 1.0
Precision for class 20: 0.5675675675675675
Precision for class 21: 0.46808510638297873
Precision for class 22: 0.9897959183673469
Precision for class 23: 1.0
Precision for class 24: 1.0
```

Figure 5.7: Precision Of Each Class

5.4.3 Sparse Principal Component Analysis (SPCA) & Multi Layer Perceptron (MLP)

A powerful method for feature extraction and classification is introduced by combining Sparse Principal Component Analysis (SPCA) and Multi-Layer Perceptrons (MLP). Information condensation is facilitated by SPCA, which recovers sparse and meaningful representations from high-dimensional data. A strong method for jobs like malware classification, where both effective data reduction and complicated connection modeling are essential, is produced by the neural network's subsequent integration with MLP, which taps into its ability to find intricate patterns. This section describes the steps we took to build a Multilayer Perceptron (MLP) model for image categorization using SPCA as a feature extractor. 30 major components were picked since they had the greatest influence on variance. The neural network model is created and trained using a module. Two hidden layers, each having 64 and 32 neurons, make up the MLP. The Rectified Linear Unit (ReLU), which is renowned for effectively controlling data non-linearities, was employed as the activation function. The "adam" solver, which employs a technique for adaptive learning rate optimization, was selected. To reproduce, a seed number of 33 was chosen at random.

5.4.3.1 Figure of confusion matrix :

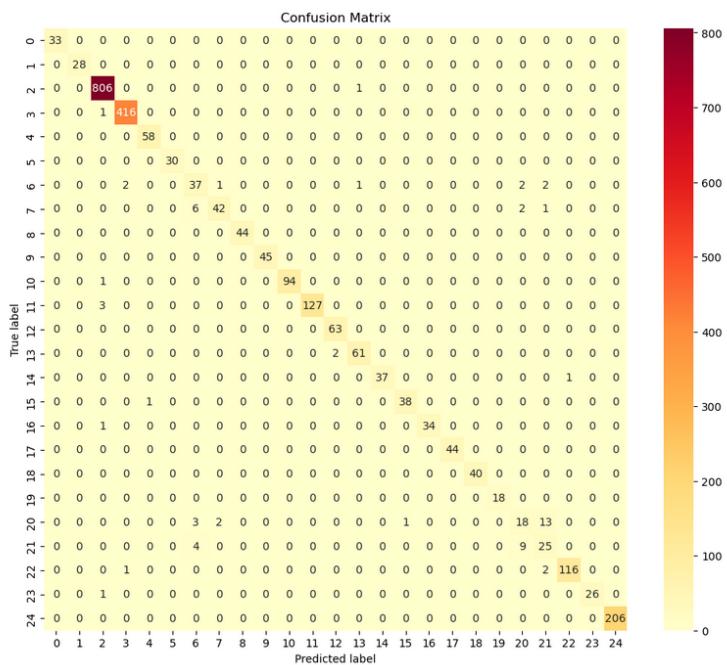


Figure 5.8: Confusion matrix of SPCA+MLP Model

5.4.3.2 Figure of ROC curve:

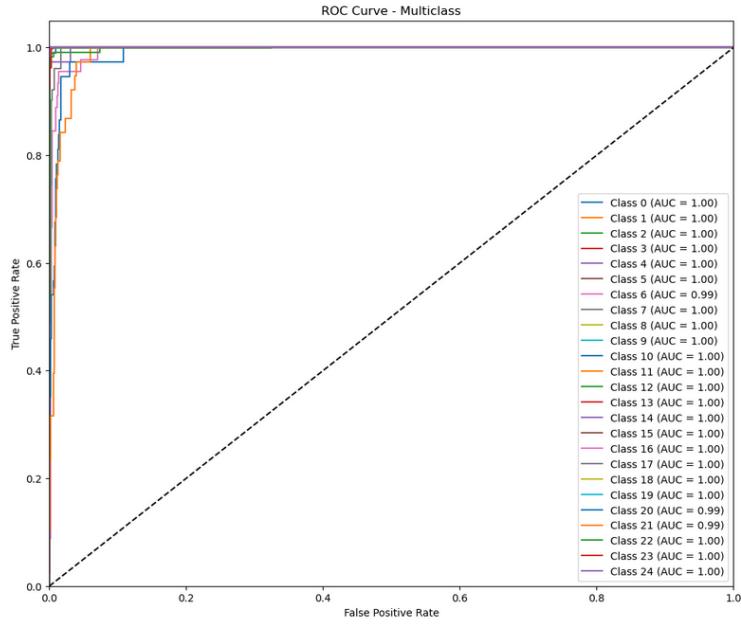


Figure 5.9: ROC Curve of SPCA+MLP Model

5.4.3.3 Precision of each class :

```
Precision for class 0: 1.0
Precision for class 1: 1.0
Precision for class 2: 0.991389913899139
Precision for class 3: 0.9928400954653938
Precision for class 4: 0.9830508474576272
Precision for class 5: 1.0
Precision for class 6: 0.74
Precision for class 7: 0.9333333333333333
Precision for class 8: 1.0
Precision for class 9: 1.0
Precision for class 10: 1.0
Precision for class 11: 1.0
Precision for class 12: 0.9692307692307692
Precision for class 13: 0.9682539682539683
Precision for class 14: 1.0
Precision for class 15: 0.9743589743589743
Precision for class 16: 1.0
Precision for class 17: 1.0
Precision for class 18: 1.0
Precision for class 19: 1.0
Precision for class 20: 0.5806451612903226
Precision for class 21: 0.5813953488372093
Precision for class 22: 0.9914529914529915
Precision for class 23: 1.0
Precision for class 24: 1.0
```

Figure 5.10: Precision Of Each Class

5.5 Conclusion

This chapter starts off by demonstrating data preparation and feature selection. In order to identify which model yields the most accurate results, the final analysis compares the output of the models to the output of the base article and displays the results visually.

Chapter 6

Conclusion and Future Works

6.1 Introduction

The results of our experiments and an assessment of how well they worked were discussed in the preceding chapter. This chapter includes an overview of all of our research activities, followed by a review of the study's limitations. The discussion of potential future works comes next, and then the chapter is over. The results of our experiments and an assessment of how well they worked were discussed in the preceding chapter. This chapter includes an overview of all of our research activities, followed by a review of the study's limitations. The discussion of potential future works comes next, and then the chapter is over.

6.2 Thesis Summary

The importance of a dataset's quality and composition in terms of its performance in the context of machine learning models cannot be understated when it comes to making judgments based on data. Before the predictive power of these models can be unleashed, it is crucial that the dataset be put through a number of preprocessing steps that increase its resilience and predictive capacity. Here are some examples of these preprocessing techniques. The current study examines a full preparation process in this setting. By not deleting any, the imputation makes it possible to use every sample in the dataset.

The classification of malware photos inside the Malimg dataset was thoroughly investigated

in the present thesis using a variety of machine learning approaches. Convolutional neural networks (CNNs), principal component analysis (PCA) [24] in conjunction with multi-layer perceptrons (MLPs), and sparse PCA (SPCA) in conjunction with MLPs are three different approaches that were investigated.

The initial method used CNNs, a deep learning method famed for image classification jobs, to harness their power. Multiple convolutional and pooling layers, followed by tightly coupled layers, were used to create a unique architecture. The CNN demonstrated its skill in capturing complex visual elements essential to malware classification with considerable accuracy.

Then, PCA[25], a dimensionality reduction method, was used with MLPs. The PCA-MLP hybrid provided a simplified representation that kept important variance by downscaling the original feature space into a lower-dimensional subspace. The model demonstrated admirable outcomes, demonstrating the effectiveness of feature extraction for challenging tasks like malware classification.

The study also explored the world of sparse representation using SPCA in conjunction with MLPs. The SPCA-MLP model displayed improved interpretability and efficiency by introducing sparsity restrictions throughout the learning phase. This was very helpful while attempting to classify virus images because it is crucial to recognize significant traits. By integrating Sparse Principal Component Analysis (SPCA) and Multi-Layer Perceptrons (MLPs), the research explored the world of sparse representation in an effort to improve the classification process. A compelling framework was introduced by this ground-breaking SPCA-MLP model, which is distinguished by its combined advantages of improved interpretability and computational effectiveness. The model made sure that only important and instructive elements were maintained by incorporating sparsity constraints during the learning process, which resulted in a more condensed representation of the data. This quality proved to be especially helpful when it came to classifying viral photos, where it is crucial to detect key characteristics. In the context of malware detection, the SPCA-MLP synergy enhanced the knowledge of the underlying patterns and reduced the computing burden, making it possible to use a classification strategy that is more precise and effective.

The potential of CNNs, PCA-MLP, and SPCA-MLP algorithms for the classification of malware images inside the Malimg dataset was thoroughly investigated in this thesis. The results underlined the importance of various techniques by illuminating their relevance and efficiency in solving actual cybersecurity concerns. The results not only advance the field of malware identification but also offer insightful information about how dimensionality reduction, neural networks, and image analysis interact.

6.3 Limitations

While the thesis progresses, some limitations are encountered. These are discussed below :

- (a) The usage of a particular dataset, the Malimg dataset, is the main constraint of this thesis. Although the dataset provides insightful information regarding the classification of malware images, the lack of diversity may limit the applicability of the suggested models to actual situations. Extending the testing to larger and more varied datasets would offer a more thorough analysis of the algorithms' effectiveness against different kinds of malware.
- (b) When dimensionality reduction methods like PCA and SPCA are used, a trade-off must be made between lowering complexity and the potential loss of important data. The number of components in this thesis was somewhat arbitrarily chosen. The ideal number of components for each technique and dataset could be determined using a more methodical approach, such as cross-validation.
- (c) While CNNs have shown excellent accuracy, deciphering their inner workings can be difficult. As a result, it may be challenging to comprehend how and why particular judgments are made by the models. Deeper understanding of feature extraction and classification decisions may be obtained by including methods for model interpretability and visualization.

6.4 Future Works

- (a) Will try to improve the outcome of classification.
- (b) Will use some different dataset to apply the model and observe the outcomes.
- (c) Will try to implement ensembled models to search for a better result on the dataset used in this work.
- (d) The real-time detection of malware is critical for effective cybersecurity. Future work could focus on optimizing the proposed models for deployment in resource-constrained environments, such as IoT devices or network appliances.

6.5 Conclusion

Overall, the work in this thesis was done carefully and dealt with challenges to reach a conclusion. The outcomes demonstrated that customized preprocessing techniques influence the performance of machine learning models. It brought to light the nuanced connection between feature choice and model performance. In order to optimize dataset value and predictive analytics potential, this study promotes a comprehensive strategy that integrates preprocessing, feature engineering, and model evaluation.

REFERENCES

- [1] “Malware.” <https://www.techtarget.com/searchsecurity/definition/malware/>. Accessed: [Date].
- [2] Ö. Aslan and A. A. Yilmaz, “A new malware classification framework based on deep learning algorithms,” *Ieee Access*, vol. 9, pp. 87936–87951, 2021.
- [3] A. not specified, “Analysis of computer vision techniques in malware classification,” Publication year not specified. Accessed: [Date].
- [4] A. not specified, “Malware statistics: Facts and trends [blog post],” Publication year not specified. Accessed: [Date].
- [5] T. Jiang, J. L. Gradus, and A. J. Rosellini, “Supervised machine learning: a brief primer,” *Behavior Therapy*, vol. 51, no. 5, pp. 675–687, 2020.
- [6] A. not specified, “An introduction to machine learning: Its importance, types, and applications [blog post],” Publication year not specified. Accessed: [Date].
- [7] P. Cunningham, M. Cord, and S. J. Delany, “Supervised learning,” pp. 21–49, 2008.
- [8] A. not specified, “How does a neural network work intuitively in code [blog post],” Publication year not specified. Accessed: [Date].
- [9] Ö. A. Aslan and R. Samet, “A comprehensive review on malware detection approaches,” *IEEE access*, vol. 8, pp. 6249–6271, 2020.
- [10] M. F. Zolkipli and A. Jantan, “A framework for malware detection using combination technique and signature generation,” pp. 196–199, 2010.
- [11] P. Khodamoradi, M. Fazlali, F. Mardukhi, and M. Nosrati, “Heuristic metamorphic malware detection based on statistics of assembly instructions using classification algorithms,” pp. 1–6, 2015.

- [12] W. Liu, P. Ren, K. Liu, and H.-x. Duan, “Behavior-based malware analysis and detection,” pp. 39–42, 2011.
- [13] Y. Jian, H. Kuang, C. Ren, Z. Ma, and H. Wang, “A novel framework for image-based malware detection with a deep neural network,” *Computers & Security*, vol. 109, p. 102400, 2021.
- [14] R. Chaganti, V. Ravi, and T. D. Pham, “Deep learning based cross architecture internet of things malware detection and classification,” *Computers & Security*, vol. 120, p. 102779, 2022.
- [15] L. Nataraj, S. Karthikeyan, G. Jacob, and B. S. Manjunath, “Malware images: visualization and automatic classification,” pp. 1–7, 2011.
- [16] A. not specified, “Malimg dataset,” Publication year not specified. Accessed: [Date].
- [17] A. not specified, “Kaggle competition: Malware classification,” Publication year not specified. Accessed: [Date].
- [18] J. Huang, Y.-F. Li, and M. Xie, “An empirical analysis of data preprocessing for machine learning-based software cost estimation,” *Information and software Technology*, vol. 67, pp. 108–127, 2015.
- [19] A. not specified, “Simplified overview of a machine learning workflow,” Publication year not specified. Accessed: [Date].
- [20] V. Nasteski, “An overview of the supervised machine learning methods,” *Horizons. b*, vol. 4, pp. 51–62, 2017.
- [21] A. not specified, “What is supervised learning?,” Publication year not specified. Accessed: [Date].
- [22] A. not specified, “Different types of cnn architectures explained with examples,” Publication year not specified. Accessed: [Date].
- [23] A. not specified, “Multi-layer perceptron in tensorflow,” Publication year not specified. Accessed: [Date].

- [24] A. not specified, “Principal component analysis (pca) explained visually with zero math,” Publication year not specified. Accessed: [Date].
- [25] A. H. Williams, “Understanding principal component analysis (pca),” March 2016. Accessed: [Date].
- [26] T. D. Science, “Confusion matrix for your multi-class machine learning model,” Date. Accessed: [Date].