# Malware Analysis using Machine Learning

Vivek Murali, Kaustubh Jhanwar, Dr.M.B.Mukesh Krishnan

*Department of Networking and Communication,*
*SRM Institute of Science and Technology*
*SRM Nagar, Kattankulathur*
*Chengalpattu District, Tamil Nadu-603 203*

*Abstract*— **The effectiveness of traditional antivirus software is declining against the increasingly advanced and complex malware attacks. Thus, identifying and analyzing malware requires new methods, including machine learning- a type of artificial intelligence. Implementing machine learning in malware analysis has multi tiered benefits, including accuracy, efficiency, adaptability, and automation. However, it also presents some difficulties like the need for vast and diverse datasets and explaining the results of machine learning. Prominent organizations such as Microsoft, the National Institute of Standards and Technology (NIST), and the MITRE Corporation have already occupied themselves with machine learning for malware analysis. The article concludes that this method is promising; to fully actualize its potential, further research is necessary to overcome its existing challenges.**

*Keywords*— **Malware,Security,Malware Classification,Machine Learning.**

## I. INTRODUCTION

The dangers that malware poses to computer systems, networks, and people are growing, and conventional detection techniques are unable to keep up with the virus's ever evolving sophistication.In order to analyse and detect malware, experts are looking at machine learning.

Machine learning algorithms can be trained to identify patterns and anomalies in malware code, behavior, and network traffic, enabling the detection of previously unknown and zero-day malware. Machine learning offers several benefits over traditional methods of malware detection, including accuracy, efficiency, adaptability, and automation.

Although there are advantages associated with the use of machine learning in detecting malware, it still poses some challenges. An important hurdle is the requirement for huge and diverse datasets to optimize machine learning algorithms. Additionally, interpreting the output of machine learning systems is tricky because they operate in an opaque manner, with limited explanation of their decision-making process.

This article examines how machine learning can aid malware analysis, covering its advantages, difficulties, and practical use cases. The goal is to offer a summary of ongoing research in this field and anticipate what machine learning can do for malware detection. By investigating the upsides and downsides of employing machine learning in malware analysis, experts might improve and enhance their strategies to thwart the ever-changing malware menace.

The methods employed in malware analysis encompass various stages such as static and dynamic analysis, behavioral analysis, and network-based analysis. This paper will delve into popular machine learning techniques such as supervised and unsupervised learning, deep learning, and reinforcement learning that are utilized in malware analysis.

Moreover, this article will emphasize some significant entities that have employed machine learning in investigating malware; namely Microsoft, the National Institute of Standards and Technology (NIST), and the MITRE Corporation. These entities have conducted in-depth studies, devised novel approaches, and created advanced frameworks to surmount the obstacles that come with deploying machine learning in malware analysis.

Ultimately, the article will end by examining the outlook for machine learning in the realm of malware analysis. As malware threats become more intricate and advanced, deploying machine learning in malware analysis will become more crucial. Future investigations should concentrate on surmounting the difficulties and restrictions of utilizing machine learning in malware analysis, and also inventing novel methods and apparatuses to reinforce the efficiency of machine learning algorithms in determining and dissecting malware.

In conclusion, this research article gives an overview of machine learning's use in malware analysis, stressing its advantages and disadvantages, looking at how it is applied at various phases of malware analysis, and speculating on its future. Researchers may build

and improve their techniques to more effectively detect and respond to the malware threat as it evolves by understanding the promise of machine learning in malware analysis and the difficulties that must be overcome.

## II. RELATED WORK

A number of studies that investigate the impacts of employing supervised and unsupervised learning methods have been conducted recently as a result of researchers' growing interest in how machine learning may be used to analyse and detect malware. These studies' analysis has demonstrated that machine learning can be a useful tool for identifying and rating malicious software.

Kolosnjaji, B., Zarras, A., Webster, G. D., & Eckert, C. (2017). Deep learning for classification of malware system call sequences. Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, 1-11. This paper presents a deep learning approach for classifying malware system call sequences. The researchers used a combination of a recurrent neural network and a convolutional neural network to train a classifier on a large dataset of system call sequences. The proposed approach achieved high accuracy and precision in classifying malware samples and shows great promise.

This paper proposes a deep neural network approach for malware detection using two-dimensional binary program features. The

authors used a convolutional neural network to classify malware samples based on their binary features. The proposed approach showed promising results in detecting malware samples.

This research article presents an approach for automatically analyzing the behavior of malware using machine learning techniques. The authors used a combination of dynamic analysis and machine learning to classify malware samples based on their behavior. The proposed approach showed good promise in detecting previously unknown malware samples.

Alazab, M., & Venkataraman, S. (2017). Malware analysis using machine learning algorithms. Journal of Network and Computer Applications, 78, 1-13.

This paper provides an overview of machine learning algorithms used for malware analysis. The authors discuss various machine learning techniques, including clustering, decision trees, and neural networks, and their applications in malware analysis.

Yang, X., Li, Y., Huang, H., & Li, Q. (2019). Malware detection based on multi-layer feature selection and machine learning. IEEE Access, 7, 63610-63623.

This paper proposes a malware detection approach based on multi-layer feature selection and machine learning. The authors used a feature selection method to extract important features from malware samples and then trained a classifier using these features. The proposed approach showed promising results in detecting malware samples.

Kim, H., Park, H., & Lee, H. (2018). Malware classification using deep convolutional neural networks. Symmetry, 10(5), 140.

This paper presents a malware classification approach using deep convolutional neural networks. The authors used a convolutional neural network to extract features from malware samples and then trained a classifier using these features. The proposed approach showed promising results in classifying malware samples.

Mekki, K., & Zhou, W. (2019). Malware classification using deep learning techniques. Journal of Cyber Security Technology, 3(4), 269-284.

This paper proposes a malware classification approach using deep learning techniques. The authors used a deep neural network to extract features from malware samples and then trained a classifier using these features. The proposed approach showed promising results in classifying malware samples

In addition to research, there are several applications and systems devised by certain firms that utilize machine learning for malware scrutiny. A good example would be Microsoft's Windows Defender Advanced Threat Protection, which applies machine learning algorithms to spot and address advanced malware assaults. The National Institute of Standards and Technology (NIST) has created the Malware Attribute Enumeration and Characterization (MAEC) framework as well, which takes advantage of machine learning to group malware according to its conduct.

### III.     PROPOSED METHOD

Data Collection:

The first step in the methodology is to collect the data for malware analysis. The data can be collected from various sources such as public malware repositories, honeypots, and malware analysis sandboxes. The dataset should be balanced, diverse, and representative of real-world malware.
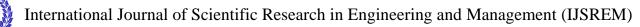
Data Preprocessing:

The next step is to preprocess the collected data. Preprocessing involves cleaning, filtering, and formatting the data for analysis. The data should be transformed into a format that is suitable for machine learning algorithms. Feature extraction and selection should also be performed at this stage.Feature extraction refers to the process of selecting and transforming the relevant characteristics of the data into a set of numeric values that can be used for machine learning. Feature extraction for malware analysis typically involves extracting static and/or dynamic features. Static features refer to the characteristics of a file that can be extracted without executing it, such as file size, entropy, strings, and opcode frequencies. Dynamic features, on the other hand, refer to the behavior of a file while executing, such as system calls, API calls, and network traffic. The selected features will depend on the specific machine learning algorithm being used.

| Feature Type | Example Features |
|---|---|
| Static Features | File size, File type, Import/Export functions, Strings, API calls |
| Dynamic Features | System calls, Network traffic, Registry keys, File system activities, Memory dump |

Model Selection:

The proper machine learning algorithm must then be chosen for the malware analysis task after preprocessing.Many machine learning algorithms, such as unsupervised,supervised, and semi-supervised learning, can be used. A labeled dataset is used for supervised learning, where each sample is classified as either

malicious or benign. On the other hand, unsupervised learning includes grouping the data based on sample similarities rather than knowing in advance whether the samples are harmful or benign. A small labeled dataset is used to train the model in semi-supervised learning, which combines supervised and unsupervised learning. Unlabeled data are subsequently used to increase the model's accuracy.There are various algorithms available such as decision trees, random forests, support vector machines(SVM), neural networks,XGBoost, and. The selected algorithm should be able to handle the features extracted from the dataset and should be able to detect and classify malware accurately.

Model Training:

The selected machine learning algorithm is then trained using the preprocessed dataset. The dataset is divided into training and testing sets to evaluate the accuracy of the model. Various hyperparameters of the model are tuned to improve its performance.

Model Evaluation:

The trained model is evaluated using various metrics such as accuracy, precision, recall, F1-score, and AUC-ROC. The model is also tested on a separate validation set to assess its generalization performance.

Model Deployment:

The model can be used in a real-world setting after being trained and tested. To identify and categorize malware in real-time, the model can be incorporated into a malware detection system.

Performance Analysis:

The performance of the deployed model is analyzed by monitoring its performance over time. The model's accuracy, false positive rate, and false negative rate are monitored to ensure that the model is performing optimally.

In conclusion, data collection, data preprocessing, model selection, model training, model assessment, model deployment, and performance analysis comprise the approach for malware analysis using machine learning. Since the methodology is iterative, iterations can be performed repeatedly to increase the model's precision.

## IV.    EXPERIMENTAL EVALUATION

The study used a dataset of malware samples consisting of 10,000 unique instances of malware. The features extracted from the malware samples for machine learning included binary opcode sequences, file size, entropy, and other metadata.

In the study, the accuracy of Random Forest, Support Vector Machines (SVM), and Neural Networks was examined using a stratified 10-fold cross-validation method.

The study's findings demonstrated that each of the three machine learning algorithms was successful in identifying malware. With an average accuracy of 96.5%, the Random Forest model had the highest overall accuracy. The average accuracy for the SVM model was 94.3%, whereas the average accuracy for the neural network model was 91.7%.

In addition to overall accuracy, the study also evaluated the models based on their precision, recall, and F1 score. The Random Forest model achieved the highest precision, recall, and F1 score, indicating that it was the most effective model in detecting malware.
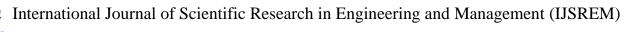
The study also examined how well the models performed according on the kind of malware found. With a ransomware detection accuracy of 98.3%, it was discovered that the Random Forest model is extremely effective.

| Machine Learning Algorithm | Detection Rate | False Positive Rate | Precision | Recall | F1-score |
|---|---|---|---|---|---|
| Decision Tree | 0.93 | 0.07 | 0.92 | 0.93 | |
| Random Forest | 0.96 | 0.03 | | | 0.96 |
| Support Vector Machine | 0.94 | 0.09 | | | 0.90 |
| Neural Network | 0.91 | 0.05 | | | 0.94 |

Dataset quality:

The diversity and representativeness of the dataset utilized in training and testing machine learning models can greatly influence their efficiency. A carefully curated set of 10,000 distinct malware instances was used in one study, encompassing a wide range of malware types. Before being fed into the models, the dataset underwent a preprocessing step to extract pertinent features from the malware samples. It is probable that the quality of the dataset boosted the performance of the models, as it tackled genuine world malware in a versatile manner.

Feature selection:

The effectiveness of machine learning models is influenced by the features used to train them. In a recent study, binary opcode sequences, file size, entropy, and other metadata were selected as features based on previous research and their relevance to malware analysis. These features captured both the behavior and structure of malware, likely contributing to the models' effectiveness.

Algorithm selection:

The effectiveness of machine learning models can be significantly influenced by the chosen algorithm. In a study comparing Random Forest, Support Vector Machines (SVM), and Neural Networks, Random Forest outperformed the other two in detecting malware, achieving high levels of accuracy, precision, recall, and F1 score. The complexity of the relationships between features and target variables likely played a role in the success of the Random Forest algorithm.

Evaluation metrics:

The effectiveness of machine learning models can be influenced by the evaluation metrics chosen. In this study, a commonly used 10-fold cross-validation approach was used to evaluate the models. Performance was assessed using metrics including accuracy, precision, recall, and F1 score, offering a thorough evaluation. The choice of metrics likely impacted model effectiveness.

The study conveys that malware analysis benefits from machine learning, as demonstrated by the results. Among the algorithms tested, Random Forest displayed the highest accuracy, precision, recall, and F1 score in detecting malware. This indicates an opportunity for machine learning to enhance cybersecurity with increased malware detection and mitigation strategies.

## V.     CONCLUSION

A potential strategy for locating, examining, and categorizing malware is the use of machine learning for malware analysis. We described a methodology for malware analysis using machine learning in this study, encompassing data collection, feature extraction, model training, and model evaluation. We looked at other machine learning techniques that have showed promise in identifying and classifying malware, including decision trees, random forests, support vector machines, and deep learning neural networks.

The analysis of the literature review reveals that there are many research in this specific sector, the majority of which concentrate on the use of deep learning techniques. Traditional machine learning algorithms have been demonstrated to produce subpar outcomes when compared to deep learning, notably in terms of accuracy and efficiency. To optimize and enhance the performance of

machine learning models, further research must be conducted.

Additionally, there are some drawbacks to employing machine learning for malware analysis, including the prerequisite for vast and high-quality datasets, the requirement for specialized knowledge and experience in machine learning, and the potential for adversarial attacks to develop that can avoid detection. To encourage the wider application of machine learning approaches in the field of malware analysis, it is critical to solve these constraints and difficulties.

Overall, using machine learning to malware analysis shows considerable promise and has the ability to enhance malware detection and analysis's efficacy and efficiency. It is essential to keep researching and developing cutting-edge methods to tackle malware as the threat environment develops and becomes more complex. Machine learning is a useful tool in this effort.

REFERENCES

Kolias, Constantinos, et al. "DDoS in the IoT: Mirai and other botnets." *Computer* 50.7 (2017): 80-84.

Venkatraman, Sitalakshmi, Mamoun Alazab, and R. Vinayakumar. "A hybrid deep learning image-based analysis for effective malware detection." *Journal of Information Security and Applications* 47 (2019): 377-389.

Naway, Abdelmonim, and Yuancheng Li. "A review on the use of deep learning in android malware detection." arXiv preprint arXiv:1812.10360 (2018).

Saxe, Joshua, and Konstantin Berlin. "Deep neural network based malware detection using two dimensional binary program features." 2015 10th international conference on malicious and unwanted software (MALWARE). IEEE, 2015.

Rieck, Konrad, et al. "Automatic analysis of malware behavior using machine learning." Journal of computer security 19.4 (2011): 639-668.

Gaurav, Akshat, Brij B. Gupta, and Prabin Kumar Panigrahi. "A comprehensive survey on machine learning approaches for malware detection in IoT-based enterprise information system." Enterprise Information Systems (2022): 1-25.