

THE EXPLAINABILITY OF A BLACK-BOX MAL-
WARE DETECTION MODEL USING TAB-
ULAR AND NLP-BASED METHODS

JUNEL ALJE B. ISANAN

SUBMITTED TO THE FACULTY OF THE INSTITUTE OF COMPUTER SCIENCE
UNIVERSITY OF THE PHILIPPINES LOS BAÑOS
IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE
DEGREE OF

BACHELOR OF SCIENCE IN COMPUTER SCIENCE
(Computer Science)

MAY 2024

The thesis hereto attached entitled THE EXPLAINABILITY OF A BLACK-BOX MALWARE DETECTION MODEL USING TABULAR AND NLP-BASED METHODS, prepared and submitted by JUNEL ALJE B. ISANAN, in partial fulfillment of the requirements for the degree of BACHELOR OF SCIENCE IN COMPUTER SCIENCE (Computer Science), is hereby accepted.

DR. RACHEL EDITA O. ROXAS
SP Adviser

RODOLFO C. CAMACLANG III
Member, Guidance Committee

Date Signed

Date Signed

RIZZA D.C. MERCADO
Member, Guidance Committee

Date Signed

Accepted in partial fulfillment of the requirements for the degree of BACHELOR OF SCIENCE IN COMPUTER SCIENCE (Computer Science).

DR. MARIA ART ANTONETTE D. CLARIÑO
Director, Institute of Computer Science

Date Signed

JAMES ROLDAN O. REYES
Dean, College of Arts and Sciences
University of the Philippines Los Baños

Date Signed

BIOGRAPHICAL SKETCH

JUNEL ALJE B. ISANAN is a candidate for the degree of Bachelor of Science in Computer Science at the University of the Philippines Los Baños, with an expected graduation date of May 2024. He has completed coursework in machine learning, data mining, and operating systems. His research focuses on exploring the explainability of black-box malware detection models using tabular and NLP-based methods. He has experience in developing machine learning models for various applications and has participated in research projects related to data analysis and cybersecurity.

Hailing from Polomolok, South Cotabato, Junel's passion for technology and its potential for positive impact led him to pursue data science. As the current Logistics Co-Director of the UP Data Science Society (AY 2023-2024), he actively fosters a community of data enthusiasts. His leadership experience also includes serving as the Scholastics Department Head of the Young Software Engineers' Society (AY 2022-2023), where he taught lessons to aspiring software engineers. Beyond his technical pursuits, Junel finds joy in exploring new places and broadening his horizons.

JUNEL ALJE B. ISANAN

ACKNOWLEDGEMENT

To Doc Rachel, for guiding me in my research. I could not have done this wonderful paper without you.

To Sir JAC, my first ICS faculty friend. You were my source of inspiration as to why I did this paper in the first place, to which I only augmented it with my interest in AI.

To Lanz, my PhD DS friend in UPD, for further igniting my passion towards data science and machine learning. As well as teaching me valuable life skills in handling organizational matters.

To Vincent, who shares life's wisdom with me, I am grateful to meet you.

To Angkol Jonas, my fellow comsci explorer, I enjoyed learning from each other, it is an honor to meet you.

To my Alpha Sigma brothers, maraming salamat sa paggabay sa akin patungo sa tamang landas. Lagi't lagi, Alay Sa Sambayanan.

To Yuan, UP DSSoc Log Co-director, my fellow directors, and other people at UP Data Science Society, let's continue to bring data science to the people.

To the Young Software Engineers' Society members, always uphold leadership, professionalism, and excellence.

Lastly, to my mom and sis. Thank you for giving me this opportunity to learn in UP.

TABLE OF CONTENTS

	<u>PAGE</u>
TITLE PAGE	i
APPROVAL PAGE	ii
BIOGRAPHICAL SKETCH	iii
ACKNOWLEDGEMENT	iv
TABLE OF CONTENTS	v
LIST OF TABLES	vi
LIST OF FIGURES	vii
ABSTRACT	viii
INTRODUCTION	1
Background of the Study	1
Statement of the Problem	2
Objectives of the Study	3
Significance of the Study	3
REVIEW OF LITERATURE	4
Introduction	4
Malware: Purpose, Functions, and Types	4
Antiviruses	4
Data Science	5
Facets of Data Science	6
The Curse of Dimensionality	6
Machine Learning	7
Multilayer Perceptrons	8
Explainable AI	8
Shapley Additive Explanations	9
Materials and Methods	10
Host Machine	11
Development Tools	11
Dataset Used	11
Methodology	12

Data Cleaning	12
Data Preprocessing	13
Machine Learning Models Used	14
XAI Methodology	14
Results and Discussion	15
Model Training and Accuracy Evaluation	15
Visualization of Features	15
Counterfactual Generation	17
Other Predictions	18
SUMMARY AND CONCLUSION	20
LITERATURE CITED	21
APPENDIX I	24
OBFUSCATED SECTION FILES	24
APPENDIX II	25
GENERATION OF COUNTERFACTUAL INFORMATION	25

LIST OF TABLES

<u>TABLE</u>		<u>PAGE</u>
1	Evaluation Metrics For Explanations	9
2	Features in Malware Dataset	12
3	Performance According to Import Size	13
4	MLP Architecture for Malware Detection	14
5	Model Performance and Explainability by Feature	15

LIST OF FIGURES

<u>FIGURE</u>		<u>PAGE</u>
1	Methodological Workflow of the Study	10
2	Sample data point predicted using the import feature	16
3	The same data point in Fig. 2 after counterfactual information was generated. VirtualAlloc was removed from the datapoint	17
4	Another datapoint predicted using the headers feature and discovered it was a packed malware (due to the pkd feature which was custom-built for the study)	19
5	Another datapoint predicted using the strings feature, but it doesn't seem to be coherent at all	19
6	Photo of obfuscated sections	24
7	Generation of counterfactual information	25

ABSTRACT

JUNEL ALJE B. ISANAN, University of the Philippines Los Baños, MAY 2024.

THE EXPLAINABILITY OF A BLACK-BOX MALWARE DETECTION MODEL USING TABULAR AND NLP-BASED METHODS

Major Professor: DR. RACHEL EDITA O. ROXAS

The exponential surge in malware types, from 100 million in 2012 to 700 million in 2018, poses a significant threat to the cybersecurity sector. While deep learning models have demonstrated impressive accuracy in malware detection and classification, even exceeding 99% in recent studies, the underlying features driving these predictions remain elusive. This lack of transparency hinders user trust in the models' decisions. To address this, this research explores the use of XAI tools to analyze the inner workings of black-box malware detection models trained on tabular and NLP-based data from the EMBER dataset. By uncovering the influential features in single data points, this study aims to enhance model transparency, foster trust in malware detection predictions, and pave the way for improved detection methods and more effective threat analysis in the ever-evolving landscape of cybersecurity.

INTRODUCTION

Background of the Study

Malware or "malicious software" is an umbrella term that describes any malicious program or code that is harmful to systems when computers execute (MalwareBytes, 2024). Malware, when attached to computer systems, can bring about serious damage to computer systems necessary for human advancement and conveniences such as registration systems, queueing systems, or even automated sensing systems, among many others. Report shows that 560,000 pieces of malware are created every day, wreaking havoc on personal and company computers alike (Palatty, 2023).

Just in recent years, there has been a surge of public interest in the field of artificial intelligence. Tools like ChatGPT, Google Bard, and Midjourney have been constantly utilized by the general public. The mechanism of how AI tools work, particularly AI that uses machine learning, involves learning based on observations made on the training data (SAS, n.d.). In recent studies, there has been malware detection and classification models that reach near 100% which leaves very little room for research in terms of improving the model accuracy for malware detection models.

However, tabular and NLP-based methods may be crucial for understanding malware detection models in general (Saxe & Sanders, 2018) especially black-box models such as multi-layer perceptrons (MLP) whose processing mechanisms are difficult to understand and which results are more difficult to interpret in general (Dobson, 2023).

Statement of the Problem

There is an ever-growing threat of malware, encompassing viruses and Trojans, among many others, posing a range of damaging consequences to individual computers, corporate networks, and the data they contain. The extent of the damage varies from subtle disruptions, where users may remain unaware, to catastrophic failures leading to network breakdowns, data loss, or even theft. These challenges people face in the field of the software industry.

According to a report, it is estimated that the financial damages brought by ransomware alone would reach up to \$250 billion by 2031 (Sausaulito, 2023).

Another report states that malware is the main cause of the loss of data integrity. This includes data loss and theft. The damage would be magnified if the data was stolen and/or encrypted, especially if it involved personal information or intellectual property (Encyclopedia, n.d.).

In recent years, applying machine learning in the field of malware detection has reached 98% to 99% accuracy rates using black-box models (M. & Sethuraman, 2023) (Yuksel & Ar, 2023). However, explainability of results in black-box malware detection models are yet to be explored fully to gain user trust and improve accuracy of the model (Li, Sun, Huang, & Chen, 2024).

Thereby, this research seeks to address the gap in this field by exploring methods to improve the explainability of black-box models for malware detection, while maintaining or even enhancing their accuracy.

Objectives of the Study

This study aims to use an XAI approach in analyzing black-box models.

Specifically, it aims to:

1. Train machine learning models from the data
2. Evaluate model accuracy as secondary methods for establishing user trust in black-box models and use XAI tools to determine which features influenced the prediction
3. Produce feature visualizations of individual predictions
4. Generate counterfactual information on some predictions to assess and verify sensitivity.

Significance of the Study

Malware continues to plague us in our day-to-day surfing of the cyberspace.

The analysis of malware data with black-box models allows the malware to be identified immediately upon retraining should a new strain of malware with similar characteristics appear. Thereby, identification of zero-day malware does not need to be constantly searched by cybersecurity experts to calculate heuristics constantly (Saxe & Sanders, 2018).

But as black-box models reach near-perfect accuracy in detecting, with one model reaching 99.8% (Yuksel & Ar, 2023). There must be a clear shift in priority of researching from improving model to exploring explainability methods to establish user trust in the model's predictions as it is in finance and healthcare (Liu, Tantithamthavorn, Li, & Liu, 2022).

REVIEW OF LITERATURE

Introduction

Today marks an era of modern technology with the widespread accessibility of mobile phones, laptops, and PCs. Technologies have advanced to generate various content, including synthetic personas and prefabricated essays on diverse topics. The convergence of technology and data has allowed mankind to emulate reality and intelligence.

Despite these advancements, questions about the safety of our data and the security of programs persist (Software, 2023). The historical encounter with the ILOVEYOU virus underscores the evolution of destructive programs, collectively known as malware.

Malware: Purpose, Functions, and Types

Malware, a portmanteau for "malicious software," is a type of software that steals sensitive information and damages or destroys computers/systems (Cisco, 2023). The purposes for creating malware include intelligence and intrusion, disruption, extortion, destruction or vandalism, unauthorized use of computer resources, and monetary gain (Cisco, 2023).

Antiviruses

Antivirus programs, developed to counter challenges brought about by viruses, use scanning, generic detection, or heuristic detection to identify virus signatures (Muchelule, Wanjala, & Misiko, 2017). However, antiviruses may not effectively contain other types of malware, such as ransomware (Mixon, 2022).

Data Science

Labelled as the "sexiest job of the 21st century", data science has garnered significant attention (Davenport & Patil, 2024). Cielen et al. defines data science as the use of methods to analyze large volumes of data and extract meaningful insights (Cielen, B., & Mohamed, 2016). On the other hand, Kudande et al. view the same field as a scientific discipline that applies scientific methods, algorithms, systems, and processes to derive valuable information from various structured and unstructured data sources (Kudande, 2020).

Cielen et al.(Cielen et al., 2016) outline a typical data science process:

1. **Setting the research goal:** Defining the objectives and questions the analysis aims to answer.
2. **Retrieving:** Gathering the relevant data from various sources.
3. **Data preparation:** Cleaning, transforming, and organizing the data for analysis.
4. **Data exploration:** Analyzing and visualizing the data to uncover patterns, trends, and relationships.
5. **Data modeling:** Developing predictive or explanatory models to gain deeper insights.
6. **Presentation and automation:** Communicating findings and potentially automating the data science process.

This process, as emphasized by Cielen et al. (Cielen et al., 2016), is iterative. It often involves revisiting earlier steps to address challenges, refine data, or improve results. This study will further elaborate on these steps in the methodology section, drawing upon insights from additional research and textbooks.

Facets of Data Science

Cielen et al. (Cielen et al., 2016) states that there are multiple facets of data. These are:

- Structured
- Unstructured
- Natural language
- Machine-generated
- Graph-based
- Audio, image, and video
- Streaming

It is important to pick the right facet of data according to the research problem as each facet requires different treatments for preprocessing, modeling, and measuring data. In the case of tabular data, cleaning data would mean imputing based on the mean, median, and mode of column values (Fan, Chen, Wang, Wang, & Huang, 2021) whereas text data would be cleaned by tokenizing, removing stop words, and stemming (Gurusamy & Kannan, 2014).

The Curse of Dimensionality

Ullah defines high-dimensional data as data with a large number of features, often represented as columns in a tabular dataset (Ullah, 2019). Working with such data presents several challenges:

- **Computational Complexity:** Analyzing high-dimensional data can be computationally expensive. The distance required to search for n nearest neighbors increases with each additional dimension.
- **Limited Human Visualization:** High-dimensional data does not lend itself well to human visualization, making it difficult to analyze the underlying phenomena visually.
- **Redundancy and Noise:** As the number of features increases, so does the potential for redundancy, which introduces noise into the dataset.
- **Exponential Search Space:** The search space within the dataset grows exponentially with added features, causing some algorithms to become unstable.

Machine Learning

Machine learning is a fusion of concepts from statistics, mathematics, and numerical analysis to gain insights from data, enabling us to summarize, visualize, group, or predict phenomena within a given dataset (Deuschle, 2019).

According to Abu-Mostafa et al. (Abu-Mostafa, Magdon-Ismail, & Lin, 2012) in their book *Learning From Data: A Short Course*, there are three types of machine learning:

1. **Supervised Learning:** This technique involves mapping an input to an output, where the output is solely determined by the input. Supervised learning finds extensive use in computer vision, with digit recognition systems trained on the MNIST dataset being a prime example.
2. **Reinforcement Learning:** This technique considers both the player input and the

environment state to assess whether transitioning to a new state is advantageous or not. Chess engines (Maharaj, Polson, & Turk, 2022) exemplify reinforcement learning by grading moves based on user input and the current game state, indicating which player has an advantage and by how much.

3. **Unsupervised Learning:** In this setting, the dataset relies solely on input without any output information. Data points learn independently by discovering patterns and structures within the dataset. An example of unsupervised learning can be seen in market basket analysis, where items are associated based on their frequency of co-occurrence.

Multilayer Perceptrons

Multi-layer Perceptrons (MLPs) are a form of feedforward neural network commonly used in supervised learning tasks. They are characterized by three distinct types of densely-connected layers: input, output, and hidden layers serving as the computational core. The output of an MLP is determined depending on the various activation functions utilized (ReLU, Tanh, etc.). Furthermore, the multi-layer perceptron is considered to be a black-box model due to its complex architecture and mysterious decision-making capabilities (Dobson, 2023).

Explainable AI

Over the past few years, there is exponential growth in AI production and research with the rise of enormous data and blazing fast computing power. That being said, black-box models can be trained for up to 99+% for a reasonable amount of time given those two conditions earlier. However, black-box models lack explainability given that the values given are only the weights without any transparent decision-making process made by white-box models. Hence, XAI on black-box models rely heavily on approximation and mathematical

nuances to calculate probable contribution of features. Below are the some of the listed metrics for XAI (Coroama & Groza, 2022):

Table 1. Evaluation Metrics For Explanations

Metric	Description
Sensitivity	The degree to which the explanation is affected by insignificant perturbations from the test point.
Faithfulness	Precisely computes which feature has the most impact on the model output when individually changed.
Simplicity	The ability to choose only the necessary and sufficient features for explaining the prediction.

Shapley Additive Explanations

SHAP or SHapley Additive exPlanations is mathematically based on game theory, the Shapley value of a feature is the average expected marginal contribution to the model's decision after all combinations are, however, these values are made possible to be calculated by assuming that the features are independent of each other (Belle & Papantonis, 2021).

Materials and Methods

This section presents the methodology workflow used for this study as presented in Figure 1. The host machine, development tools, and dataset used for this study are discussed in detail. Initial stages include data cleaning and data preprocessing in preparation for training the black-box malware detection model. After training, the XAI methodology and visualization are outlined.

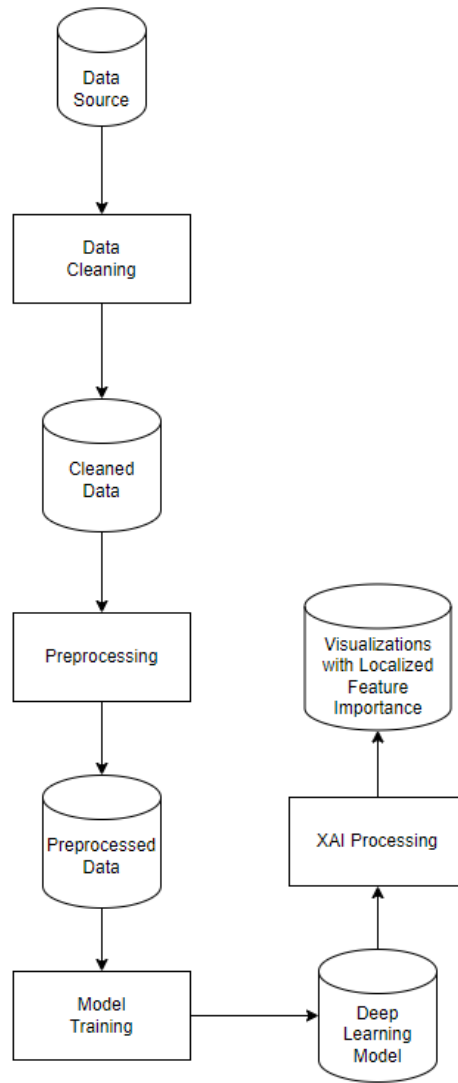


Figure 1. Methodological Workflow of the Study

Host Machine

- **Operating System:** Windows 10 Pro 64-bit (10.0, Build 19045) (19041.vb_release.191206-1406)
- **System Manufacturer:** Colorful Technology And Development Co.,LTD
- **System Model:** H410M-T PRO
- **Processor:** Intel(R) Core(TM) i3-10100F CPU @ 3.60GHz (8 CPUs), 3.6GHz
- **PC Memory:** 16 GB
- **Google Colab Pro Memory:** 52 GB

Development Tools

- **Jupyter Notebook:**

A modular computing environment/text editor where machine learning frameworks / data science related stuff

- **Git:**

This software provides version control, it is responsible for releasing which versions of the project is stable.

- **Python:** The primary programming language used for developing machine learning models. The language is embedded within Jupyter Notebook.

Dataset Used

The dataset used for conducting this study was the EMBER dataset. The EMBER dataset provides eight groups of stringified features that allow for static analysis (Anderson

& Roth, 2018). These features are categorized as follows:

Table 2. Features in Malware Dataset

Feature Group	Description
sha256 (id)	Unique identifier for each sample
appeared	Date of appearance
label	Target variable (malicious/benign)
general	<i>(Discarded)</i>
header	<i>(Discarded)</i>
section	<i>(Section information about the executable)</i>
imports	<i>(Import functions used by the executable)</i>
exports	<i>(Discarded)</i>
histogram	<i>(Discarded)</i>
byteentropy	<i>(Discarded)</i>
strings	<i>(String information of the executable)</i>

For the purposes of this study, the features **section**, **imports**, and **strings** were primarily considered as they were tabular and natural language in nature (Saxe & Sanders, 2018).

Methodology

Data Cleaning

The feature columns header, imports, and strings came as stringified JSON objects. The study used Jupyter Notebook, pandas, and the JSON library in Python to parse and clean the dirty data in order for these source dataset to become useful.

Retrieval of section names in the were limited as some as the values in the section were obfuscated or packed. The technique was to obtain the usual entry points and section names. Such examples would be .text or .data section. For which code and variable names are written on those sections.

Data Preprocessing

Imports were treated as natural language and was therefore tokenized according to frequency using CountVectorizer. The output is a tabularized feature of keywords. Other features are tabularized already, so normalization of values were done for features in Table 2.

Due to memory and time limitations, only the top 500 out of the 200k+ most-frequently called imports were utilized.

The basis for cutting down to 500 imports was term frequency of each keywords in order to maximize information capture and at the same time reduce dimensionality for training to be more feasible and faster (Ganesan, 2021). Below are the training accuracies gathered after training some varied number of imports.

Table 3. Performance According to Import Size

Import Size	Training Accuracy
500	88.02%
250	86.47%
100	80.58%

Based on the given values, 500 imports would yield better results than the other sizes.

The features are converted using one-hot encoding since they are categorical in nature. This is to ensure that each unique values are represented numerically and are understood by the model.

Machine Learning Models Used

The study used a multilayer perceptron (MLP) in order to allow for a thorough learning of pattern in between datasets (Alzubaidi et al., 2021). An n number of nodes in the input layer is implemented, followed by 64 nodes of hidden layer and lastly, 2 nodes for the output layer, classifying benign and malicious probability values. The corresponding activation functions used for the hidden layer were ReLU and sigmoid was implemented for the output layer. It is the purpose of the study to use black-box models such as multi-layer perceptrons to extract insights due to their high accuracy (Abu-Mostafa et al., 2012) but opaque decision-making methodology (Dobson, 2023)

Table 4. MLP Architecture for Malware Detection

Layer Type	Configuration	Output Shape
DenseFeatures	Depending on the size of feature	(Num. Features)
Dense	64 neurons, ReLU activation	(64)
Dropout	0.2 dropout rate	(64)
Dense	2 neurons, Sigmoid activation	(2)

XAI Methodology

The XAI Methodology used in this study is Shapley Additive Explanations (SHAP) and the prediction is displayed on a bar chart in order of descending individual feature importance. SHAP cannot attribute the totality of the feature on a global scale. The prediction is done post-hoc, meaning that the prediction will be done after the model has been trained and is fed a number of individual data points for features to be predicted.

Results and Discussion

This section addresses the objectives of the study and is subdivided into three parts:

1. Model Training and Accuracy Evaluation
2. Visualization of Features
3. Counterfactual Generation
4. Other Predictions

Model Training and Accuracy Evaluation

On Table 4, the post-hoc explainability results of the model is shown. This means that the explanation comes only after a single datapoint has been predicted.

Table 5. Model Performance and Explainability by Feature

Feature	Model Used	Accuracy	XAI Method
Imports	MLP	88.02%	SHAP
String	MLP	84.26%	SHAP
Header	MLP	76.04%	SHAP

Since the study focuses on explainability, Table V is merely for analytics purposes. Most of the study results will delve into the explainability of individual predictions.

Visualization of Features

This subsection presents the visualizations of the XAI model. Visualizations are presented in a descending order based on magnitude of local feature importance.

To interpret the values on each visualization, features leaning to the orange color imply that these features contribute to the prediction of a single data point malware. On the other hand, features leaning to the blue color imply that these features contribute to the benignity of the data point. The longer the bar, the greater the share of the feature's prediction towards a certain value.

A sample data point predicted using the imports feature as shown in Figure 2.

Actual: 1 | Predicted: 0.926

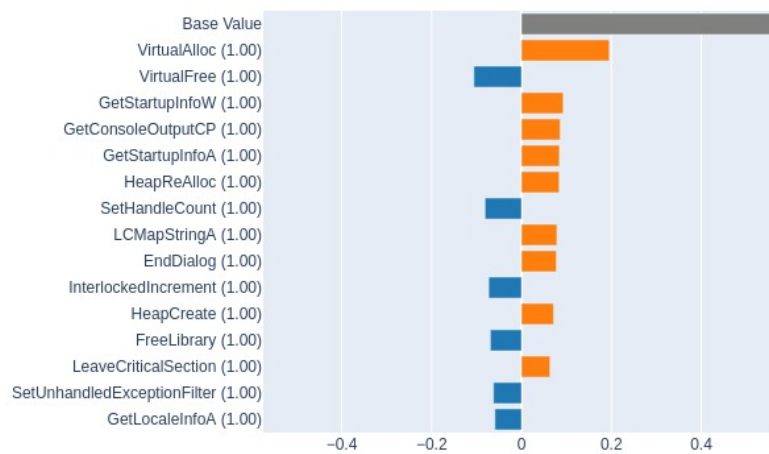


Figure 2. Sample data point predicted using the import feature

In Figure 2, it is shown that 18.5% of the prediction has been attributed to VirtualAlloc alone. It implies that on the absence of this single feature, the prediction of this malware will change drastically.

Counterfactual Generation

In this subsection, a counterfactual generation data point is created to test the sensitivity and faithfulness of the data. When the most salient feature is removed, it should significantly lower the prediction of the data point by the model.

Based on Figure 2, it is 92.6% certain that the data point is malware. On Figure 3, however, the absence of VirtualAlloc is simulated while the other features remain intact. As a result, the prediction drastically changes to 40.6% as evidenced in Figure 3.

Actual: 1 | Predicted: 0.406



Figure 3. The same data point in Fig. 2 after counterfactual information was generated. VirtualAlloc was removed from the datapoint

By changing just one of the most important features of the local prediction, not only

did it have a drastic change in prediction, but also the model predicts benignity in the code despite the ground truth stating otherwise.

Other Predictions

Regarding other predictions, Figures 4 and 5 are samples of individual predictions for other features.

As shown on Figure 4, the prediction implies that most of the features seem to connive that the said data point is malicious based on the given graph.

On the other hand, Figure 5 doesn't seem to be coherent at all since it violates the XAI principle of simplicity. Each feature should have semantically coherent implications.

Furthermore, the study also produced 10 predictions per feature as a way to explain other local predictions.

Actual: 1 | Predicted: 0.934

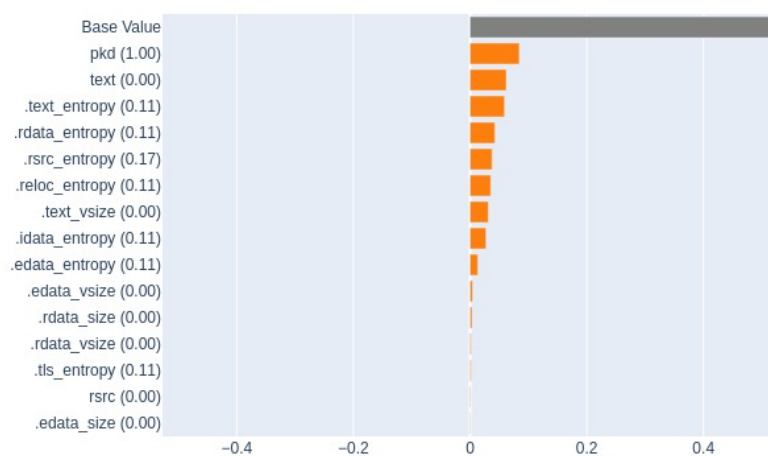


Figure 4. Another datapoint predicted using the headers feature and discovered it was a packed malware (due to the pkd feature which was custom-built for the study)

Actual: 1 | Predicted: 0.973

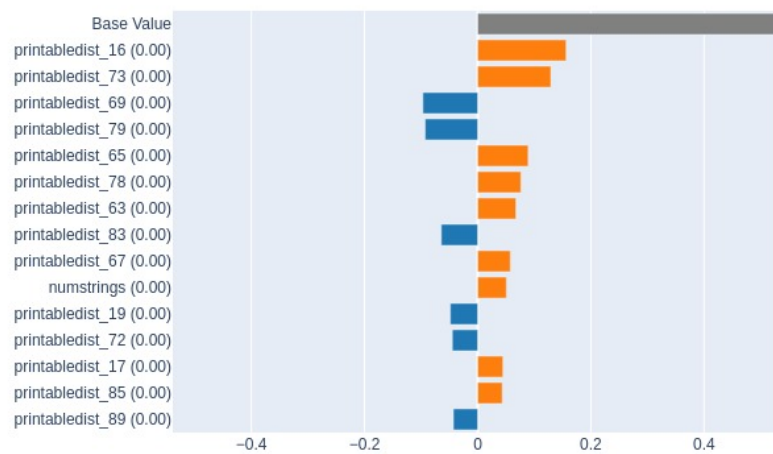


Figure 5. Another datapoint predicted using the strings feature, but it doesn't seem to be coherent at all

SUMMARY AND CONCLUSION

With this study, black-box malware detection models with tabular and NLP-based features were interpreted locally using SHAP. This will allow cybersecurity experts to comment and improve on existing machine learning models.

As for the future works, with Kolmogorov-Arnold Networks (KAN) as the emerging trend for deep learning, where activation functions are placed on edges instead of nodes. This type of neural network may offer more explainability due to the individual weights placed on edges instead of nodes (Vaca-Rubio, Blanco, Pereira, & Caus, 2024). For the meantime, while the capabilities of Kolmogorov-Arnold Networks are not yet fully researched, the study could be enriched further with neural networks with the help of domain experts to determine whether such features predicted by the model using game-theoretic approaches are relevant using counterfactual generation.

Moreover, the explainability of machine learning models can easily be transferred to other sectors such as healthcare, agriculture, social sciences, and humanities. XAI will help machine learnings improve by counterchecking on features / pixels that influenced such decisions. The source code can easily be tweaked according to need.

And lastly, validate the truthfulness of the features assigned by the model with a domain expert.

LITERATURE CITED

- ABU-MOSTAFA, Y. S., MAGDON-ISMAIL, M., & LIN, H.-T. (2012). *Learning from data: A short course*. AMLBook.com.
- ALZUBAIDI, L., ZHANG, J., HUMAIDI, A. J., AL-DUJAILI, A., DUAN, Y., AL-SHAMMA, O., ... FARHAN, L. (2021, Mar). Review of deep learning: Concepts, cnn architectures, challenges, applications, future directions. *Journal of Big Data*, 8(1). doi: 10.1186/s40537-021-00444-8
- ANDERSON, H. S., & ROTH, P. (2018). *Ember: An open dataset for training static pe malware machine learning models*.
- BELLE, V., & PAPANTONIS, I. (2021, Jul). Principles and practice of explainable machine learning. *Frontiers in Big Data*, 4. doi: 10.3389/fdata.2021.688969
- CIELEN, D., B., M. A. D., & MOHAMED, A. (2016). *Introducing data science: Big data, machine learning, and more, using python tools*. Manning.
- CISCO. (2023, Jul 21). *What is malware? - definition and examples*. Retrieved from <https://www.cisco.com/site/us/en/products/security/what-is-malware.html>
- COROAMA, L., & GROZA, A. (2022). Evaluation metrics in explainable artificial intelligence (xai). *Communications in Computer and Information Science*, 401–413. doi: 10.1007/978-3-031-20319-0_30
- DAVENPORT, T. H., & PATIL, D. (2024, Mar). *Data scientist: The sexiast job of the 21st century*. Retrieved from <https://hbr.org/2012/10/data-scientist-the-sexiest-job-of-the-21st-century>
- DEUSCHLE, W. J. (2019). *Undergraduate fundamentals of machine learning* (Bachelor's thesis). Harvard College.
- DOBSON, J. E. (2023, Nov). On reading and interpreting black box deep neural networks. *International Journal of Digital Humanities*, 5(2–3), 431–449. doi: 10.1007/s42803-023-00075-w
- ENCYCLOPEDIA, K. I. (n.d.). Damage caused by malware. Retrieved from <https://encyclopedia.kaspersky.com/knowledge/damage-caused-by-malware/>
- FAN, C., CHEN, M., WANG, X., WANG, J., & HUANG, B. (2021, Mar). A review on data preprocessing techniques toward efficient and reliable knowledge discovery

from building operational data. *Frontiers in Energy Research*, 9. doi: 10.3389/fenrg.2021.652801

GANESAN, K. (2021, Aug). *What is term frequency?* Retrieved from [https://www.opinosis-analytics.com/knowledge-base/term-frequency-explained/#:~:text=Term%20frequency%20\(TF\)%20means%20how,about%20how%20you%20define%20it.](https://www.opinosis-analytics.com/knowledge-base/term-frequency-explained/#:~:text=Term%20frequency%20(TF)%20means%20how,about%20how%20you%20define%20it.)

GURUSAMY, V., & KANNAN, S. (2014, 10). Preprocessing techniques for text mining..

KUDANDE, P. D. (2020). A review paper on the formation of data science. *International Journal of Engineering Applied Sciences and Technology*. Retrieved from <https://tinyurl.com/muxs43ud>

LI, M., SUN, H., HUANG, Y., & CHEN, H. (2024, Feb). Shapley value: From cooperative game to explainable artificial intelligence. *Autonomous Intelligent Systems*, 4(1). doi: 10.1007/s43684-023-00060-8

LIU, Y., TANTITHAMTHAVORN, C., LI, L., & LIU, Y. (2022, Oct). Explainable ai for android malware detection: Towards understanding why the models perform so well? *2022 IEEE 33rd International Symposium on Software Reliability Engineering (ISSRE)*. doi: 10.1109/issre55969.2022.00026

M., G., & SETHURAMAN, S. C. (2023, Feb). A comprehensive survey on deep learning based malware detection techniques. *Computer Science Review*, 47, 100529. doi: 10.1016/j.cosrev.2022.100529

MAHARAJ, S., POLSON, N., & TURK, A. (2022, Apr). Chess ai: Competing paradigms for machine intelligence. *Entropy*, 24(4), 550. doi: 10.3390/e24040550

MALWAREBYTES. (2024, Apr). Retrieved from <https://www.malwarebytes.com/malware>

MIXON, E. (2022, Apr). *Ransomware vs. malware: What's the difference?* Retrieved from <https://www.blumira.com/ransomware-vs-malware/>

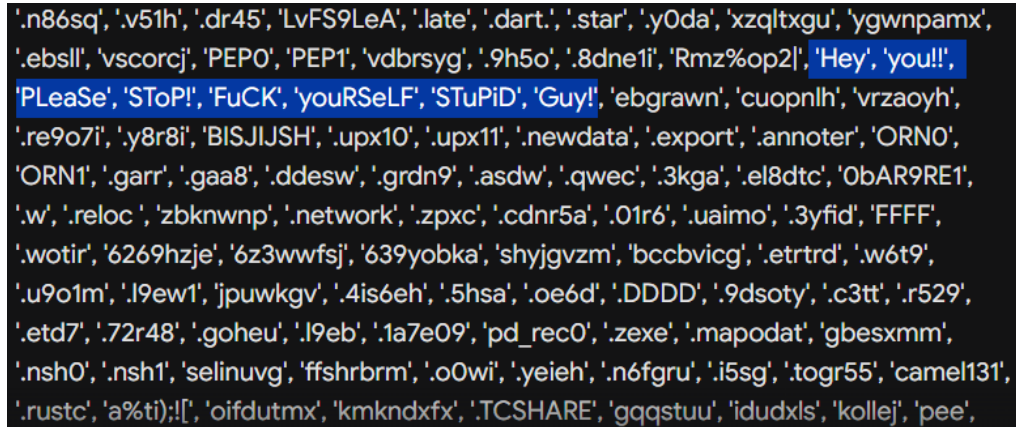
MUCHELULE, Y., WANJALA, N., & MISIKO, J. (2017, January). Review of viruses and antivirus patterns. *Journal of Computer Science and Technology*.

PALATTY, N. J. (2023, Dec). *30+ malware statistics you need to know in 2024*. Retrieved from <https://www.getastra.com/blog/security-audit/malware-statistics/>

- SAS. (n.d.). Machine learning: What it is and why it matters. Retrieved from https://www.sas.com/en_us/insights/analytics/machine-learning.html
- SAUSAULITO, C. (2023). Global ransomware damage costs predicted to exceed \$265 billion by 2031. *Cybercrime Magazine*. Retrieved from <https://cybersecurityventures.com/global-ransomware-damage-costs-predicted-to-reach-250-billion-usd-by-2031/>
- SAXE, J., & SANDERS, H. (2018). *Malware data science: Attack detection and attribution*. No Starch Press.
- SOFTWARE, C. P. (2023). Biggest cyber security challenges in 2023. *Check Point Software*. Retrieved from <https://www.checkpoint.com/cyber-hub/cyber-security/what-is-cybersecurity/biggest-cyber-security-challenges-in-2023/>
- ULLAH, F. (2019). Retrieved from https://web.lums.edu.pk/imdad/pdfs/CS5312_Notes/CS5312_Notes-11-Curse_of_Dimensionality.pdf
- VACA-RUBIO, C. J., BLANCO, L., PEREIRA, R., & CAUS, M. (2024). Kolmogorov-arnold networks (kans) for time series analysis.
- YUKSEL, A. K., & AR, Y. (2023, Aug). A machine learning approach to malware detection using application programming interface calls (mdapi). *Traitement du Signal*, 40(4), 1511–1520. doi: 10.18280/ts.400419

APPENDIX I

OBFUSCATED SECTION FILES



```
'n86sq', 'v51h', 'dr45', 'LvFS9LeA', 'late', 'dart.', 'star', 'y0da', 'xzqltxgu', 'ygwnpamx',  
'ebsll', 'vscorcj', 'PEP0', 'PEP1', 'vdbrsyg', '9h5o', '8dne1i', 'Rmz%op2|', 'Hey', 'you!!',  
'PLeaSe', 'SToP!', 'FuCK', 'youRSeLF', 'STuPiD', 'Guy!', 'ebgrawn', 'cuopnlh', 'vrzaoyh',  
're9o7i', 'y8r8i', 'BISJIJSH', 'upx10', 'upx11', 'newdata', 'export', 'anoter', 'ORN0',  
'ORN1', 'garr', 'gaa8', 'ddesw', 'grdn9', 'asdw', 'qwec', '3kga', 'el8dtc', 'ObAR9RE1',  
'w', 'reloc', 'zbknwnp', 'network', 'zpxc', 'cdnr5a', '01r6', 'uaimo', '3yfid', 'FFFF',  
'wotir', '6269hzje', '6z3wwfsj', '639yobka', 'shyjgvzm', 'bccbvicg', 'etrtrd', 'w6t9',  
'u9o1m', 'l9ew1', 'jpuwkgv', '4is6eh', '5hsa', 'oe6d', 'DDDD', '9dsoty', 'c3tt', 'r529',  
'etd7', '72r48', 'goheu', 'l9eb', '1a7e09', 'pd_rec0', 'zexe', 'mapodat', 'gbesxmm',  
'nsh0', 'nsh1', 'selinuv', 'ffshrbrm', 'o0wi', 'yeieh', 'n6fgru', 'i5sg', 'togr55', 'camel131',  
'rustc', 'a%ti);!]', 'oifdutmx', 'kmkndxfx', 'TCSHARE', 'gqqstuu', 'idudxls', 'kollej', 'pee',
```

Figure 6. Photo of obfuscated sections

Figure 6 represents some of the sections that were obfuscated out of the thousands retrieved from the dataset. Indicating that the section texts might be hidden to prevent from being reverse-engineered.

APPENDIX II

GENERATION OF COUNTERFACTUAL INFORMATION

```
✓ [216] 1 sixth_row = X_predict.head(6).tail(1)
0s

✓ [221] 1 sixth_row['VirtualAlloc']
0s

⇒ 5 1
   Name: VirtualAlloc, dtype: int64

[135] 1 sixth_row['VirtualAlloc'] = 0
```

Figure 7. Generation of counterfactual information

Figure 7 represents the method on how to remove the most salient feature in Figure 2, which is VirtualAlloc.