# Malware Detection in Portable Executable (PE) Files Using Supervised Machine Learning (Applied Artificial Intelligence in Cyber Security)

1 author:

Muhammad Qasim Raza
University of South Wales
**11** PUBLICATIONS   **133** CITATIONS

Faculty of Computing, Engineering and Science

National Cyber Security Academy

University of South Wales

MSc Applied Cyber Security

"Malware Detection Using Supervised Machine Learning"

(Applied Artificial Intelligence in Cyber Security)

IY4S734

Submitted To:

Muhammad Abdullah Awais

Submitted By:

Muhammad Qasim Raza

30115377

# Table of Contents

# Table of Figures:

# 1  Introduction to AI in Cybersecurity

Artificial intelligence (AI) is transforming cybersecurity by delivering innovative tools and approaches for mitigating online attacks. Finding and stopping malware, malicious software intended to damage networks, steal data, or interfere with services, is one of the biggest issues in cybersecurity. Because AI makes threat identification quicker, more accurate, and more flexible, it is proving to be extremely useful in the fight against such threats.

## 1.1  How AI Helps in Cybersecurity

AI plays a crucial role in enhancing cybersecurity by analyzing vast amounts of data and identifying patterns that humans might miss. In the context of malware detection, AI can spot unusual behaviors in files, network traffic, and system activities that might indicate a cyberattack. Machine learning, a branch of AI, is particularly effective in detecting new, unknown malware, often referred to as zero-day threats.

For instance, traditional cybersecurity systems rely on databases of known malware signatures to identify threats. However, cybercriminals constantly create new malware variants that escape from these systems. AI overcomes this limitation by recognizing suspicious activities or anomalies, even if the malware is entirely new. This proactive approach allows organizations to detect and respond to threats much faster than traditional methods.

## 1.2  Challenges in Identifying Malware

Identifying malware is not an easy task, as cybercriminals use sophisticated techniques to disguise their attacks. Some malware is designed to remain dormant until triggered by a specific event, making it harder to detect. Others use polymorphic techniques to change their code structure constantly, evading signature-based detection methods.

Another challenge is the sheer volume of data that cybersecurity systems must analyze. Large organizations often process millions of files, emails, and network logs daily. Manually reviewing this data for signs of malware is nearly impossible. Moreover, distinguishing between legitimate activities and malicious ones can be tricky, as attackers often mimic normal user behavior to avoid detection.

## 1.3 Supervised Learning for Malware Detection

Supervised learning, a type of machine learning, can be a powerful tool in detecting malicious files. In supervised learning, a model is trained on labeled data, which means it learns from examples where the outcome (e.g., whether a file is malicious or safe) is already known. This training process enables the model to recognize patterns and features associated with malware. (Priyanka Jadav, April 27, 2024)

For example, a supervised learning model might analyze features such as the size of a file, its code structure, or the type of system calls it makes. By comparing these features with its training data, the model can predict whether a new file is likely to be malicious. This approach is highly effective because it can learn from past examples and improve over time.

Supervised learning also reduces the likelihood of false positives, which occur when a system incorrectly flags safe files as malicious. By continuously refining its understanding of what constitutes malware, the model becomes more accurate and reliable.

AI, and particularly supervised learning, is transforming the way cybersecurity professionals detect and respond to malware. Despite the challenges posed by evolving cyber threats, AI provides tools to identify even the most sophisticated attacks. As technology advances, AI-driven cybersecurity solutions will continue to play a vital role in protecting individuals and organizations from malicious actors.

## 2   Understanding PE File Structure

Portable Executable (PE) files are a fundamental part of Windows operating systems, serving as the format for executable files, dynamic link libraries (DLLs), and other binary programs. These files contain essential information for the system to load and execute programs. Understanding the structure of PE files is critical in malware analysis because many malicious programs disguise themselves as legitimate PE files to bypass security measures.

### 2.1   What is the PE File Structure?

The PE file structure is like a roadmap that tells the operating system how to execute a program. It is divided into sections, each serving a specific purpose, such as holding executable code, data, or information about external libraries. Key components of a PE file include: (Microsoft, 02/29/2024)

- **Headers:** Contain metadata about the file, such as its type (executable or DLL), timestamp, and entry point address.

- **Sections:** Hold the main content of the program, including code, resources, and data. Common sections include .text (executable code), .data (initialized data), and .rdata (read-only data).

- **Imports and Exports:** Lists of functions and libraries the file depends on (imports) or provides to other programs (exports).

These structural elements are crucial for malware analysis because they provide insights into the file's behavior and purpose.

### 2.2   Key Features for Malware Analysis

Several features of PE files can be analyzed to determine whether a file is malicious or legitimate. These features are used by machine learning models to detect patterns indicative of malware. Some of the most common features include:

### 2.2.1  Number of Sections:

    a.   Legitimate PE files typically have a standard number of sections.

b. Malware authors often add unusual sections or modify existing ones to hide malicious code.

### 2.2.2 Section Sizes:

a. Normal PE files have sections with predictable sizes.

b. Suspiciously large or tiny sections can indicate packed or encrypted malware.

### 2.2.3 Imports:

a. Malware often imports functions related to suspicious activities, such as CreateRemoteThread or VirtualAlloc, which are used to inject code into other processes.

b. The absence of expected imports (e.g., standard Windows API functions) can also be a red flag.

### 2.2.4 Exports:

a. Files with unusual or excessive exports may be designed to interact with multiple programs in an unintended way, often seen in malicious DLLs.

### 2.2.5 Entropy:

a. Entropy measures the randomness of data within a file.

b. High entropy sections might indicate that the file is packed or encrypted, a common tactic to avoid detection.

### 2.2.6 File Metadata:

a. Information like compile timestamps and file signatures can be checked against known malware databases.

b. Invalid or unusual metadata can hint at tampering.

## 2.3 How These Features Help Machine Learning Models

Machine learning models use these features to identify patterns that distinguish malware from normal files. For instance:

- By analyzing the combination of section sizes and entropy, the model can detect files that use packing techniques to hide malicious code.

- Unusual imports or exports can signal the presence of functions commonly used by malware, prompting further investigation.

- The number and structure of sections provide clues about whether the file was altered from its original state.

By feeding the model large amounts of labeled data (e.g., known malware and legitimate files), it learns to recognize the subtle differences and make predictions about new, unseen files. This automated analysis is especially valuable in detecting polymorphic or previously unknown malware.

Understanding the PE file structure is a cornerstone of malware analysis. Features like section sizes, imports, exports, and entropy provide vital clues about a file's behavior, enabling machine learning models to make informed decisions. By leveraging these insights, cybersecurity professionals can detect and mitigate threats more effectively, ensuring a safer computing environment.

# 3   Data Set

Building a dataset for malware analysis requires finding reliable sources for both malware and normal PE (Portable Executable) files. Ensuring ethical practices and maintaining safety are crucial throughout this process, as mishandling malware can cause unintentional harm. (Iryna Sydorenko, December 21, 2023)

## 3.1   Where and How to Find Malware and Normal PE Files

### 3.1.1   Malware Samples:

a. **Online Repositories:** Reputable repositories like VirusTotal, MalwareBazaar, and CIC Andmal2020 provide malware samples for research purposes.

b. **Access Requirements:** These platforms usually require you to register and agree to ethical usage policies. Some may only be accessible to verified researchers or organizations.

c. **Safety Measures:** Always download and handle malware samples within a controlled environment, such as a virtual machine (VM) isolated from your main network.

### 3.1.2   Normal PE Files:

a. **Windows System Files:** Obtain legitimate PE files from your own Windows operating system or open-source software projects like Sysinternals.

b. **Software Repositories:** Platforms like GitHub and SourceForge provide executable files from various open-source projects.

c. **Safety Measures:** Verify the authenticity and integrity of the files to ensure they are not tampered with.

## 3.2   Step-by-Step Guide to Extract Features from PE Files Using Python

Feature extraction is essential for analyzing PE files. Below is a step-by-step guide to extract useful features:

### 3.2.1   Set Up the Environment:

a. Install Python and required libraries:

- import pandas as pd

- import numpy as np

- from sklearn.model_selection import train_test_split

- from sklearn.ensemble import RandomForestClassifier

- from sklearn.metrics import classification_report, accuracy_score

- from sklearn.pipeline import Pipeline

- from sklearn.preprocessing import StandardScaler

- import pefile

- import os

### 3.2.2  Write a Python Script to Extract Features:

```python
# Feature extraction from PE files

def extract_pe_features(file_path):

    """Extract basic features from a PE file."""

    try:

        pe = pefile.PE(file_path)

        features = {

            "NumberOfSections": len(pe.sections),

            "ImageBase": pe.OPTIONAL_HEADER.ImageBase,

            "EntryPoint": pe.OPTIONAL_HEADER.AddressOfEntryPoint,

            "SizeOfCode": pe.OPTIONAL_HEADER.SizeOfCode,

            "SizeOfInitializedData": pe.OPTIONAL_HEADER.SizeOfInitializedData,

            "SizeOfUninitializedData": pe.OPTIONAL_HEADER.SizeOfUninitializedData,

            "NumberOfImports": len(pe.DIRECTORY_ENTRY_IMPORT) if hasattr(pe,
"DIRECTORY_ENTRY_IMPORT") else 0,
```

```
        }

    return features

except Exception as e:

    print(f"Error processing {file_path}: {e}")

    return None
```

### 3.2.3  Validate the Extracted Features:

a.  Check the dataset for consistency and correctness.

b.  Visualize data distributions to identify anomalies or patterns.

## 3.3   Issues and Solutions

### 3.3.1  Small Dataset:

a.  **Problem:** A limited number of samples can lead to poor model generalization.

b.  **Solution:**

i.    Use data augmentation techniques, such as adding noise to features or using oversampling methods like SMOTE (Synthetic Minority Oversampling Technique).

ii.   Leverage open-source datasets and collaborating with cybersecurity organizations to access more data.

### 3.3.2  Unknown Extensions:

a.  **Problem:** Some files might not have the .exe or .dll extensions but are still valid PE files.

b.  **Solution:** Use file signature checks to identify PE files based on their internal structure rather than extensions.

### 3.3.3  Imbalance Between Normal and Malicious Files:

a.  **Problem:** Having more normal files than malware can bias the model.

b.  **Solution:**

i. Use under sampling for normal files or oversampling for malicious files to balance the dataset.

ii. Apply advanced techniques like weight loss functions in machine learning models to account for the imbalance.

Building a dataset of PE files and extracting features is a foundational step in malware analysis. By ensuring ethical and safe practices when handling malware, using Python scripts to extract meaningful features, and addressing issues like data imbalance, researchers can create effective models for detecting malicious software.

# 4  Model Training and Testing

Training a machine learning model to detect malware involves selecting an appropriate supervised learning algorithm, training it on extracted features, and evaluating its performance. Let's explore these aspects in detail.

## 4.1  Supervised Learning Algorithm

In our Python script we used the **Random Forest Classifier**, which is a supervised learning algorithm from the Decision Trees family.

**How Random Forest Works:**

- Random Forest is an ensemble learning algorithm that creates multiple decision trees during training. (Niklas Donges, Nov 26, 2024)

- Each tree makes a prediction, and the majority vote determines the final classification.

- This method helps reduce overfitting and improves accuracy.

We have used **RandomForestClassifier(n_estimators=100)**, which means it creates 100 decision trees and takes a vote to decide if a file is malware or normal.

**Why Random Forest?**

Random Forest is a method of group learning that uses several decision trees to produce predictions. It functions by dividing the data into decision trees, which are smaller sections that each concentrate on distinct patterns. Either the majority vote (classification) or average (regression), the outputs of all the trees are combined to provide the final forecast.
Even if the dataset contains noise or unimportant features, this approach performs well and is resistant to overfitting.

## 4.2   Accuracy and F1-Score

**Accuracy:** Basically Accuracy measures the percentage of correctly classified instances out of all predictions. (Huilgol, Aug 24, 2019)

It is calculated as

*Accuracy= (TruePositives+TrueNegatives)/ TotalSamples*

A high accuracy (close to 1.0 or 100%) means the model is correctly detecting malware and normal files most of the time.

where:

- **TP (True Positives)** = Malware correctly classified as malware

- **TN (True Negatives)** = Normal files correctly classified as normal

- **FP (False Positives)** = Normal files incorrectly classified as malware

- **FN (False Negatives)** = Malware incorrectly classified as normal

**F1-Score:** F1-score is a balance between **precision** and **recall**:

- **Precision:** Out of all files predicted as malware, how many were actually malware?

- **Recall:** Out of all actual malware files, how many were correctly detected?

F1-score helps when the dataset is **imbalanced** (e.g., more normal files than malware). It ensures the model isn't just biased towards the majority class.

**Accuracy & F1-Score in our Script**

our script evaluates the model using:

print(classification_report(y_test, y_pred))

print("Accuracy:", accuracy_score(y_test, y_pred))

After training, the model prints accuracy and an F1-score based on the test dataset. As accuracy is **too high (1.0 or 100%)**, it might indicate **overfitting**, so it means the model is memorizing rather than generalizing.

## 4.3 How the Model Classifies Files

Our script extracts features from PE (Portable Executable) files and uses them to classify the file.

**Steps the Model Follows to Classify a File:**

### 4.3.1 Feature Extraction

In machine learning and data analysis, feature extraction is the process of finding and removing pertinent features from unprocessed data. A more informative dataset is then produced using these attributes, and it can be used for several purposes, including: (Domino, 2025)

- Classification
- Prediction
- Clustering

- Extracts information like NumberOfSections, ImageBase, SizeOfCode, etc., from the file.

- These are stored as numerical values in a DataFrame.

### 4.3.2 Training the Model

- The dataset is split into **training (80%)** and **testing (20%)** using:

    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

- The RandomForestClassifier is trained on these extracted features.

### 4.3.3 Making Predictions

- When a new file is tested:

    prediction = model.predict(feature_df)

- The model uses patterns learned during training to decide if the file is **Malicious (1) or Normal (0)**.

### 4.3.4 Output the Result

return "Malicious" if prediction[0] == 1 else "Normal"

- If the model predicts 1, it labels the file as **Malware**.

- If the model predicts 0, it labels the file as **Normal**.

# 5 Evidence:

## 5.1 Complete python script in our report

### 5.1.1 Step 1: Feature extraction from PE Files

Import Libraries and **Feature extraction from PE Files**



*Figure 5.1.1: Feature Extraction from PE Files*

## 5.1.2  Step 2: Dataset Preprocessing:



*Figure 5.1.2: Dataset Preprocessing*

## 5.1.3 Step 3: Training and Evaluation Pipeline:



*Figure 5.1.3: Training and Evaluation Pipeline*

## 5.1.4 Step 4: Malware Detection Function:



*Figure 5.1.4: Malware Detection Function*

### 5.1.5  Step 5: Testing Normal and Malicious file:



*Figure 5.1.5: Testing Normal and Malicious file*

## 5.2   Malware Samples:

### 5.2.1  Malware Sample Files

We have collected 5 samples from the malware bazaar:

- Malware_1.exe

- Malware_2.exe

- Malware_3.exe

- Malware_4.exe

- Malware_5.exe

*Figure 5.2.1: Malware Samples*

## 5.3 Normal Samples:

### 5.3.1 Normal Sample Files

We have created 5 different programs in python language as normal files.

- add.exe
- division.exe
- input.exe
- multiply.exe
- sub.exe

*Figure 5.3.1: Normal Samples*

## 5.4   Malicious file as a malware

### 5.4.1  Output:

We have provided a malicious file (malware_6.exe) to our trained model for testing purpose, Here is output.

*Figure 5.4.1: Malware File Output*

## 5.5　Normal file as a Normal

### 5.5.1　Output:

We have provided a normal file (list.exe) to our trained model for testing purpose, Here is output.

*Figure 5.5.1: Normal File Output*

# 6   Model as a Cybersecurity Tool

Machine learning models, like the Random Forest-based malware detection system, can be an integral component of larger cybersecurity tools. By automating the analysis and classification of potential threats, such a model enhances the efficiency, accuracy, and adaptability of cybersecurity measures.

## 6.1   Integration into a Cybersecurity System

Here's how the model could fit into a comprehensive cybersecurity tool:

### 6.1.1   Malware Detection Module:

a. The model serves as the core engine for identifying malicious files.

b. It processes incoming files in real-time, extracting features like entropy, imports, and section sizes, then classifies them as either malware or normal.

### 6.1.2   Preprocessing and Feature Extraction:

a. A preprocessing module prepares the input data for the model by extracting relevant features from PE files using tools like pefile.

b. This ensures consistent and accurate feature representation, critical for the model's performance.

### 6.1.3   Threat Mitigation Pipeline:

a. When a file is flagged as malware:

   i. The system quarantines files to prevent execution.

   ii. It alerts the cybersecurity team or triggers automated responses, such as blocking the file on all connected systems.

b. For normal files, the system allows safe execution, minimizing disruptions to users.

### 6.1.4   Behavioral Monitoring:

a. The model could work alongside behavioral analysis tools that monitor runtime actions of executables.

b.  If the model misses a threat initially, unusual behavior during execution can trigger additional scans.

### 6.1.5  Integration with Threat Intelligence:

a.  The model can leverage external threat intelligence feeds to update its feature set, improving detection of new malware variants.

b.  For example, if a new malware family is detected globally, its features can be incorporated into the model's training pipeline.

## 6.2   Advantages in a Larger Cybersecurity Ecosystem

### 6.2.1  Real-Time Threat Detection:

a.  The model can analyze files as they are downloaded, opened, or transferred, providing immediate feedback on potential risks.

### 6.2.2  Scalability:

a.  The model can handle large volumes of data, making it suitable for enterprise-level deployments where thousands of files are analyzed daily.

### 6.2.3  Automation:

a.  Automating malware detection reduces the workload for cybersecurity teams, allowing them to focus on more complex threats.

### 6.2.4  Adaptability:

a.  As new malware emerges, retraining or fine-tuning the model ensures it remains effective.

## 6.3   Enhancements and Features for Practical Deployment

### 6.3.1  User-Friendly Dashboard:

a.  A graphical interface displays detection statistics, flagged files, and threat levels in real-time.

### 6.3.2  Integration with Endpoint Protection:

a.  The model can be embedded into endpoint protection tools like antivirus software or host intrusion detection systems (HIDS). (Redscan, 2025)

### 6.3.3  Regular Updates:

a.  A mechanism to regularly update the model with new malware data ensures it remains effective against evolving threats.

### 6.3.4  Cross-Platform Support:

a.  The model can be extended to analyze files on different platforms, such as macOS or Linux, for broader protection.

### 6.3.5  Incident Reporting:

a.  When malware is detected, the system generates detailed reports, including feature analysis and reasons for classification, aiding in forensic investigations.

## 6.4   Example Workflow in a Cybersecurity Tool

### 6.4.1  File Upload or Detection:

a.  A user downloads or opens a file.

b.  The file is intercepted by the cybersecurity tool and sent to the model for analysis.

### 6.4.2  Feature Extraction:

a.  Features like section count, entropy, imports, and exports are extracted and passed to the model.

### 6.4.3  Classification:

a.  The model predicts whether the file is malware or normal.

### 6.4.4  Response:

a.  If malware: The file is quarantined, and a system-wide alert is issued.

b.  If normal: The file is allowed to be executed.

### 6.4.5 Post-Incident Learning:

a. Feedback from the cybersecurity team or additional behavioral analysis helps refine the model.

Integrating the malware detection model into a larger cybersecurity tool creates a robust, automated defense mechanism capable of adapting to new threats. By combining real-time analysis, automation, and continuous learning, this model can significantly enhance the overall security posture of individuals and organizations.

# 7 Conclusion

## 7.1 Summary of Key Points

In this report, we explored how Artificial Intelligence (AI) and machine learning can enhance cybersecurity, specifically in detecting malware. Below are the key takeaways:

1. **AI in Cybersecurity:** AI helps in identifying malware by analyzing large datasets, recognizing patterns, and detecting anomalies that traditional methods might miss. It is particularly effective in identifying zero-day threats and polymorphic malware.

2. **Understanding PE File Structure:** The Portable Executable (PE) file structure provides critical information about executable files. Features such as section count, section sizes, imports, and entropy offer valuable insights for malware detection.

3. **Dataset Preparation:** Safe and ethical practices are essential when collecting malware samples and legitimate PE files. Extracting meaningful features using Python scripts ensures that the dataset is suitable for machine learning applications.

4. **Model Training and Testing:** The Random Forest algorithm was chosen for its ability to handle complex datasets and prevent overfitting. Metrics like accuracy and F1-score were used to evaluate the model's performance, ensuring reliable and balanced predictions.

5. **Model Integration:** The model can be integrated into a larger cybersecurity tool, providing real-time malware detection, scalability, automation, and adaptability to evolving threats.

## 7.2 Reflection and Real-World Applications

This project provided valuable insights into the intersection of AI and cybersecurity, emphasizing the role of machine learning in automating and enhancing threat detection. The process of collecting data, extracting features, and training a model highlighted the importance of meticulous preparation and ethical considerations in cybersecurity research.

Key lessons learned include:

1. **Feature Importance:** Understanding the significance of PE file attributes helped in designing an effective detection model. This knowledge can be applied to analyze other file formats and extend protection to different platforms.

2. **Algorithm Selection:** Choosing the right machine learning algorithm and evaluation metrics ensures robust and reliable solutions. This approach can be generalized to solve other cybersecurity problems, such as phishing detection or network intrusion prevention.

3. **Handling Imbalanced Data:** Techniques to address class imbalance, such as oversampling and class weighting, are crucial in real-world datasets where malicious files are relatively rare.

4. **Automation and Scalability:** The integration of machine learning models into larger tools demonstrates how AI can streamline cybersecurity operations, reducing the manual workload for analysts while improving detection rates.

This knowledge has practical applications in real-world scenarios, such as:

- **Enterprise Security:** Organizations can deploy similar models to scan incoming files, emails, or network traffic for potential threats.

- **Incident Response:** Security teams can use the model's outputs to prioritize and respond to high-risk incidents.

- **Threat Intelligence:** Continuous updates to the model using global threat intelligence can ensure resilience against emerging malware.

This report underscores the transformative potential of AI in cybersecurity. By leveraging machine learning models, we can build proactive and adaptive systems that keep pace with the ever-evolving threat landscape. The knowledge gained from this project can contribute to designing innovative solutions that protect individuals and organizations from sophisticated cyberattacks.

# 8 References

- Breiman, L. (2001) 'Random forests', *Machine Learning*, 45(1), pp. 5–32. Available at: https://link.springer.com/article/10.1023/A:1010933404324 (Accessed: 12 February 2025).

- Domino, 2025. *Feature Extraction*. [Online]
  Available at: https://domino.ai/data-science-dictionary/feature-extraction
  [Accessed 11 02 2025].

- Huilgol, P., Aug 24, 2019. *Accuracy vs. F1-Score*. [Online]
  Available at: https://medium.com/analytics-vidhya/accuracy-vs-f1-score-6258237beca2
  [Accessed 09 02 2025].

- Iryna Sydorenko, December 21, 2023. *What Is a Dataset in Machine Learning:*.
  [Online]
  Available at: https://labelyourdata.com/articles/what-is-dataset-in-machine-learning
  [Accessed 11 02 2025].

- Microsoft, 02/29/2024. *PE Format*. [Online]
  Available at: https://learn.microsoft.com/en-us/windows/win32/debug/pe-format
  [Accessed 08 02 2025].

- Niklas Donges, Nov 26, 2024. *Random Forest: A Complete Guide for Machine Learning*. [Online]
  Available at: https://builtin.com/data-science/random-forest-algorithm
  [Accessed 09 02 2025].

- Priyanka Jadav, April 27, 2024. *Malware Detection Using Machine Learning Techniques*. [Online]
  Available at: https://www.einfochips.com/blog/malware-detection-using-machine-learning-techniques/
  [Accessed 11 02 2025].

- Redscan, 2025. *HIDS - Host-based intrusion detection*. [Online]
  Available at: https://www.redscan.com/services/hids/
  [Accessed 12 02 2025].