

# Building a Static Malware Detector: Comprehensive Report

## Introduction:

The proliferation of malware poses significant threats to digital security, making malware detection an essential component of cybersecurity measures. In this report, we present the process of building a static malware detector using machine learning techniques. The detector aims to differentiate between benign software and malicious code based on static features extracted from executable files.

## Background:

Static malware detection involves analyzing the characteristics of executable files without executing them. This approach is particularly valuable for detecting unknown malware variants and has become increasingly important alongside traditional signature-based methods.

## Data Collection and Preprocessing:

The dataset used for training and testing the malware detector comprises samples collected from directories labeled as benign and malware. Each sample is represented as an executable file, and labels are assigned based on the source directory. To ensure model generalization, the dataset is split into training and testing sets, maintaining class balance.

## Feature Extraction:

Feature extraction plays a crucial role in capturing the distinguishing characteristics of malware and benign software. In our approach, we employ two main types of features:

- N-gram Features:** Byte sequences of the samples are converted into N-grams, capturing patterns within the binary code. The frequency of each N-gram is computed to construct feature vectors.

- Metadata Features:** Additional features such as DLL imports and PE sections are extracted from the samples, providing contextual information about the executable files.

### **Model Training:**

A Random Forest classifier is chosen for its ability to handle high-dimensional data and robustness against overfitting. The classifier is trained using the feature vectors generated from the training dataset. However, it's important to note that the high training accuracy observed may indicate potential overfitting, which requires further investigation.

### **Model Evaluation:**

The trained model is evaluated on the testing dataset to assess its performance in classifying malware and benign samples. Despite encountering errors during feature extraction, the model achieves a high accuracy on the test set, demonstrating its effectiveness in distinguishing between the two classes.

### **Diagnosis of Results:**

The exceptional accuracies observed in both training and testing phases can be attributed to various factors, including the balanced dataset, the robustness of the random forest classifier, and the quality of feature extraction. However, there is a potential risk of overfitting, which necessitates further validation and fine-tuning of the model.

### **Deployment:**

The trained model is deployed using Amazon SageMaker, allowing for real-time predictions on new samples. The model endpoint enables seamless integration into existing cybersecurity systems, enhancing detection capabilities against evolving malware threats.

### **Conclusion:**

In conclusion, the static malware detector demonstrates promising results in accurately identifying malicious code and benign software based on static features. However, ongoing refinement and validation are essential to address potential overfitting and enhance the detector's performance in real-world scenarios. By leveraging machine learning techniques, we

contribute to the advancement of cybersecurity measures, safeguarding digital assets against malicious attacks.