

Malware Classification Using Deep Learning

Statement

In today's digital landscape, malware, an abbreviation for malicious software, represents a pervasive threat. Its primary goal is to disrupt, damage, or gain unauthorized access to computer systems. Given the ever-evolving nature of cyber threats, the risk of sensitive information theft or manipulation is ever-present, highlighting the urgent need for robust defense mechanisms.

Antivirus and anti-malware software offer essential protection against malware attacks. However, to stay ahead of the curve, it's imperative to maintain an updated repository of past malware instances. This approach enables the prediction of future attack patterns and ensures a swift and effective response to emerging threats, safeguarding critical digital assets.

Dataset

Maling Dataset

The Maling dataset is a public benchmark dataset commonly used for training and evaluating machine learning models for malware classification. The Maling Dataset contains 9339 malware images, belonging to 25 families/classes.

Content:

The dataset consists of images representing malware byteplots. These byteplots are visual representations of the byte sequences that make up the malware code.

Each image in the dataset corresponds to a specific malware sample belonging to a particular malware family. Thus, our goal is to perform a multi-class classification of malware.

Here is the information regarding the dataset :

	Family/Class	Type
0	Adialer.C	Dialer
1	Agent.FYI	Backdoor
2	Allapple.A	Worm
3	Allapple.L	Worm
4	Alueron.gen!J	Worm
5	Autorun.K	Worm:AutoIT
6	C2LOP.P	Trojan
7	C2LOP.gen!g	Trojan
8	Dialplatform.B	Dialer
9	Dontovo.A	Trojan Downloader
10	Fakerean	Rogue
11	Instantaccess	Dialer
12	Lolyda.AA1	PWS
13	Lolyda.AA2	PWS
14	Lolyda.AA3	PWS
15	Lolyda.AT	PWS
16	Malex.gen!J	Trojan
17	Obfuscator.AD	Trojan Downloader
18	Rbot!gen	Backdoor
19	Skintrim.N	Trojan
20	Swizzor.gen!E	Trojan Downloader
21	Swizzor.gen!I	Trojan Downloader
22	VB.AT	Worm
23	Wintrim.BX	Trojan Downloader
24	Yuner.A	Worm

Maling Dataset

The Maling dataset is a public benchmark dataset commonly used for training and evaluating machine learning models for malware classification.

Content:

The dataset consists of images representing malware byteplots. These byteplots are visual representations of the byte sequences that make up the malware code.

Each image in the dataset corresponds to a specific malware sample belonging to a particular malware family.

ApproachProblem

1. **Feature Extraction:** Convolutional Neural Networks (CNNs) excel at extracting features from spatial data like images. Byteplot images, which represent the malware code visually, contain spatial information. The convolutional layers in a 2D CNN can automatically learn relevant features from these images, such as patterns, edges, and textures, that differentiate between different malware families.
2. **Image Recognition Expertise:** CNNs have been highly successful in various image recognition tasks. They can leverage this expertise to analyze byteplot images effectively. The existing knowledge gained from image recognition tasks can be transferred to this domain through techniques like transfer learning, further improving performance.
3. **Efficiency:** 2D CNNs are computationally efficient for processing 2D data like images. This allows for faster training and evaluation compared to other deep learning architectures that might be less suitable for image data.

Model :

2D CNN

A CNN is used for analyzing images like byteplots of malware code.

Layers in CNN : A CNN is like a layered stack of filters. Each layer learns to detect specific features in the image.

Convolution: The first layers apply small filters that slide across the image, detecting basic shapes, edges, and textures.

Pooling: These layers summarize the information from the previous layer, reducing the image size while keeping important details.

Deeper Layers: As you go deeper, the filters become more complex, combining the simpler features into more intricate patterns specific to malware families.

Why is a 2D CNN good for Malware Classification?

Byteplots: Malware byteplots are essentially 2D images, where each pixel represents a value in the malware code.

Feature Extraction: CNNs excel at automatically learning features from images. These features could be specific patterns or textures that differentiate between different malware families.

Image Recognition Expertise: CNNs have been very successful in recognizing objects and patterns in images. This knowledge can be transferred to analyzing byteplots for malware classification.

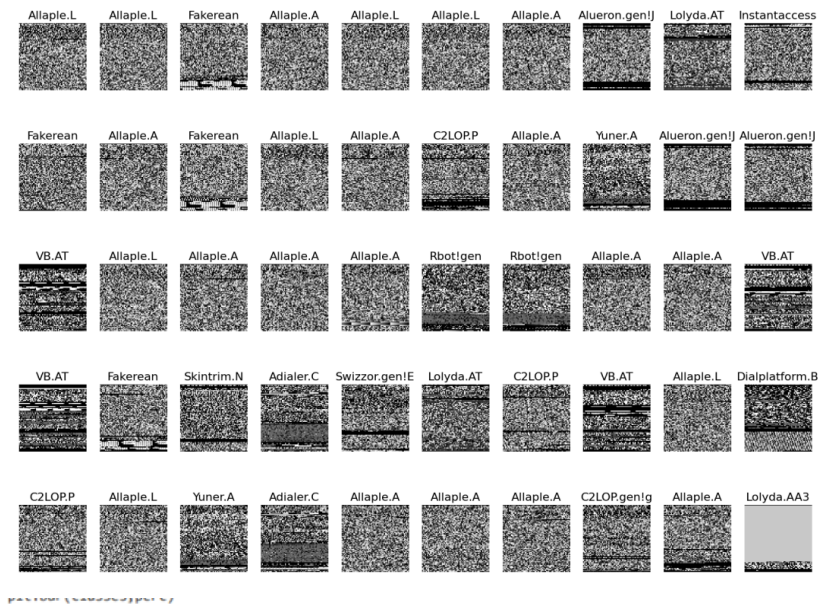
Efficiency: 2D CNNs are computationally efficient for processing 2D data like images, allowing for faster training and evaluation compared to other deep learning architectures.

1) Data Management

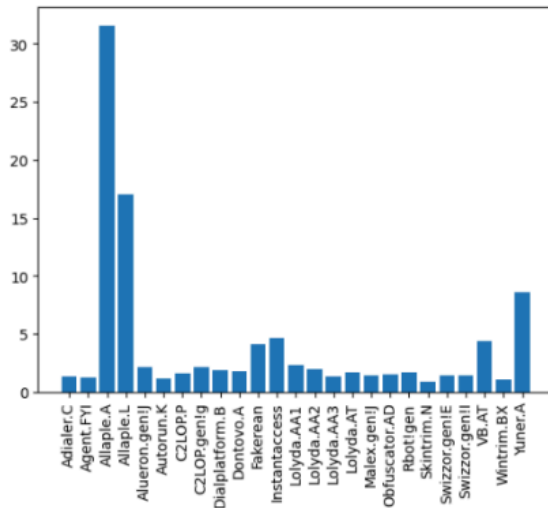
Python libraries like NumPy, Pandas, and Matplotlib can be used for data manipulation, analysis, and visualization.

Deep learning frameworks like TensorFlow or PyTorch often have built-in functionalities for data loading and preprocessing.

2) Visualize



<BarContainer object of 25 artists>

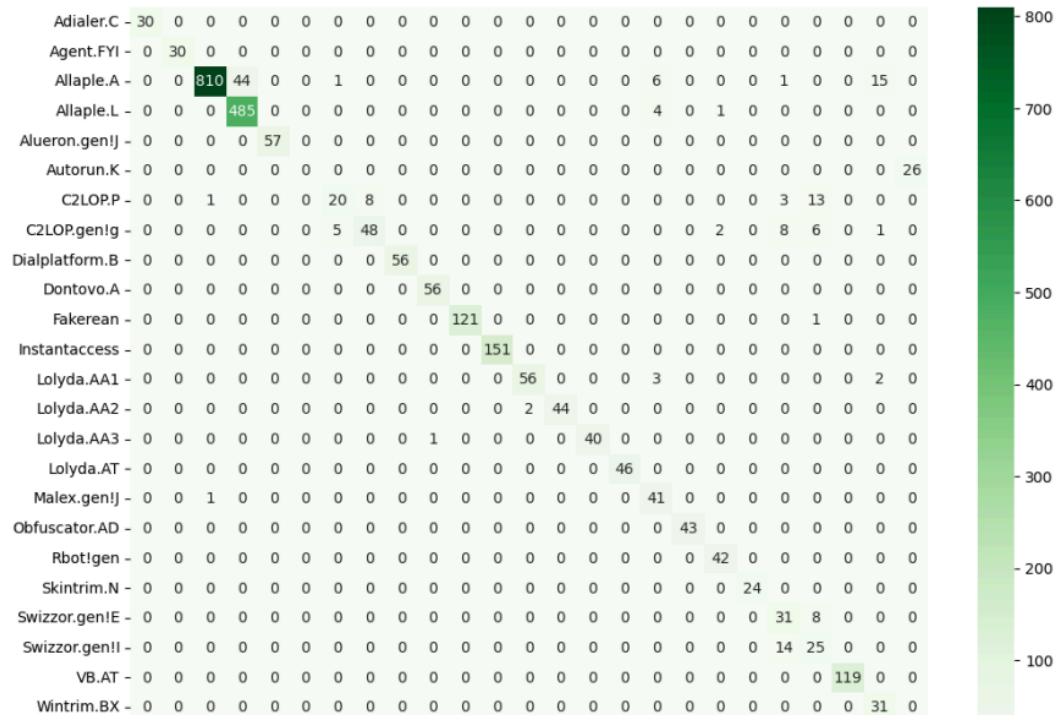


```

import seaborn as sns
import pandas as pd
def confusion_matrix(confusion_matrix, class_names, figsize = (12,9), fontsize=16):

    df_cm = pd.DataFrame(confusion_matrix, index=class_names, columns=class_names)
    fig = plt.figure(figsize=figsize)
    .....
    class_names= batches.class_indices.keys()
    confusion_matrix(c_matrix, class_names, figsize = (12,9), fontsize=16)

```



3) Computing Weights

4) Splitting of Data

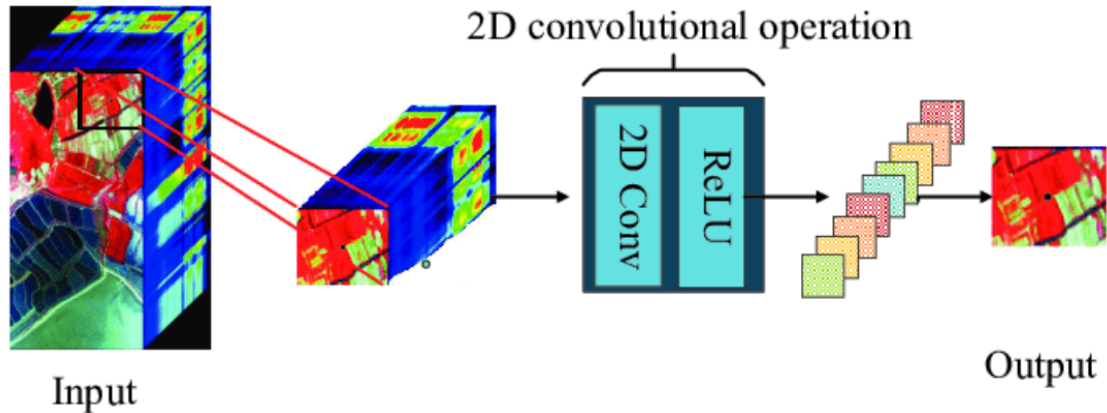
Divide your dataset into training, validation, and test sets.

Training Set (70-80%): Used to train the CNN model. The model learns to identify features and patterns in the data.

Validation Set (10-15%): Used to monitor the model's performance during training and prevent overfitting. Overfitting happens when the model memorizes the training data too well and performs poorly on unseen data. The validation set helps you adjust hyperparameters (learning rate, number of layers) to achieve optimal performance.

Test Set (10-15%): Used to evaluate the final performance of the trained model on unseen data. This provides an unbiased estimate of how well the model will generalize to real-world malware samples.

5) CNN Architecture



The 2D CNN takes the byteplot image of the malware code as input. It then processes the image through its layers, automatically learning features that distinguish between different malware families. Finally, the model outputs a prediction for which malware family the code belongs to.

6) Accuracy:

Accuracy is calculated as the number of correct predictions divided by the total number of test samples.

The final Accuracy that we got after using 2D CNN model is

Final CNN accuracy: 0.9368308186531067

7) Confusion Matrix

Future Scope

1. **Beyond Images:** Deep learning will move beyond analyzing just image data (like byteplots). It will incorporate other sources like network traffic patterns or data from emulated malware execution for a more comprehensive understanding.
2. **Explainable AI:** There will be a focus on making deep learning models more transparent. Techniques like Explainable AI (XAI) will help us understand why the model classifies something as malware, increasing trust and identifying potential biases.
3. **Integration and Automation:** Deep learning models will be integrated with security tools like sandboxes and Security Information and Event Management (SIEM) systems. This will enable real-time analysis of suspicious files, network traffic, and system events to identify and respond to malware threats more effectively.
4. **Zero-Day Malware:** A significant challenge is detecting previously unseen malware (zero-day attacks). Future models will explore techniques like anomaly

detection or one-class classification to identify suspicious outliers that might be new malware variants.

5. **Finer-grained Classification:** Instead of just classifying malware families, the future holds promise for pinpointing specific functionalities (ransomware, spyware) or even specific malware strains. This will require more complex models and well-annotated datasets.

Conclusion

This project investigated using a 2D Convolutional Neural Network (CNN) to classify malware families. The model analyzes byteplot images, a visual representation of the malware's code, to identify different malware types. This approach leverages deep learning's ability to automatically learn features from the images, potentially even detecting new and unseen malware variants. In essence, the project explores using deep learning to create a "malware identifier" based on how the malware code looks visually. This has the potential to be a powerful tool in the fight against ever-evolving cyber threats.