

# アンサンブル学習を用いたマルウェア判定システムの検討

小寺 建輝<sup>†</sup> 房安 良和<sup>†</sup> 泉 隆<sup>†</sup>

<sup>†</sup> 日本大学

## 1. はじめに

近年、マルウェアの亜種生成が高速化され、ウイルス定義ファイルの作成・配布が追いつかない現状となっている。このような亜種検知に関する問題を解決するため、機械学習により亜種を検知・分類する研究が現在取り組まれている。その中でも、マルウェアを画像化し、画像特徴量をもとに亜種を該当するファミリーに分類する先行研究[1]では、高い識別精度でマルウェアを分類できたことが報告されている。これは、亜種が元のコードの一部のみを改変して作成されるため、元のマルウェアとその亜種、つまり同一ファミリーでは視覚的に類似した画像(テキスト画像)が得られるためである。しかし、先行研究では、マルウェアと正常ファイルを識別することが検討されていない。そこで本研究では、画像特徴量の類似したマルウェアの画像をクラスタリングによりグループ化してグループごとにマルウェア検知モデルを構築する。そして、グループごとに構築した複数のマルウェア検知モデルを組み合わせたアンサンブル学習により、マルウェアと正常ファイルを識別するマルウェア判定システムを検討する。

## 2. マルウェア画像化

マルウェアを画像化する手法[1]を以下に示す。また実際にマルウェアを画像化した例を図1に示す。

- (1) 対象ファイルを1Byte(8bit)ずつ読み込み1次元配列に格納する
- (2) ファイルサイズ(配列の要素数)に応じて幅を決定し、2次元配列に変換する
- (3) 配列の要素の値は8bitであり、0-255の範囲であるため、その値を画素値として256階調のグレースケール画像を生成する

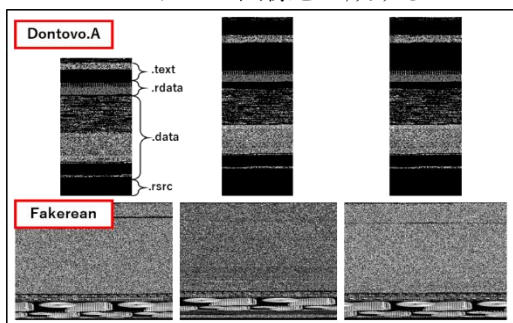


図1.マルウェア画像化の例

本研究では、これらのマルウェアの画像から、テキスト画像認識に有効な Gist 特徴量, LBP 特徴量, HLAC 特徴量を抽出し、学習や判定に利用する。

## 3. マルウェア判定システム

複数のマルウェア検知モデルを組み合わせたアンサンブル学習を用いて、マルウェアと正常ファイルを識別するマルウェア判定システムを以下の5つのフェーズに分けて説明する。

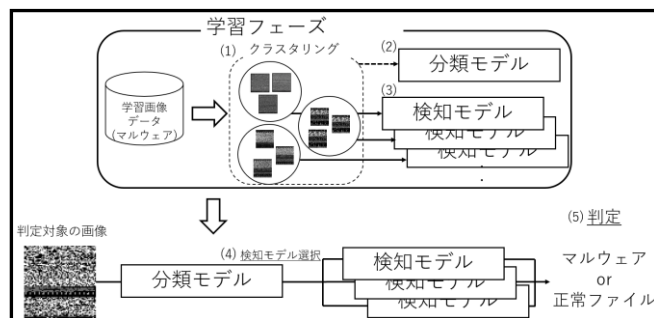


図2.マルウェア判定システム

- (1) 類似した画像のグループ化  
学習データであるマルウェアをクラスタリングによりグループ化し、各グループにグループIDを割り当てる。クラスタリング手法にはクラスタ数を自動推定可能なX-means[2]を用いる。
- (2) 分類モデルの構築  
グループIDを教師ラベルにし、SVMにより分類モデルを構築する。
- (3) 検知モデルの構築  
判定対象が(1)で作成したグループに属しマルウェアか否かを判定する、検知モデルを各グループで構築する。また、検知モデルの構築には異常検知アルゴリズムであるIsolation Forest[3]を用いる。
- (4) 判定対象の分類・検知モデルの選択  
判定対象を分類モデルに入力する。このとき、分類モデルが出力するクラスごとのスコアからロジスティック回帰式によりクラスごとの所属確率を求める。そして、所属確率が任意の閾値を超えたクラス(グループ)の検知モデルを選択する。
- (5) マルウェアの判定  
選択された検知モデルに対し、判定対象を入力し、各検知モデルで出力されたスコアをもとに、マルウェアか否かの判定を行う。このとき、全ての検知モデルにおいてマルウェアでないと判定された画像を正常ファイルと識別する。

## 4. まとめ

本研究では、グループごとの複数のマルウェア検知モデルを組み合わせたアンサンブル学習を用いて、マルウェアと正常ファイルを識別するマルウェア判定システムについて検討した。

今後は、実際にマルウェア判定システムを実装し、検知率や誤検知率等の精度についての検証を行う。

## 参考文献

- [1] L. Nataraj, et al.: "Malware Images: Visualization and Automatic Classification", VizSec'11(2011-07)
- [2] Dan Pelleg, Andrew W. Moore: "X-means: Extending K-means with Efficient Estimation of the Number of Cluster", ICML, pp.727-734(2000)
- [3] Fei Tony Liu, et al.: "Isolation-Based Anomaly Detection", ACM Transactions on Knowledge Discovery from Data (TKDD), Vol.6, No.1, pp.1-39(2013-03)