

データ圧縮アルゴリズムによる マルウェア分類および亜種判定について

武智 聡平¹ 甲斐 博¹ 森井 昌克²

概要: マルウェアを詳細な解析をすることなく、できるだけ短時間で、かつ自動的に分類することは重要な課題となっている。従来から、その動的解析より観測される API コールの特徴によって分類が行われ、その精度が競われている。本研究では API コールの時系列発生順序に着目し、その特徴によって分類を試みる。独創的な点はその特徴の計算のためにユニバーサル圧縮手法として知られる LZ77 アルゴリズムを用いる点である。LZ77 はスライド辞書法の一つであり、スライド辞書のサイズを変更することにより、API コールの出現特性に応じた圧縮率となる。この特性を用いて、種々の辞書サイズによる圧縮率によって分類を試みる。分類するための方法としては SVM を用い、分類の精度を求め、本提案手法の有効性を示す。

キーワード: マルウェア解析, 分類, 亜種, データ圧縮アルゴリズム, SVM

Malware Classification Based on Data Compression Algorithm

SOHEI TAKECHI¹ HIROSHI KAI¹ MASAKATU MORII²

Abstract: The method to classify automatically the malwares as quickly as possible without their detailed analysis has become an important research subject. Traditionally, the malware classification has been performed by the characteristics of the APIs observed from the dynamic analysis, and the researchers compete the accuracy. In this research, we focus on the time order of the API call, and try to classify the malware variant by its characteristics. The originality of our research is to use the LZ77 algorithm, known as an universal data compression algorithm, to derive the characteristics. The LZ77 is a slide dictionary method, and its data compression ratio is depend on the size of the slide dictionary and the time order of the API call. We classify the malware variants by the characteristics derived from the data compression ratios, which are computed by the different kind of sizes of the slide dictionary. We perform the malware variant classification based on the SVM algorithm, and examine the accuracy of our method. Experimental results show that the proposed method is effective.

Keywords: Malware Analysis, Classification, Malware Variant, Data Compression Algorithm, SVM

1. はじめに

インターネットの利用者にとってマルウェアの増加は問題になっており、AV-TEST のレポートによると、新種のマルウェアは近年では年間に 1 億種類以上発見されてい

る [1].

マルウェア対策のためのアンチウィルスソフトでは「ウイルス定義データベース」を利用してマルウェアを検出する。ウィルス定義データベースによる検出から逃れるため多くの亜種が作成されており、マルウェアの発見数が膨大になる要因の 1 つである。

マルウェアの亜種が増加する中で、マルウェアを詳細な解析をすることなく、できるだけ短時間で、自動的に分類することは重要な研究課題の一つである。従来から、動的

¹ 愛媛大学大学院理工学研究科
Graduate School of Science and Engineering, Ehime University

² 神戸大学大学院工学研究科
Graduate School of Engineering, Kobe University

解析より観測される API コールの特徴によって分類することが行われ、その精度が競われている。例えば、多次元ベクトルでマルウェアの挙動を表現しハミング距離で亜種判定を行う研究 [10] や、API コール列を Paragraph Vector を用いて表現し亜種判定を行う手法 [5] などが提案されている。

本稿では、API コールの時系列発生順序に着目し、その特徴によって分類を試みる。独創的な点はその特徴の計算のためにユニバーサル圧縮手法として知られる LZ77 アルゴリズムを用いる点である。LZ77 はスライド辞書法の一つであり、スライド辞書のサイズを変更することにより、API コールの出現特性に応じた圧縮率となる。この特性を用いて、種々の辞書サイズによる圧縮率によって分類を試みる。分類するための方法としては SVM を使い、分類の精度を求め、本提案手法の有効性を示す。FFRI dataset を用いた亜種判定の実験を行い、平均として約 82% で亜種分類が成功することを示す。

2. マルウェアの動的解析と FFRI dataset

マルウェア解析手法としては大きく分類して表層解析、静的解析、動的解析がある。表層解析はファイルの情報などの収集を目的として行われ、静的解析では逆アセンブラなどによりバイナリファイルを用いてマルウェアの特徴を解析する。一方で、動的解析は検体を実際に実行し、プロセスの起動、ファイルの書き換えやダウンロード、通信パケットの記録などの不正な動作を記録し、マルウェアの振る舞いを解析するものである。

動的解析はその動作を検証するために監視ツールを用いて記録を行い、隔離された環境下で実施する必要がある。そのため動的解析のためのシステムが存在し、オープンソースマルウェア解析システムである Cuckoo Sandbox などが知られている。

また、研究用データセットとして配布されている MWS Dataset 2016 には、マルウェアの動的解析結果の FFRI Dataset 2013~2016 が含まれている。このデータセットは株式会社 FFRI が独自に収集したマルウェアの動的解析ログである。

FFRI Dataset 2013 は 2641 検体、FFRI Dataset 2014 および 2015 はそれぞれ 3000 検体、FFRI Dataset 2016 は 8243 検体に対する動的解析のログファイルが存在している。FFRI Dataset の動的解析は Cuckoo Sandbox を用いて行われている [2]。また、FFRI Dataset 2014 では Cuckoo Sandbox での解析に加えて FFRI yarai analyzer Professional の解析結果ログも含まれている [3]。

動的解析結果には virustotal, behavior, network など 13 種類の項目に分類された結果が記録されている。

例えば virustotal には、VirusTotal の検査履歴との照合結果が示されており、Kaspersky, ESET などアンチウィル

```
LdrLoadDll LdrGetDllHandle DeviceIoControl NtAllocateVirtualMemory LdrLoadDll IsDebuggerPresent LdrLoadDll LdrLoadDll LdrLoadDll LdrLoadDll NtCreateMutant NtOpenSection SetWindowsHookExA SetWindowsHookExA CreateThread NtAllocateVirtualMemory NtAllocateVirtualMemory NtAllocateVirtualMemory NtAllocateVirtualMemory NtOpenSection NtClose SetWindowsHookExA SetWindowsHookExA RegOpenKeyExW RegQueryValueExW RegCloseKey RegCloseKey LdrGetDllHandle ...
```

図 1 “Trojan.Win32.Agent.rxht” の API コール列
Fig. 1 The API call sequence of “Trojan.Win32.Agent.rxht”

スベンダーが分類した亜種名が記載されている。

また behavior には、実行された API 名や時刻などが記載されており、実行された時系列順に並んでいる。例えば、“Trojan.Win32.Agent.rxht” について、behavior の中に現れる API コール列を取得した例を図 1 に示す。

ここで “Trojan.Win32.Agent.rxht” は Kaspersky によるマルウェア分類名である。Kaspersky の命名規則は次の通りである。

[Prefix:]Behaviour.Platform.Name[.Variant]

ここで、Prefix は対象を検知したサブシステム、Behaviour は検知された対象、Platform はプログラムコードが実行される環境、Name は検知された対象のファミリー名、Variant は亜種を指す [4]。

3. マルウェアの動的解析結果とそれに基づく分類

API コール列を利用したマルウェアの動的解析結果に基づくマルウェア分類では、呼び出された API 名の数による比較や、呼び出された API 名の出現パターンによる比較でマルウェアが分類される。

例として “Backdoor.Win32.Androm.bknn” と “Backdoor.Win32.Androm.affh” の API コール列を図 2 と図 3 に示す。これらの同じマルウェアファミリーに属する亜種では、“LdrGetProcedureAddress” と呼ばれる API がどちらも 17 回出現している。また、1 番目に呼び出された API の “NtOpenDirectoryObject” から、22 番目に呼び出された API の “LdrGetProcedureAddress” までの API の出現するパターンが全く同じであるという特徴を持つ。

一方で、全く別のマルウェアファミリーの亜種 “Trojan-Spy.Win32.Zbot.hpyq” の API コール列の一部を図 4 に示す。このマルウェアには “LdrGetProcedureAddress” と呼ばれる API の出現回数は 407 回出現しており、先ほど述べたマルウェアファミリー “Backdoor.Win32.Androm” における API の出現回数とは全く異なる。マルウェアファミリー “Backdoor.Win32.Androm” において観測できた API コール中の特徴的なパターンはマルウェアファミリー “Trojan-Spy.Win32.Zbot.hpyq” では見られない。

以上の例に見られるように、マルウェアファミリーが異

スライド窓と呼ばれるバッファを用意し、このバッファを辞書として利用することで符号化を行う。符号化における辞書の参照は、現在符号化の対象となっている位置から始まる記号列と一致する（スライド窓中の記号列で最も長い）ものを探すことを基本としている。スライド窓の長さは固定長であり、アルゴリズム実行時に決定する。

圧縮の対象とする入力データが1文字単位の語で構成される場合を考えてアルゴリズムの説明を行う。まずはアルゴリズムの中で用いる語句を説明する。

入力ストリーム データ圧縮の対象となる入力データ列

符号化位置 符号化を行う入力ストリーム上の位置（先読みバッファの先頭）

先読みバッファ 入力ストリームから記号を読み込む固定長のバッファ

辞書 先読みバッファで処理されたあとの入力ストリームが書き込まれる固定長のバッファ（初期状態は空である。辞書と先読みバッファをあわせてスライド窓と呼ばれる。）

ポインタ 先読みバッファの先頭から始まる文字列と最長一致するスライド窓中の文字列の先頭の位置とその文字列の長さ

このとき LZ77 のアルゴリズムの手順は以下のようにまとめられる。

- (1) 入力ストリームの開始位置に符号化位置を設定する。
- (2) スライド窓の中で先読みバッファと最長一致する文字列を見つける。
- (3) そのような文字列が見つければ、ポインタを出力し、符号化位置を最長一致した文字列の長さ分だけシフトさせて先読みバッファの最初の1文字を出力する。もし見つからなければ、空ポインタを出力し1文字分だけ符号化位置をシフトして先読みバッファの最初の1文字を出力する。
- (4) 先読みバッファが空でなければステップ(2)に戻る。つまり処理の手順としては、スライド窓に入力ストリームから文字列を読み込んで、辞書を参照しながら先読みバッファの中で繰り返されるパターンを見つけて符号化を行う処理をする。

ここで、動的解析結果より抽出した API コール列を保存したファイルを A とし、スライド窓サイズ s の LZ77 により圧縮したファイルを $f_s(A)$ とする。またファイル A のバイト数を $B(A)$ とすると、スライド窓サイズ s に関する圧縮率 $\rho_s(A)$ を次のように定義する。

$$\rho_s(A) = \frac{B(f_s(A))}{B(A)} \times 100.$$

この時、API コール列 A に対して n 種類のスライド窓のサイズを $s_i, 1 \leq i \leq n$ とした時、本稿で用いる n 次元特徴ベクトル $V(A)$ を

$$V(A) = (\rho_{s_1}(A), \rho_{s_2}(A), \dots, \rho_{s_n}(A))$$

と定義する。

本研究では、特徴ベクトル $V(A)$ の計算を次のような手順で求める。

- (1) FFRI dataset からの API コール列の取得
- (2) API コール列を2バイトコード列に変換
- (3) 特徴ベクトルの作成

API コール列は図1に示したように API が呼び出された順に API の名前を並べることにより得られる。

API 名のままではパターンを利用した圧縮を行うことができない。そのため API の名前に対して一意に対応がつく ID (整数値) を割り当てる。FFRI dataset 2013~2016 では現れる API の種類数が 256 以上あるため 2 バイトで整数値を表現しそれを 1 文字とする。

以上のように作成したファイルに対して LZ77 を用いた圧縮を行う。例えば、 $n = 4$ の場合の亜種名 “Worm.Win32.WBNA.hae” の特徴ベクトルを示す。

(77.95, 77.51, 76.86, 74.45)

ここで左からスライド窓のサイズが 2B, 3B, 4B, 5B とした時のデータである (B はバイトを表す)。計算には Python を用いてデータの表示に str 関数を用いたため圧縮率は 12 桁の浮動小数点数で得られるが、ここでは 4 桁で丸めて表示している。同じ亜種の “Worm.Win32.WBNA.bul” では、

(77.92, 76.80, 76.13, 74.32)

となり “Worm.Win32.WBNA.hae” と非常に近い特徴ベクトルが得られていることがわかる。また、マルウェアのファミリーが異なる例として、“Trojan.Win32.Llac.ddxv” に対する特徴ベクトルを示す。

(80.15, 63.74, 63.74, 63.74)

このように “Worm.Win32.WBNA” のファミリーの特徴ベクトルとは大きく異なる結果が得られ、本研究で提案する圧縮率を用いた特徴ベクトルによりマルウェア亜種の特徴が表現できていることを実験的に確認することができる。

5. SVM と特徴ベクトルを用いたマルウェア亜種分類の方法

本研究で用いる亜種分類手法は佐藤らの手法 [5] と同じであるが、Paragraph Vector の代わりに前節で定義した特徴ベクトルを用いる点が異なる。すなわち、入力とする特徴ベクトルが、あるマルウェアファミリーに属するかそうでないかを分類することを目的とし、学習用特徴ベクトルを SVM に入力しマルウェアの特徴を学習させる。

SVM の分類ラベルとして、特徴ベクトルがマルウェア

91.2510936133	89.4138232721	88.3639545057	86.0017497813
83.3770778653	80.927384077	77.4278215223	70.5161854768
69.3788276465	69.1163604549	68.7664041995	66.0542432196
65.0043744532	65.0043744532	64.1294838145	63.9545056868
63.3420822397	62.5546806649	62.2922134733	61.854768154
61.7672790901	61.154855643	60.3674540682	59.4925634296
59.3175853018	59.1426071741	58.0927384077	57.8302712161
57.4803149606	57.1303587052	91.2510936133	89.4138232721
88.188976378	85.5643044619	83.1146106737	81.8022747157
76.8153980752	70.2537182852	69.4663167104	69.4663167104
68.5914260717	65.8792650919	64.47944007	64.1294838145
63.2545931759	63.0796150481	62.9921259843	62.5546806649
62.3797025372	61.7672790901	61.5923009624	60.8923884514
60.542432196	59.6675415573	59.4925634296	59.230096238
57.7427821522	57.7427821522	57.217847769	56.605424322

図 5 “Trojan.Win32.Agent.rxht” の特徴ベクトル

Fig. 5 The characteristics vector of “Trojan.Win32.Agent.rxht”

ファミリーに属する場合は 1、そうでない場合は-1 を、特徴ベクトルの先頭に付加し、各マルウェアファミリー毎に学習用特徴ベクトルを作成し、分類器を作成する。

次節で述べる実験では、スライド窓のサイズは 2B, 3B, 4B, 5B, 6B, 7B, 8B, 9B, 10B, 12B, 14B, 16B, 18B, 20B, 22B, 24B, 26B, 28B, 30B, 33B, 36B, 39B, 42B, 45B, 50B, 55B, 60B, 65B, 70B, 80B と取り、これら 30 種類のスライド窓のサイズにおける圧縮率を求める。同様に API コール列を逆順に並べたものに対しても同じ処理を行い、各スライド窓のサイズにおける圧縮率を求める。求めた圧縮率を API コール列の正順、逆順という順番で、スライド窓サイズの小さい順に並べ、60 次元の特徴ベクトルを用いる。例えば “Trojan.Win32.Agent.rxht” の特徴ベクトルを図 5 に示す。マルウェアファミリー “Trojan.Win32.Agent” の分類器を作成する場合は、学習用特徴ベクトルとして図 5 のベクトルの先頭に 1 を付加したベクトルを作成し SVM で利用する。

6. 提案手法によるマルウェア亜種推定の実験

本実験では、動的解析結果として FFRI dataset 2013～2016 のデータを対象として実験を行う。但し、前述したように、FFRI Dataset 2014 では Cuckoo Sandbox での解析に加えて FFRI yarai analyzer Professional の解析結果ログも含まれているが、解析方法を統一するため、Cuckoo Sandbox による解析結果のみを用いる。また、FFRI Dataset 2016 には Windows10 と Windows8.1 の解析結果が含まれているが、本稿では Windows10 の動的解析結果のみを用いる。また、亜種判定における正解として、Kaspersky が命名した名称を基準として分類を行う。

API のコールパターンに着目し SVM を用いた判定を行うため、全 16887 検体のうち、API コールが取得できるマルウェアであり、亜種の数 が 100 検体以上ある 22 種類の

表 1 マルウェアファミリー名

Table 1 The malware family names

マルウェアファミリー名	亜種検体数
Trojan.Win32.WaldeK	687
Trojan-Spy.Win32.Zbot	581
Trojan.Win32.Yakes	577
Backdoor.Win32.Androm	445
Trojan.Win32.Agent	434
Trojan.Win32.Inject	351
Trojan.Win32.WBNA	308
Hoax.Win32.ArchSMS	304
Trojan-PSW.Win32.Fareit	272
Trojan-PSW.Win32.Tepfer	218
Backdoor.Win32.DarkKomet	206
Trojan-Ransom.Win32.Foreign	204
Trojan-Dropper.Win32.Injector	185
Trojan.Win32.Jorik	171
Packed.Win32.Tpyn	163
Trojan.Win32.Kovter	145
Trojan.Win32.Scar	142
Downloader.Win32.LMN	123
Worm.Win32.Vobfus	120
Backdoor.Win32.Matsnu	110
Trojan-Downloader.Win32.Upatre	108
Trojan.Win32.Llac	100

マルウェアファミリーを判定対象とした。22 種類のマルウェアファミリーの名前とマルウェアファミリーに含まれる検体の数の一覧を表 1 に示す。

佐藤らの手法 [5] の結果と条件を揃えるため、判別するマルウェアファミリーとそうではないマルウェアファミリーの検体数が 21:1 となるように検体数を揃えて、マルウェアファミリー毎に 22 種類の分類器を作成する。

また亜種判定には 4 分割交差検証を行う。つまり 4 分割したデータセットのうちの一つをテストセットとし、残りをトレーニングセットとする。4 個に分割されたデータセットそれぞれをテストセットとして 4 回検証を行う。

ここで、SVM には LIBSVM[9] を用い、グリッドサーチを行い最適な数値を求めた。

提案手法の推定精度の結果を表 2 に示す。また同様の実験を Paragraph Vector を用いて行っている佐藤らの手法 [5] による実験結果 (API 名のみを扱った手法 1 の結果) も示す。

表 2 より 22 種類のマルウェアファミリーのうち 9 種類のマルウェアファミリーについて、佐藤らの手法より推定精度が向上したことがわかる。また、推定率が 90% を超えたマルウェアは 22 種類中 6 種類となった。この中で最も推定率が高いものは “Worm.Win32.WBNA” の 94.05% であり、最も推定率が悪いものは “Trojan.Win32.Agent” の 63.69% となった。

圧縮し、30種類のスライド窓で圧縮した場合の圧縮率を特徴ベクトルとする方法を提案した。機械学習の手法の一種であるSVMを用いることで、平均で約82%の精度でマルウェアの亜種を推定できることを示した。

今後の課題として、推定率が70%以下になるようなマルウェアファミリーに対して有効な手法の検討を行うことがあげられる。

参考文献

- [1] AV-TEST, Statistics, Malware, 入手先 (<https://www.av-test.org/en/statistics/malware/>) (2017-08-19)
- [2] 神菌雅紀, 畑田充弘, 寺田真敏, 秋山満昭, 笠間貴弘, 村上純一: マルウェア対策のための研究用データセット MWS Datasets 2013, 情報処理学会シンポジウムシリーズ (コンピュータセキュリティシンポジウム 2013 論文集), Vol.2013, No.4, pp.1-8(2013).
- [3] 株式会社 FFRI: FFRI Dataset 2014 のご紹介, 入手先 (http://www.iwsec.org/mws/2014/files/FFRI_Dataset_2014.pdf) (2017-08-05).
- [4] Securelist: Rules for naming, 入手先 (<https://securelist.com/threats/rules-for-naming/>) (2017-08-08).
- [5] 佐藤拓未, 後藤滋樹, 武部高礼: Paragraph Vector を用いたマルウェアの亜種推定法, 情報処理学会シンポジウムシリーズ (コンピュータセキュリティシンポジウム 2016 論文集), Vol.2016, No.2, pp.298-304(2016).
- [6] 岡野靖, 神谷和憲, 谷川真樹: データ圧縮アルゴリズムを用いたマルウェア感染通信ログの判定, 情報処理学会シンポジウムシリーズ (コンピュータセキュリティシンポジウム 2016 論文集), Vol.2016, No.2, pp.640-646, 2016.
- [7] 衛藤公希, 小野廣隆, 山下雅史, 竹内純一, 文字列圧縮を用いたネットワークセキュリティにおけるインシデント検出, 情報科学技術フォーラム講演論文集, Vol.11, No.1, pp.9-16, 2012.
- [8] Ziv J. and Lempel A.: A Universal Algorithm for Sequential Data Compression, IEEE Trans. Inf. Theor., Vol.23, No.3, pp.337-343(1977).
- [9] LIBSVM – A Library for Support Vector Machines, 入手先 (<https://www.csie.ntu.edu.tw/~cjlin/libsvm/>) (2017-08-05).
- [10] 堀合啓一, 今泉隆文, 田中英彦: マルウェア亜種の動的挙動を利用した自動分類手法の提案と実装, 情報処理学会論文集, Vol.50, No.4, pp.1321-1333(2009).