

OLAP(多次元データ分析)を利用した攻撃元データの分析と検体の自動解析

永田大 堀合 啓一 田中 英彦
情報セキュリティ大学院大学

あらまし OLAP とは、BI などでも利用されるデータを多次元に解析する手法である。本論では、研究用データセット CCC DATASET 2008 の攻撃元データを OLAP による分析を行い、分析の方法とその結果を考察する。また、筆者らはマルウェアの自動解析システムを構築し運用しているが、このシステムの概要と研究用データセット CCC DATASET 2008 のマルウェア検体の自動解析を行った結果についても報告する。

Analysis of CCC DATASET 2008 using OLAP and Automated malware variant classification

Dai Nagata, Keiichi Horiiai, Hidehiko Tanaka
INSTITUTE of INFORMATION SECURITY

Abstract The OLAP is a method to analyze multi-dimensional data in the area such as Business Intelligence. In this report, we describe the method of applying OLAP to the analysis of CCC DATASET 2008 and its results. Authors have been developing an automated malware analyzing system. We have adopted the system to analyze the specified malware in the CCC DATASET 2008 and we describe the outline of the system and the obtained results.

1 はじめに

OLAP (On-line Analytical Processing) とは、ビジネスインテリジェンス等で利用されるデータを多次元に解析する手法である。本論では、研究用データセット CCC DATASET 2008 の攻撃元データ(以降、攻撃元データ)を分析するに当たり OLAP を使用した。多次元でデータを扱うことにより、マルウェアの傾向分析を行った。また、筆者らは定点観測と連動したマルウェアの自動解析システムを構築し運用しており、このシステムを拡張し、研究用データセット CCC DATASET 2008 のマルウェア検体(以降、マルウェア検体)の挙動解析を行った。

2 OLAP の概要

OLAP は、ビジネスインテリジェンス (BI) と呼ばれる経営分析等で用いられる。企業の業

務システム等で蓄積されたデータの分析・加工を行い、企業経営の意思決定に活用しようとする手法である。このように OLAP は、一般的には経営計画や企業戦略に利用される。BI ツールでは、論理的に 3 次元以上の視点を持つ分析用データベースから 2 次元のクロス集計表としてデータを取り出しデータの多次元的な分析を行うことができる。

本論では、攻撃元データについて多次元データとして解析することを試み、その手法として OLAP を利用することとした。OLAP にてデータを扱うことにより、基礎数値の把握と傾向分析を簡潔に行うことができるのではないかと考えた。

3 分析システムの概要

本研究では、OLAP のソフトウェアとして、

OpenStandia[1]を利用した。OpenStandia は、OpenOLAP をベースにしたフリーウェアである。OpenStandia では、データベースに MySQL が使われている。図 1 に構成を示す。

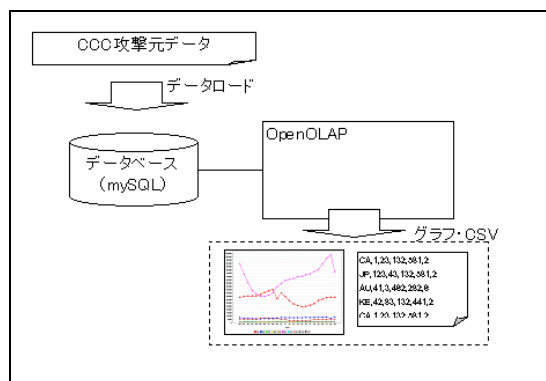


図 1 システムの構成

OLAP では、リレーショナルデータベースのデータを特定のタイミングで多次元データとして再構成する。再構成されたデータをキューブと呼ぶ。キューブに対して問い合わせを行うことによりグラフや CSV の元となる分析情報を取得することができる。またキューブは、複数のディメンジョンで構成される。ディメンジョンとは多次元的にデータを集計する際の軸である。

3.1 テーブル構造

OLAP でのデータ解析にあたり、MySQL にログのデータをロードするが、そのテーブル構成を表 1 に示す。ログの値と同じものに加え、マルウェアの Family 名や CountryCode を追加した。マルウェア名称についてカンマより前の部分を Family 名として扱う。例えば PE_BOBAX.AH であれば、PE_BOBAX を Family 名とする。追加した項目は、他のカラムからの導出が可能であるが、性能を考慮し、あらかじめデータを設定しておく方針とした。

3.2 ディメンジョンの構成

今回の分析では、表 2 に示すディメンジョン

を作成した。これらの各ディメンジョンを組み合わせて多次元的にデータを分析した。

表 1 ログ管理テーブル

カラム名	説明
pk	プライマリキー (追加項目)
datetime	日時
ip	攻撃元 IP アドレス
port	ポート番号
direction	攻撃方向
hash	ハッシュ値
virusname	マルウェア名
filename	ファイル名
virusnamegroup	マルウェアの Family 名 (追加項目)
countrycode	CN (追加項目)
24h	24 時間表記での受信時刻 (追加項目)
24h_timezone	24 時間表記での受信時刻で日本との時差を考慮したもの (追加項目)

表 2 使用したディメンジョン

名称	説明
virusgroupname	マルウェアの Family 名を扱う
countrycode	攻撃元 IP アドレスから導出した国コードを扱う
24 時間	24 時間表記での時刻時間を扱う
ハッシュ値	検体のハッシュ値を扱う

4 分析結果

本システムを使用し次の分析を行った。

- ・全体基本情報の把握
- ・国別、時間帯別の分析
- ・ハッシュ値、マルウェア名称の分析

4.1 全体基本情報の把握

マクロ的な概要をつかむため、全体の基本と

なる統計情報を取得した。本システムでは、データをテーブルにロードしているため、全体の種類数や種類別出現数の情報は、SQLにて容易に取得可能である。表3にSQLで取得した件数情報を記述する。また表4にマルウェア名称についての種類別出現件数の情報を記述する。

表3 全体基本情報 種類数

全ログレコード数	2,942,221 件
マルウェア名称による種類数	1,082 種類
ハッシュ値による種類数	52,465 種類
攻撃元 IP 種類数	258,711 種類
port/プロトコルによる種類数	63,820 種類
攻撃国別種類数	150 種類

表4 マルウェア名称別出力件数（上位5）

UNKNOWN	651,492
PE_BOBAX.AH	130,708
BKDR_VANBOT.AX	97,176
BKDR_VANBOT.AD	68,785
TROJ_POEBOT.AGU	66,649

4.2 国別、時間帯別グラフ

鬼頭[2]らは、不正ホストからの攻撃を時間軸での分析を行い、その中でそれぞれの国における日本時間を現地時間に修正し、時間帯別による攻撃数の変化のグラフ化を行っている。本システムで同様に、国別・時間帯別の集計を行った。図2は、出現数の多い国から10カ国（JP,CA,CN,US,TW,CZ,KR,DE,IN,PH）に対して、時差を考慮せずにグラフ化したものである。また、図3は、同じ国について時差を考慮してグラフ化したものである。

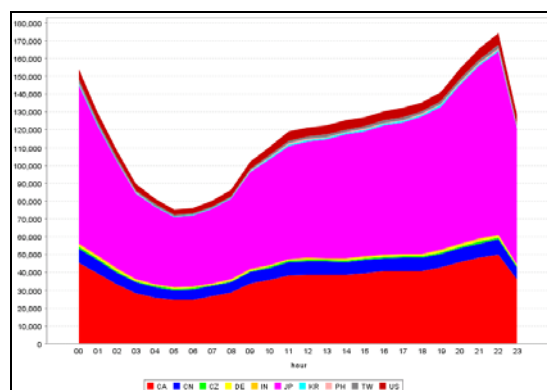


図2 時間帯別攻撃数

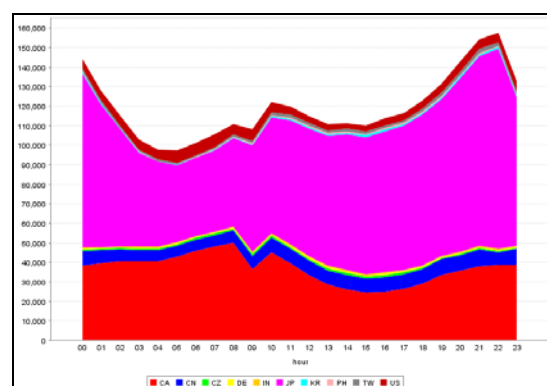


図3 時間帯別攻撃数（時差補正）

今回の分析では、時差補正を行った場合に顕著な特徴見られなかった。

4.3 ハッシュ値、マルウェア名称の分析

ハッシュ値とマルウェア名称の関係を分析した。図4は、ウィルスパターンが更新されるまでの期間を視覚化したものである。横軸は月を示し、縦軸はマルウェアの検体を示す。UNKNOWNの期間が赤色で、名前が付くと黄色、さらに名前が変わると別の色で表現している。図5は、UNKNOWNの検体へマルウェアの名称が付与されるまでの期間を示したものである。

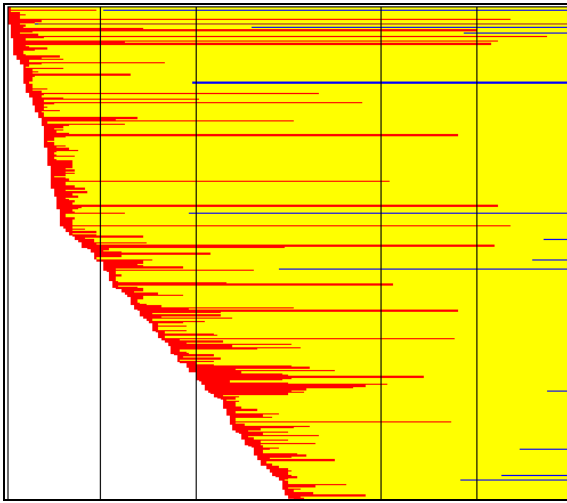


図4 ウィルスパターン更新の可視化

ハッシュ値	Family名 (1番目)	1番目の名称がはじ められた日時	Family名 (2番目)	2番目の名称がはじ められた日時	経過時間 (日)
004f9e4571f1e0bde437a71783e7f0235846c	PE_VIRUT	2008.02.23 19:31	PE_VIRUT	2008.02.29 09:47	34.0
01726e04844e48108c3c203e9ad508f0c2c2c	PE_VIRUT	2008.06.11 04:38	PE_VIRUT	2008.05.25 10:05	28.6
2249f1c0e38e70e0e0e0e0e0e0e0e0e0e0e0	PE_VIRUT	2007.11.20 00:01	PE_VIRUT	2007.11.11 18:04	1.0
008a280200000000000000000000000000	PE_VIRUT	2008.02.07 17:26	PE_VIRUT	2008.02.29 09:56	17.7
008f9e4571f1e0bde437a71783e7f0235846c	PE_VIRUT	2007.11.03 23:21	PE_VIRUT	2007.11.14 16:58	12.4
0a72e0e072726e0e0e0e0e0e0e0e0e0e0e0	PE_VIRUT	2007.11.04 02:04	PE_VIRUT	2007.11.11 16:02	1.9
00a58348b0c0f07339f46f0a0a0a0a0a0a0a0	PE_VIRUT	2008.07.09 00:02	PE_VIRUT	2008.02.29 20:10	88.8
f47e00884c000000000000000000000000	PE_VIRUT	2007.12.06 18:56	PE_VIRUT	2008.01.29 20:58	58.1
11e189f457a00000000000000000000000	PE_VIRUT	2008.02.20 19:54	PE_VIRUT	2008.04.15 22:47	56.1
1130e0a81e000000000000000000000000	PE_VIRUT	2007.11.29 17:44	PE_VIRUT	2007.12.01 03:08	11.0
10000a00101000100010001000100010001	PE_VIRUT	2008.02.02 09:34	PE_VIRUT	2008.04.19 22:41	74.1
10000a00101000100010001000100010001	EMUL_DRM_EXEC	2007.11.03 00:04	EMUL_DRM_EXEC	2007.12.23 17:46	58.5
1a0f20a7a0000000000000000000000000	PE_VIRUT	2008.04.12 21:34	PE_VIRUT	2008.04.29 19:16	107.9
1a0f20a7a0000000000000000000000000	PE_VIRUT	2008.07.23 20:53	PE_VIRUT	2008.04.19 21:48	89.1
1f10e0c078e0d0e0000000000000000000	PE_VIRUT	2008.04.23 21:26	PE_VIRUT	2008.04.23 21:52	1.0
2081e0c078e0d0e0000000000000000000	EMUL_DRM_EXEC	2007.11.01 00:04	EMUL_DRM_EXEC	2007.12.23 17:58	88.7
2081e0c078e0d0e0000000000000000000	PE_VIRUT	2008.02.23 23:39	PE_VIRUT	2008.02.29 03:51	8.3
2081e0c078e0d0e0000000000000000000	PE_VIRUT	2008.02.23 02:42	PE_VIRUT	2008.02.29 00:15	2.9
2081e0c078e0d0e0000000000000000000	PE_VIRUT	2007.11.22 23:05	PE_VIRUT	2007.12.11 05:25	18.3
330e0c078e0d0e00000000000000000000	PE_VIRUT	2008.03.15 14:02	PE_VIRUT	2008.02.29 03:51	35.3
330e0c078e0d0e00000000000000000000	PE_VIRUT	2007.11.10 06:40	PE_VIRUT	2007.11.18 00:04	5.7
330e0c078e0d0e00000000000000000000	PE_VIRUT	2007.12.02 02:17	PE_VIRUT	2007.12.09 22:26	6.8
330e0c078e0d0e00000000000000000000	PE_VIRUT	2007.12.02 00:04	PE_VIRUT	2007.12.01 00:18	2.4
330e0c078e0d0e00000000000000000000	PE_VIRUT	2007.11.14 18:00	PE_VIRUT	2007.12.06 20:00	22.0
330e0c078e0d0e00000000000000000000	PE_VIRUT	2008.02.15 15:34	PE_VIRUT	2008.02.29 14:14	13.9
330e0c078e0d0e00000000000000000000	PE_VIRUT	2008.02.24 14:16	PE_VIRUT	2008.02.29 11:07	30.8
330e0c078e0d0e00000000000000000000	PE_VIRUT	2008.02.23 23:39	PE_VIRUT	2008.04.05 00:06	36.0
330e0c078e0d0e00000000000000000000	PE_VIRUT	2007.12.06 00:00	PE_VIRUT	2007.12.07 22:42	11.8

図5 マルウェア名称付与までの期間

これらの図によって、マルウェア名称の変化（定義体を更新される期間）を見ることが出来る。マルウェア名称が途中で変化しているもので、UNKNOWNを除外したハッシュ値は、247件であった。247件の中でFamily名の変化に着目して確認したところ、Family名が変わった検体が109件、Family内だけで変わった検体が138件（内126件はPE_VIRUT Family）という結果であった。マルウェアの名称が途中で変わった247件のうち、ほぼ半数にあたる126件がPE_VIRUT Familyの検体であった。

PE_VIRUTは、2008年1月のCCC分析レポート[3]によると、2007年11月末より増えている。1日300種の異なる検体が発生すると報告されている。マルウェア名称が変更されるということは、はじめの定義ファイル提供時の分析

の時点から分析の結果が変わっているの、解析が難しいマルウェアや、新しいタイプのマルウェアと考えられる。

4.4 Excelによるグラフと簡易アニメーション化

OpenOLAPでは、Viewerから分析結果をCSV形式のファイルに出力することができる。その出力データを利用して、Excelによる簡易アニメーションを行った。時系列に複数のシートにデータを配置した。グラフの元データを連続的に切り替えることにより、簡易アニメーションとして、傾向の変化を見る事ができる。図6は、国とマルウェア名称を軸に取り、アニメーション切り替えを時間帯によって実施した場合の表示例である。

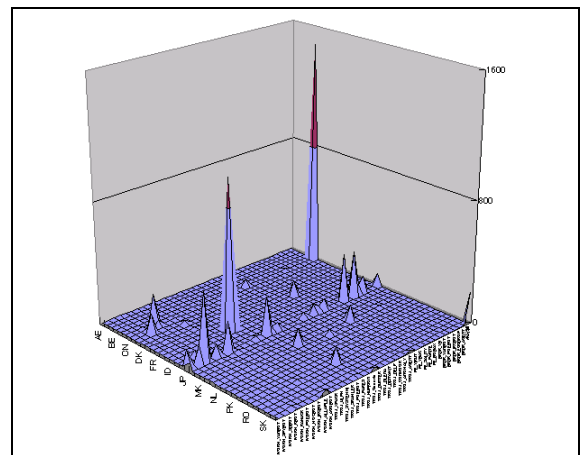


図6 アニメーション表示例

5 マルウェア検体の自動解析

筆者らは定点観測と連動したマルウェアの自動解析システムを構築し運用している[4,5]. このシステムは、ハニーポットを定点観測のセンサーとして利用し、捕獲したマルウェアを一般のネットワークから隔離した安全な環境で実行して、その挙動を自動的に解析するものである。解析の対象は、マルウェアの実行前後のWindowsのファイルやプロセス等の変化と、観測される通信の状況である。また、このシステ

ムは解析済みデータをDBへ蓄積し、新たに捕獲したマルウェアの挙動とDB内のデータとの類似性を自動的に算出して、その名称を推定する機能を有している。マルウェアの実行環境は、感染後に必要となるWindowsシステムの復旧の容易性から仮想マシン上で構築している。

このシステムを利用して、指定されたマルウェア検体の解析を試みたが、arpパケットが観測されるのみで、その他の挙動は観測されなかった。近年のマルウェアの中には、自身が解析されることを妨害する機能を持ったものの存在が知られていることから、今回の検体がこれに該当すると想定し、仮想マシンを利用しない環境を構築して解析を試みた。

5.1 解析環境の概要

自動解析は解析全般の処理を制御するPC(制御PC)とマルウェアを実行するPC(感染PC)の2台で構成した。制御PCはLinuxをOSとし、解析全般の制御を行う他、次の処理を行う。

- ・解析環境に必要なネットワークの模擬
- ・DNS,IRCなど関連サーバの模擬
- ・マルウェア実行中の通信パケットの記録
- ・感染PCから取得したマルウェア実行前後のログを比較して得られる挙動の変化の抽出とHTTPサーバと連動した解析結果の出力

一方の感染PCは、マルウェアを実行するWindows OS起動用のパーティションと、このパーティションのファイルを感染前の状態へ復旧させるための環境を備えたLinux起動用のパーティションのデュアルブート構成とし、以下のステップで処理を行う。

(1) Linuxのパーティションから起動し、感染前の正常な状態のWindowsのディスク・イメージ・ファイルを、Windows用のパーティションへコピー。

(2) ブートマネージャGrub[6]の機能を利用し、当該PCをリブートした際に、Windowsのパーティションから1回だけ起動するように設定して感染PCをリブート。

(3) 起動したWindowsは、マルウェア実行前の状態を記録し、その後に制御PCから解析対象のマルウェアファイルを受信して、このファイルを実行する。

(4) マルウェア実行後の状態を記録し、実行前のログとともに制御PCへ転送する。

(5) 感染PC自身をリブートし、(1)のステップから繰り返す。

基本的には制御PCと連携し、以上のステップを繰り返すことによって、全自動でマルウェアの解析が可能であるが、マルウェアの実行に伴うWindowsシステムのクラッシュ等の原因で(4),(5)が機能しない場合が発生する。この状態を検出し、制御PCからの指令によって感染PCをハード的にリブート(電源のOFF/ON)させる仕組みを持たせている。

5.2 解析結果

解析結果は、Webブラウザで閲覧可能なシステムとして実装しているが、マルウェア検体の解析結果の表示例を図7に示す。図において②レジストリの改ざんが1件、③システム関連ファイルの生成と削除がそれぞれ1件、④NVCOMプロセスの起動などが記録されている。また、⑤には起動したプロセスと、このプロセスが待ち受けているポート番号、⑥にはHOSTSファイルの改ざんの状態が表示される。解析対象の検体は、HOSTSファイルの改ざんを行うことが知られているが[7]、本自動解析ではファイルの最後に1バイト(¥x0a)が追加されていることを確認している。これらの解析結果は、ウィルス対策製品ベンダーなどがWeb上で公開している情報[7]とほぼ一致している。⑦にマルウェアの実行に伴って発生した通信パケットの情報を示す。マルウェアは解析環へ実装した模擬DNSサーバへホスト名rx7.teensmutbox.comのFQDNの問い合わせを行い、模擬DNSが指定した仮のIPアドレス4.3.2.234のIRCサーバへログインする様子が示されている。

(172.16.50.128は感染PCのIPアドレスである。)

このように、既知のプロトコルを利用するマルウェアについては、動的な挙動を観測することで、ボットネットの指令サーバを自動的に特定し、指令文字列を自動的に抽出することが可能である。

本自動解析システムはこれまでに述べた挙動をカテゴリデータへ変換してベクトル化し、ベクトル間のハミング距離を指標とした挙動の類似性を基に、マルウェアの名称を推定する機能を有している[5]。筆者らが仮想マシン上の環境で自動解析を行って蓄積した約 7,600 検体のデータとの挙動の類似性から推定したマルウェアの名称が①に示されている。この例では、推定した名称が BKDR_VANBOT.NM となり、実際と一致しなかった。その原因としては、実行環境の違い（仮想マシン上と実マシン上）に加えて、蓄積している解析済の DB 内に、今回の検体と挙動が類似しているデータが含まれていないことが原因と推定している。

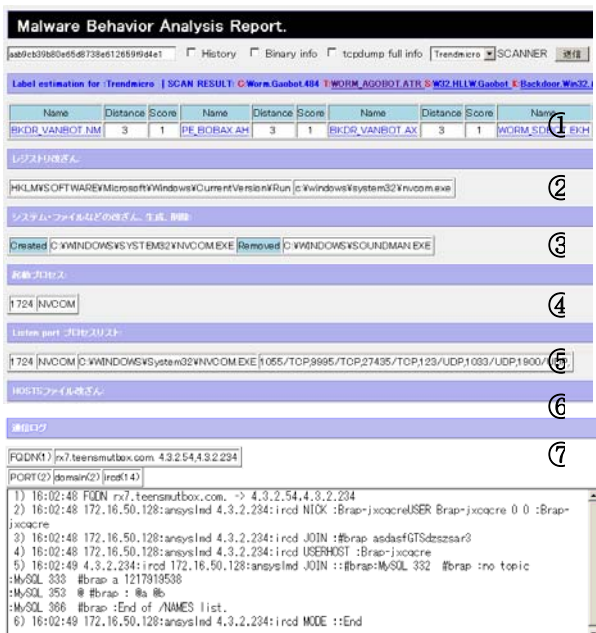


図 7 検体の自動解析結果

6 まとめ

OLAP を利用することにより、準備されたディメンジョンに対しては、容易にグラフ化・CSV

化をすることができるようになった。しかしながら、事前にディメンジョンやキューブを作成する必要があるためリアルタイムでの分析は難しい。

マルウェアの自動解析については、挙動の解析については所期の目的を達することができた。ただし、挙動の類似性を基にしてマルウェアの名称を推定する機能については、DB の充実やアルゴリズムの改善による精度の向上が今後の課題である。

参考文献

[1] OpenStandia

<http://www.nri-aitd.com/openstandia/>

[2] 鬼頭哲郎, 仲小路博史, 寺田真敏, 菊池浩明「インターネット上の不正ホスト分布に関する社会的レイヤからの考察」情報処理学会研究報告. CSEC Vol.2007, No.71(20070719), Jul.2008

[3] 2008 年 1 月 CCC 分析レポート

<https://www.ccc.go.jp/report/200801/0801monthly.html>

[4] 堀合啓一, 今泉隆文, 田中英彦「定点観測によるボットネットの観測と Malware の動的挙動解析システムの提案」情報処理学会論文誌 Vol.49 No.4 1-12, Apr.2008

[5] 堀合啓一, 今泉隆文, 田中英彦「ハミング距離によるマルウェア亜種の自動分類」CSEC Vol.2008, No.45(20080523), May.2008

[6] GNU GRUB

<http://www.gnu.org/software/grub/>

[7] WORM_AGOBOT.ATR

http://www.trendmicro.co.jp/vinfo/virusencyclo/default5.asp?VName=WORM_AGOBOT.ATR&VSection=T