

分散ハニーポット観測からのダウンロードサーバ間の相関ルール抽出

大類 将之† 菊池 浩明† 寺田 真敏††

† 東海大学情報理工学部情報メディア学科

259-1292 神奈川県平塚市北金目 1117 yama, kkn@cs.dm.u-tokai.ac.jp

†† 日立製作所 Hitachi Incident Response Team (HIRT)

212-8567 神奈川県川崎市幸区鹿島田 890 日立システムプラザ新川崎

あらまし 本研究では、マルウェアをダウンロードするサーバ間の連携活動に着目し、データベースから価値ある相関ルールを抽出するデータマイニング手法であるアソシエーション分析を CCC DATASET 2009 攻撃元データに適用する。分析の結果明らかになった関連性が強いサーバ間の組み合わせやダウンロードホストとマルウェアの関連性を示した相関ルールを報告する。

Mining Association rules consisting of Download Servers from Distributed Honeypot Observation

Masayuki Ohru† Hiroaki Kikuchi† Masato Terada††

† School of Information Science and Technology, Tokai University,
1117 Kitakaname, Hiratsuka-shi, Kanagawa 259-1292

†† Hitachi, Ltd. Hitachi Incident Response Team (HIRT),
890 Kashimada, Saiwai-ku, Kawasaki-shi, Kanagawa 212-8567

Abstract This paper aims to find interested association rules, known as data mining technique, out of the dataset of downloading logs by focusing on the coordinated activity between downloading servers for malware. The result of the analysis shows the association rules of the relation among the downloading servers and that of the malwares.

1 はじめに

ポットにおける感染に関して、複数のサーバを連携して感染する特徴が報告されている [1]。表 1 は、CCC DATASET 2009 攻撃通信データ [2] の中から、PE_VIRUT.AV に感染したとき、数分後に TROJ_BUZUS.AGB と WORM_SWTYMLAI.CD が同時に感染している例である。

これより、ダウンロードサーバ（以下、DL サーバ）が異なる IP アドレスであっても、マルウェア（以下、MW）の感染活動には類似性がある事が分かる。このようなポットネットの複数の DL サーバによる感染を連携活動と呼ぶ事とする。

しかし、このような連携活動を発見するためには、多量のデータの中から共通に生じるパターンを抽出する必要があり、非常に困難である。例えば、CCC

DATASET 2009 攻撃元データ [2] を対象に、その出現数トップ 4 の DL ホスト IP アドレスのダウンロード数の推移を示した図 1 について考えよう。

IP アドレスや時期によるダウンロード数の変動は大きい。全く観測しない期間も存在し、ハニーポットによっても結果は異なる。そのため、連携活動だと考えられる特徴を発見しても、その真偽を測るのは難しい。1 年間で 1,335 種類の MW が観測され [2]、先の例のように 3 種類の異なる連携を考えるだけでも、 ${}_{1,335}C_3 = 395,654,395$ 個組み合わせがある。しかも、 $365 \text{ 日/年} \times 94 \text{ 台} \times 24 \text{ 時間/日} \times 3 \text{ スロット/時間} = 2,470,320$ の全てのスロットについてそれらが成立するかを総当りで調べるのは現実的ではない¹。

そこで本研究では、大規模データからのマイニング

¹ スロットの定義は 3.1 節にて後述する。

表 1: 攻撃通信データから抽出した MW の連携活動例

時刻	DL ホスト IP アドレス	Dst Port	プロトコル	MW 名
0:02:11	124.86.***.111	47556	TCP	PE_VIRUT.AV
0:03:48	67.215.*.206	80	TCP	TROJ_BUZUS.AGB
0:03:48	72.10.***.195	80	TCP	WORM_SWTYMLAI.CD
0:36:46	124.86.**.109	33258	TCP	PE_VIRUT.AV
0:36:52	72.10.***.195	80	TCP	WORM_SWTYMLAI.CD
0:36:52	67.215.*.206	80	TCP	TROJ_BUZUS.AGB
0:46:56	124.86.**.109	33258	TCP	PE_VIRUT.AV
0:48:52	67.215.*.206	80	TCP	TROJ_BUZUS.AGB
0:48:52	72.10.***.195	80	TCP	WORM_SWTYMLAI.CD

表 2: 攻撃元データから抽出した DL ホスト IP アドレス上位 10 個のデータ

順位	DL ホスト IP アドレス	DL 回数	平均 DL 回数	MW 数	ハニーポット数
TOP1	72.10.***.74	462246	3884.4	119	91
TOP2	72.10.***.195	399562	8324.2	48	92
TOP3	85.114.***.2	33283	1147.7	29	82
TOP4	85.114.***.207	32202	870.3	37	78
TOP5	67.215.*.206	26780	3825.7	7	59
TOP6	211.95.**.6	19641	198.4	99	85
TOP7	72.10.***.26	14951	287.5	52	82
TOP8	92.48.**.63	11699	117.0	100	69
TOP9	67.18.***.250	10060	76.8	131	68
TOP10	72.8.***.164	5099	127.5	40	81

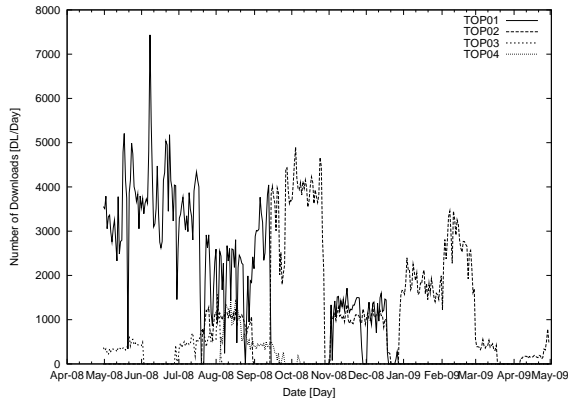


図 1: 1 年間で観測されたダウンロード数の推移

技術であるアソシエーション分析を適用する事で、連携活動に関する価値ある相関ルールの抽出を試みる。アソシエーション分析(相関ルール)には、Agrawal らによって提案された Apriori アルゴリズムがあり、支持度 (Support) と確信度 (Confidence) の最小値を設定することで、膨大な組み合わせを効果的に枝狩りして効率よくルールを抽出できる技術として広く知られている [3]。

本稿では、その実装として Christian Borgelt による Apriori Program [4] を活用して、攻撃元データに適用し、1 年間に頻出している共通の攻撃パターン

を抽出する。攻撃パターンには、ダウンロードサーバ間の連携によるものと、MW の種類の組み合わせによるものがあり、それぞれについて、観測装置(ハニーポット)による差や観測時期による差が生じるか検証する。

2 要素技術

2.1 相関ルール

アソシエーション分析とは、データベースの中から X (前件部) \Rightarrow Y (結論部) という相関ルールを抽出するデータマイニング手法である。

支持度とは、ルールの出現率を表し、全トランザクション N のうち、ルールの条件 X と結論 Y を共に含む確率

$$Supp(X \Rightarrow Y) = \frac{X \cap Y}{N}$$

で定義する。

確信度とは、ルールの関連性の強さを表し、ルールの条件 X が発生するトランザクションのうち、条件 X と結論 Y を共に含む確率

$$Conf(X \Rightarrow Y) = \frac{X \cap Y}{X}$$

で与える。

表 3: トランザクションの例

TID	A	B	C	D	E
1	1		1	1	
2		1	1		1
3	1	1	1		1
4		1			1

2.2 Apriori アルゴリズム

Apriori アルゴリズムは, Agrawal らが提案した代表的な相関ルール抽出アルゴリズムである. 支持度と確信度に最小値を与え, 数多く抽出される相関ルールの中からルールを絞る事によって, 価値の低いルールを除き, 価値あるルールを発見できる.

表 3 で与えるトランザクションデータを用いて, $B, C \Rightarrow E$ の相関ルールの評価例を示す.

$$Supp(B, C \Rightarrow E) = 2/4 = 0.5$$

$$Conf(B, C \Rightarrow E) = 2/2 = 1$$

これから, $B, C \Rightarrow E$ のルールは支持度 50%, 確信度 100% である事が分かる. つまり, このルールは 50% の確率で出現し, B, C が発生した場合, 100% の確率で E が発生する事を示している.

3 調査方法

3.1 実験データ

実験データとして, 攻撃通信データと攻撃元データを使用する. 攻撃通信データは定期的にクリーンな状態にリセットされるため, Windows XP が送信する NTP パケットを利用して, 2 日分のデータをタイムスロット (以下, スロット) として 145 個に分割した. スロットを 1 つのトランザクションとし, その間にダウンロードされた MW の種類と DL サーバをそのトランザクションに生じるアイテムとして, 相関ルールを抽出する. 同様に, 攻撃元データもスロットに分割して, 相関ルールを抽出する [5]. なお, 本実験では, Windows XP のスロットを使用する.

3.2 調査項目

調査項目を以下に示す.

1. 攻撃通信データに対する MW の相関ルール抽出

2. 攻撃通信データに対する DL サーバの相関ルール抽出
3. 攻撃元データから抽出する MW の相関ルールの観測地点間の差
4. 攻撃元データから抽出する MW の相関ルールの観測時期の差

攻撃通信データは攻撃元データとマッチングを行う事により, MW 名を判定する事が可能である. そのため短期間ではあるが, 2009 年 03 月 13 日, 14 日に関しては正確なデータが得られ, 精度の高い相関ルールが期待できる. 本調査では, スロットごとに MW 名及び DL サーバの IP アドレスを抽出し, 2 つの観点からアソシエーション分析を行う事で, 連携活動の関連性の強さを調査する.

一方, 攻撃元データは 94 台のハニーポットを使用し, 1 年間観測したデータである. このデータを使用する事により, 抽出された相関ルールがどの位一般的であるか, すなわち, 異なるハニーポットで共通に観測されるかどうか, あるいは, 観測期間による差異は生じるかを明らかにする事ができる. すなわち, ハニーポット間及び長期間での 2 つの観点からアソシエーション分析を行う.

しかし, Honey003 及び 004[2] を除いて, 正確にスロットに分割するのは困難である. 誤差を許容して, 単純に 20 分ごとにスロットを分割して分析を行う.

4 調査結果

4.1 攻撃通信データに対する MW の相関ルール抽出

本節では, MW の連携活動の関連性を調査する. タイムスロットごとの MW リストの一部を表 4 に示す. 全 145 スロットのうち, 58 スロットが感染しており, 最大 1 スロットで 11 回感染している.

表 4 をアソシエーション分析した結果を表 5 に示す. 最小支持度 10% 以上, 最小確信度 80% 以上の全ての相関ルールを示している². ここで, 支持度は 145 個中その相関ルールが生じたスロットの割合, 確信度は前件部の MW を含むスロットのうち, 結

²関連性が高いルールを抽出するため, 支持度を低く, 確信度を高く設定している.

表 4: 攻撃通信データから抽出したスロットごとの MW 名リスト (一部)

スロット	MW 名			
0	PE_VIRUT.AV	TROJ_BUZUS.AGB	WORM_SWTYMLAI.CD	
2	WORM_ALLAPLE.IK	PE_VIRUT.AV	WORM_SWTYMLAI.CD	TROJ_BUZUS.AGB
3	PE_VIRUT.AV	TROJ_BUZUS.AGB	WORM_SWTYMLAI.CD	PE_VIRUT.AV
14	BKDR_POEBOT.GN	TROJ_BUZUS.AGB	WORM_SWTYMLAI.CD	
15	BKDR_MYBOT.AH	PE_VIRUT.AV		
⋮				
141	PE_BOBAX.AK	WORM_SWTYMLAI.CD	WORM_AUTORUN.CZU	WORM_IRCBOT.CHZ

論部の MW もダウンロードしているものの割合を表している。

表 1 の連携した攻撃パターン PE_VIRUT.AV ⇒ TROJ_BUZUS.AGB, WORM_SWTYMLAI.CD というルールは発見されなかったが, 類似した相関ルール PE_VIRUT.AV, TROJ_BUZUS.AGB ⇒ WORM_SWTYMLAI.CD 及び PE_VIRUT.AV, WORM_SWTYMLAI.CD ⇒ TROJ_BUZUS.AGB は高い確信度で発見された。また, 表 5 の結果から連携活動として, TROJ_BUZUS.AGB と WORM_SWTYMLAI.CD は特に強い関連性があると考えられる。

4.2 攻撃通信データに対する DL サーバの相関ルール

本節では, 4.1 節と同様に DL サーバの観点から関連性を調査する。DL サーバの IP アドレスで, アソシエーション分析を行った結果を表 6 に示す。最小支持度 10% 以上, 最小確信度 50% 以上である³。対応順位は, 表 2 から割り出している。ここで, 支持度は相関ルールの DL サーバと通信していたスロットの割合を, 確信度は前件部の DL サーバからダウンロードしていたもののうち, 結論部の DL サーバも通信しているスロットの割合を表している。

分析の入力データは, 表 4 の MW 名を対応する IP アドレスに置き換えたものだが, MW 名と DL サーバは 1 対 1 で対応している訳ではない。例えば, PE_VIRUT.AV は 16 種類の IP アドレスからダウンロードされている。

No.1, 2 の IP アドレスは PE_VIRUT.AV のものだが, 114.145.**.166 は 12 回, 122.18.**.123 は 21 回使用されており, 16 種類のうち上位の 2 種類である。また, No.3, 4 から IP アドレスの観点から見た場合でも連携活動として, TROJ_BUZUS.AGB と WORM_SWTYMLAI.CD の間には強い関連性がある事が

³確信度を 80% から 50% に下げたのは, 抽出されたルールが少なかったためである。

分かる。しかし, 複数の DL サーバを使用する MW, すなわち, PE_VIRUT.AV を含む表 1 に関するルールは抽出できなかった。

4.3 攻撃元データから抽出する MW の相関ルールの観測地点間の差

本節では, 4.1 節の相関ルールが異なるハニーポットで共通に観測されるかを調査する。全 94 のハニーポット ID でアソシエーション分析を行った結果を表 7 に示す。2009 年 3 月 13 日のみのデータで分析を行っており, 閾値はスロット数 3 以上, 確信度 80% 以上である⁴。例えば, No.1 のルールのハニーポット数 32 は, 全 94 台のハニーポットのうち, このルールを 1 回でも抽出したハニーポットが 32 台あることを表している。このとき, 支持度と確信度の違いは無視している。

上位の相関ルールの出現数 (相関ルールが成立するスロットの数) を図 2 に示す。X 軸はハニーポット数を k とし, Y 軸は k 以上のハニーポットで観測された異なる相関ルールの個数 $N(k)$ を表す。

結果は, TROJ_BUZUS.AGB と WORM_SWTYMLAI.CD の相関ルールが上位に抽出されたことを示しており, 概ね 1 台のハニーポットで観測した結果と矛盾がなかった。3 分の 1 以上のハニーポットにて, TROJ_BUZUS.AGB と WORM_SWTYMLAI.CD 間のルールが得られている。

また, 図 2 から, 共通したルールばかりでなく, 多様なルールが数多く生じることが分かる。広範囲で観測されたルールは特に連携活動を行っている可能性が高いと考えられ, これらルールには MW に偏りがみられる事から, 特定の MW のみが連携活動を行っていると考えられる。

⁴支持度の代わりにスロット数にしたのは, スロットごとに抽出されるルール数が違うため, ルール総数で正規化する支持度で揃えては共通のルールを抽出するのに不都合だったからである。

表 5: 攻撃通信データに対する MW の相関ルール

No.	前件部	結論部	支持度	確信度
1	TROJ_BUZUS.AGB	⇒ WORM_SWTYMLAI.CD	41.4	100
2	WORM_SWTYMLAI.CD	⇒ TROJ_BUZUS.AGB	46.6	88.9
3	TROJ_BUZUS.AGB BKDR_POEBOT.GN	⇒ WORM_SWTYMLAI.CD	10.3	100
4	BKDR_POEBOT.GN BKDR_POEBOT.GN	⇒ TROJ_BUZUS.AGB	10.3	100
5	PE_VIRUT.AV TROJ_BUZUS.AGB	⇒ WORM_SWTYMLAI.CD	29.3	100
6	PE_VIRUT.AV WORM_SWTYMLAI.CD	⇒ TROJ_BUZUS.AGB	29.3	100

表 6: 攻撃通信データに対する DL サーバの相関ルール

No.	前件部	結論部	支持度	確信度	対応 MW	対応順位
1	114.145.**.166	⇒ 122.18.**.123	12.1	85.7	PE ⇒ PE	
2	122.18.**.123	⇒ 114.145.**.166	15.5	66.7	PE ⇒ PE	
3	67.215.*.206	⇒ 72.10.**.195	46.6	100	TROJ ⇒ WORM	TOP5 ⇒ TOP2
4	72.10.**.195	⇒ 67.215.*.206	46.6	100	WORM ⇒ TROJ	TOP2 ⇒ TOP5

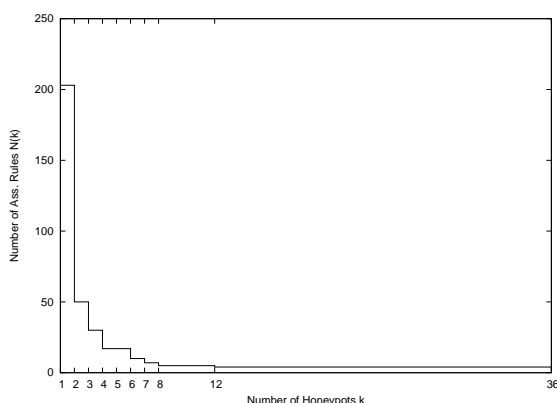


図 2: 上位 10 位の相関ルールの出現スロット数

4.4 攻撃元データから抽出する MW の相関ルールの観測時期の差

本節では、観測期間による相関ルールに差異が生じるかを調査する。1 年間の Honey003 でアソシエーション分析を行った結果を表 8 に示す。閾値はスロット数 3 以上、確信度 80 % 以上で 365 日全ての相関ルールを抽出し、該当した MW 名を含むルール数を月ごとに纏めたものである。4.3 節と同様に支持度と確信度の違いは無視している。また、抽出した MW 名はこれまでの調査で強いと考えられる PE_VIRUT.AV (PE), TROJ_BUZUS.AGB (TROJ), WORM_SWTYMLAI.CD (WORM) とした。

以下は、1 年間で観測された PE, TROJ, WORM が含まれる相関ルールのうち、上位 3 個である。

1. PE_VIRUT.AV WORM_SWTYMLAI.CD ⇒ TSPY_KOLABC.CH
2. TROJ_BUZUS.AGB ⇒ WORM_SWTYMLAI.CD

3. TSPY_KOLABC.CH ⇒ WORM_SWTYMLAI.CD

表 8 より、PE_VIRUT.AV は年間を通して観測されており、関係する多くのルールが抽出できている。その中でも 1 のルールより、WORM_SWTYMLAI.CD と関係するルールが一番多かった事から、長期間でも連携活動として、強い関連性があると考えられる。また、2, 3 のルールより、TSPY_KOLABC.CH とのルールが上位に来ており、4.1 節では見られなかった関連性、すなわち、観測期間による差異が伺える。

1 年間で観測された UNKNOWN を含まない相関ルールのうち、以下の上位 3 個のルールの推移を図 3 に示す。

1. BKDR_VANBOT.HI ⇒ BKDR_SDBOT.BU
2. BKDR_POEBOT.AHP ⇒ TROJ_QHOST.WT
3. TSPY_KOLABC.CH ⇒ WORM_SWTYMLAI.CD

これは、年間で最も観測されたルールであるが、ルールの出現期間は短く、1 年間を通して観測されていない。この事から、連携活動期間は短い事が分かる。理由としては、新たな MW が出現したり、更新されたりするために、特定の MW 間の連携期間が短くなってしまふ事が考えられる。

5 おわりに

連携して MW をダウンロードしているサーバや MW に関する規則を、機械的に抽出する方式を提案した。また、通信データを詳細に分析して得られる結果とほぼ同じ結果が抽出できる事を実証した。共通

表 7: 2009 年 3 月 13 日に観測された MW に関する相関ルールのハニーボット数

No.	前件部	結論部	ハニーボット数
1	TROJ_BUZUS.AGB	⇒ WORM_SWTYMLAI.CD	36
2	WORM_SWTYMLAI.CD	⇒ TROJ_BUZUS.AGB	36
3	TROJ_BUZUS.AGB BKDR_VANBOT.AHH	⇒ WORM_SWTYMLAI.CD	12
4	WORM_SWTYMLAI.CD BKDR_VANBOT.AHH	⇒ TROJ_BUZUS.AGB	12
5	TROJ_DLOADR.CBK	⇒ UNKNOWN	8
6	TROJ_BUZUS.AGB PE_VIRUT.AV	⇒ WORM_SWTYMLAI.CD	7
7	WORM_SWTYMLAI.CD PE_VIRUT.AV	⇒ TROJ_BUZUS.AGB	7
8	PE_VIRUT.AV TROJ_BUZUS.AGB	⇒ WORM_SWTYMLAI.CD	6
9	TROJ_AGENT.ANDF	⇒ UNKNOWN	6
10	PE_VIRUT.AV WORM_SWTYMLAI.CD	⇒ TROJ_BUZUS.AGB	6

表 8: 1 日ごとに観測された PE, TROJ, WORM を含むルール数

月	PE	TROJ	WORM
2008/05	31	0	0
2008/06	76	0	0
2008/07	111	0	0
2008/08	5	0	0
2008/09	8	0	0
2008/10	44	0	0
2008/11	27	0	0
2008/12	35	0	0
2009/01	135	0	0
2009/02	125	0	226
2009/03	79	53	74
2009/04	30	0	0

に現われたのは, PE_VIRUT.AV, TROJ_BUZUS.AGB, WORM_SWTYMLAI.CD 間の強い関係である.

本実験を総合して, TROJ と WORM の関連性が強い相関ルールが最も頻出していた. 広範囲で観測されたルールは, 連携活動を行っている可能性が高い事を示し, 長期間で観測した結果では, 連携活動期間は短い事が分かった. DL サーバ間の相関ルールを抽出する事は, 複数の DL サーバを使用する PE_VIRUT.AV 等があるため難しいが, 4.1 節と 4.2 節の結果から, MW の相関ルールを抽出し, それを IP アドレスに適用する事が有効であると考えられる.

Honey003 及び 004 を正確に分割した場合と単純に分割した場合で比較した結果, 支持度及び確信度の違いはあったが, 抽出された上位の相関ルールは同じだったため, 分析結果は妥当であるとする.

謝辞

本研究を遂行するにあたり, ご助言及びご協力を下さった日立製作所の藤原将志氏, 鬼頭哲郎氏, 仲

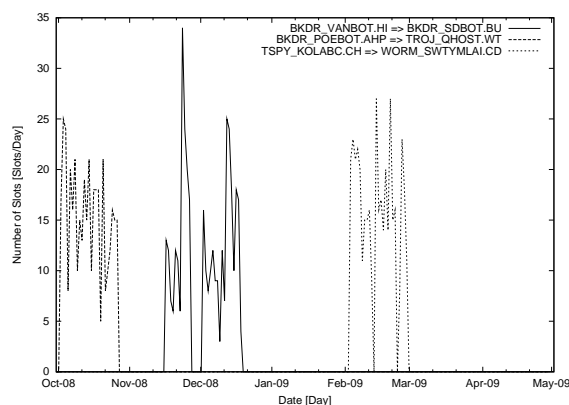


図 3: 1 年間の UNKNOWN を含まない上位 3 個の相関ルール数の推移

小路博史氏, 東海大学大学院の松尾峻治氏に深く感謝する.

参考文献

- [1] 桑原, 他, “パケットキャプチャーから感染種類を判定する発見的手法について”, マルウェア対策研究人材育成ワークショップ 2009 (MWS2009), 2009 (発表予定).
- [2] 畑田, 他, “マルウェア対策のための研究用データセットとワークショップを通じた研究成果の共有”, マルウェア対策研究人材育成ワークショップ 2009 (MWS2009), 2009 (発表予定).
- [3] Agrawal R, Imielinski T, Swami AN, “Mining Association Rules between Sets of Items in Large Databases”, Proceedings of ACM SIGMOD-93, pp. 207-216, 1993.
- [4] Christian Borgelt, Apriori - Association Rule Induction, <http://www.borgelt.net/apriori.html>
- [5] 小堀, 他, “マルウェアの通信履歴と定点観測の関連について”, マルウェア対策研究人材育成ワークショップ 2008 (MWS2008), 2008.