

明治大学総合数理学部
2018年度

卒業研究

ドメイン情報とHTTPレスポンスヘッダに基づく
フィッシングサイトの識別と評価

学位請求者 先端メディアサイエンス学科

桜井啓多

目次

| | |
|-----------------------------------|----|
| 1. はじめに | 1 |
| 1.1 研究背景 | 1 |
| 1.2 研究目的 | 1 |
| 2. 収集したデータについて | 2 |
| 2.1 概要 | 2 |
| 2.2 ドメインに関するデータ | 3 |
| 2.3 国に関するデータ | 4 |
| 2.4 HTTPレスポンスヘッダに関するデータ | 6 |
| 2.5 トップレベルドメインに関するデータ | 6 |
| 3. フィッシングサイト検出支援システム | 8 |
| 4. おわりに | 11 |
| 参考文献 | 12 |
| 謝辞 | 13 |
| 付録A 実際に収集したウェブサイトのデータ | 14 |
| 付録B FFRIデータセットのマルウェアのAPI情報による識別実験 | 16 |
| B.1 はじめに | 16 |
| B.1.1 研究背景 | 16 |
| B.1.2 研究方法 | 16 |
| B.2 収集したデータについて | 17 |
| B.2.1 概要 | 17 |
| B.2.2 エリア毎のデータ | 17 |
| B.2.3 国毎のデータ | 17 |
| B.2.4 パスワードのデータ | 18 |
| B.2.5 ユーザ名のデータ | 19 |
| B.3 FFRIのデータについて | 20 |
| B.3.1 概要 | 20 |
| B.3.2 送信先の国について | 21 |
| B.3.3 API | 21 |
| B.4 まとめ | 22 |
| 参考文献 | 23 |

1.はじめに

1.1 研究背景

フィッシングとはユーザから情報を盗む詐欺行為のことを指す。一般的なフィッシングの手口は[1]によると、初めに攻撃者が標的のユーザに偽のサイトのリンクを添付したメールを正規の送信元を装い送信する。そのユーザが偽のサイトでIDやパスワードを入力し、攻撃者の元に送信される。

具体的なフィッシング詐欺の被害として、[2]によると2017年下半期にフィッシング対策協議会に届け出されたフィッシングサイトのURLの件数は約5000件にのぼっている。また、[3]によると、日本からフィッシングサイトへ誘導された件数は2018年上半期において過去最大の290万件を超え、2017年下半期の約106万件に比べ、前期比2.7倍となった。盗む対象の情報としてクレジットカード情報とApple IDやマイクロソフトアカウントなどの複数のサービスが利用可能なクラウドサービスのアカウントが狙われた。個人を狙ったフィッシング詐欺以外にも法人組織に標的を絞った攻撃も発生している事も述べられている。

1.2 研究目的

そこで本研究では、フィッシング詐欺で利用されるフィッシングサイトのIPアドレスの割り当てられた国、ドメインと情報HTTPレスポンスヘッダの特徴に着目し、フィッシングサイトの情報を提供しているphishtank[4]内のフィッシングサイト100件と株式会社のホームページやウェブサービス等の正規サイト100件、計200件分のIPアドレスの割り当てられた国、ドメインに関する情報とサイト側が設定しているHTTPレスポンスヘッダの情報の収集とデータの分析を行う事でフィッシングサイトの特徴を求め、収集した特徴量によって、フィッシングサイトのドメイン情報とHTTPレスポンスヘッダに関する情報が本物のサイトと似たコンテンツを持ち人手では判別が困難なフィッシングサイトを検知する事において有効なものであるかを確認すると共に、フィッシングサイトを自動検出するシステムを構築する事である。

2. 収集したデータについて

2.1 概要

本研究で収集したフィッシングサイトの主要な特徴量を表1に示す。表1の特徴量を正規サイトとフィッシングサイト100件ずつ計200件求めた。

表1 ウェブサイトの特徴量

| 特徴 | 内容 | 例 |
|-------------------------|---|----------------|
| country | サイトを運用しているIPアドレスの国コード | US |
| domain interval | アクセス日時 - ドメイン作成日時 | 98 |
| domain lifetime | ドメイン期限日 - ドメイン作成日時 | 6 |
| x-xss-protection | HTTPレスポンスヘッダの一種。XSS攻撃を防止するための設定 | TRUE/ FALSE |
| x-frame-options | HTTPレスポンスヘッダの一種。クリックジャッキング攻撃を防止する設定 | TRUE/ FALSE |
| x-content-type-options | HTTPレスポンスヘッダの一種。コンテンツの内容を見ない様にする事でXSS攻撃のリスクを減らす設定 | TRUE/ FALSE |
| content-security-policy | HTTPレスポンスヘッダの一種。UAが読み込みを許可されたリソースを管理出来る様にするための設定 | TRUE/ FALSE |
| FR | 対象のサイトが正規サイトかフィッシングサイトか | TRUE/ FAKE |

実際に上記の特徴量を格納したデータフォーマットを以下に示す。

```
{
  "URL": {
    "country": "US",
    "domain interval": 28,
    "domain lifetime": 1,
    "x-xss-protection": True,
    "x-frame-options": False,
    "x-content-type-options": False,
    "content-security-policy": False,
    "fr": Fake
  }
}
```

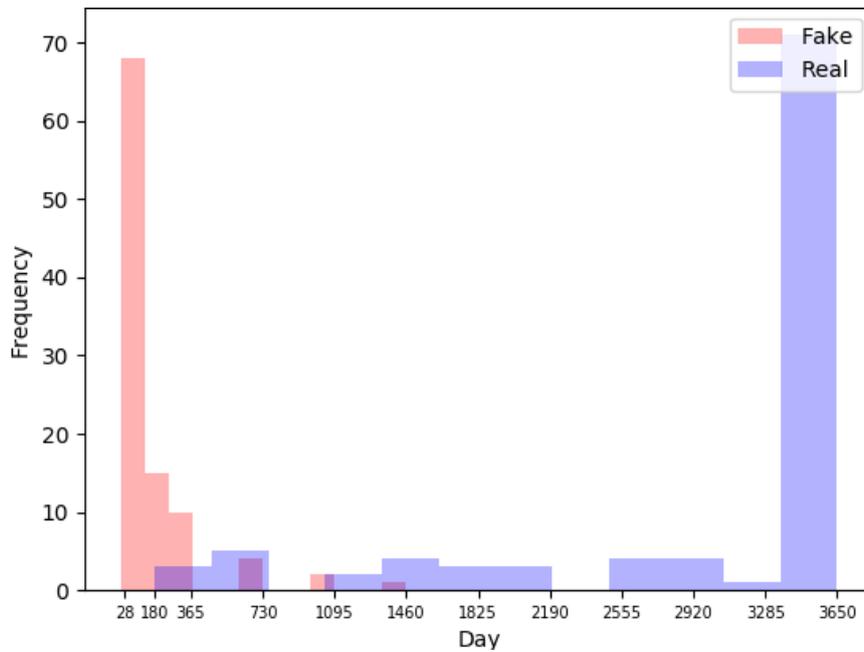
2.2 ドメインに関するデータ

本サイトアクセス時にサイトのドメインが作成されてから、何日間経過しているかを表す”domain interval”の散布図を図1に示す。その中央値・平均値・最頻値を表2に示す。

表2 domain interval

| | 正規サイト[日] | フィッシングサイト[日] |
|-----|----------|--------------|
| 平均値 | 3089 | 135 |
| 中央値 | 3650 | 14 |
| 最頻値 | 3650 | 7 |

図1 domain interval-散布図



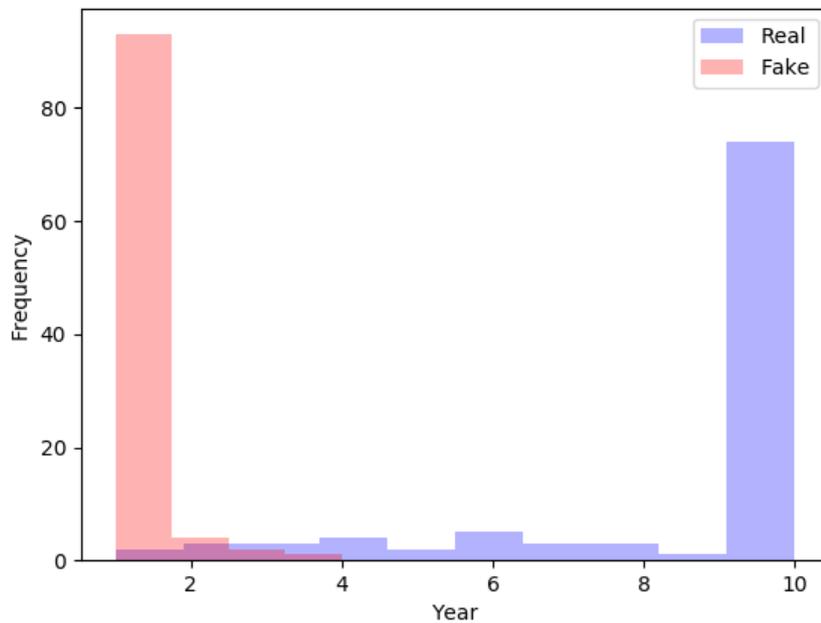
“domain interval”は先行研究[3]を参考にし、区間で分割しており10年以上の場合は10年として扱う。正規サイトの最頻値が3650日であるのに対し、フィッシングサイトは7日と非常に期間が短い。散布図からもフィッシングサイトの多くは7日と非常に短い。

次に、ドメインが作成された日時からドメインの有効期限日までの期間を”domain lifetime”と呼ぶデータの散布図を図2に示す。その中央値・平均値・最頻値を表3に示す。

表3 domain lifetime

| | 正規サイト[年] | フィッシングサイト[年] |
|-----|----------|--------------|
| 平均値 | 8.6 | 1.1 |
| 中央値 | 10 | 1 |
| 最頻値 | 10 | 1 |

図2 domain lifetime-散布図



“domain interval”がドメインの新しさを示す指標であるのに対して，“domain lifetime”はドメインがどの位の期間利用するつもりであるかを知るための指標である。

“domain lifetime”も先行研究[5]を参考にした。区間で分割をしており10年以上の場合は10年として扱う。最頻値と中央値が正規サイトが10年なのに対して、フィッシングサイトは1年と非常に短い。

2.3 国に関するデータ

サイトのIPアドレスの国分布を調べ、表4に示した。国コードはOtherは1件のみの国コードをまとめたものとなっている。

表4 正規サイトの国分布

| 国コード | 正規サイト | フィッシングサイト |
|---------|-------|-----------|
| US | 56 | 40 |
| JP | 40 | 2 |
| IE | 2 | 0 |
| TW | 0 | 20 |
| NL | 0 | 6 |
| UK | 0 | 4 |
| PA | 0 | 4 |
| RU | 0 | 4 |
| CA | 0 | 3 |
| KR | 0 | 2 |
| MA | 0 | 2 |
| FR | 0 | 2 |
| Other | 2 | 9 |
| Unknown | 0 | 2 |

表4よりどちらもアメリカ(US)が最も多い事は同じだが、一方で台湾(TW)、オランダ(NL)、ウクライナ(UK)、パナマ(PA)、ロシア(RU)などがフィッシングサイトのみで使用されている違いがある。フィッシングサイトは正規サイトと比較して様々な国のIPアドレスが利用されている。

トレンドマイクロセキュリティブログの記事[4]によると、オランダとルーマニアの小規模なサーバホスティング事業者が2015年に標的型攻撃などに使用されていた事があり、その他のサイバー犯罪でも利用されていたという事例がある。

台湾のIPアドレスで運用されていたフィッシングサイトはほぼ日本のサービスを狙っているものであり、これは地理的に近い事が影響していると思われる。

これらの事から正規サイトに含まれている国コードのUS,JPを正規サイトの特徴として扱い、それ以外の国コードは評価しない様にした。しかしUSに関してはフィッシングサイトも多く利用していたので、国コードという特徴はあまり重要な特徴として扱ってない。

2.4 HTTPレスポンスヘッダに関するデータ

求めたHTTPレスポンスに関する項目は4つあり、それぞれ正規サイトとフィッシングサイトで各々100件中何件設定されているかを求めたものを表5に示す。

表5 各HTTPレスポンスヘッダの件数

| | 正規サイト | フィッシングサイト |
|-------------------------|-------|-----------|
| x-frame-options | 70 | 4 |
| x-content-type-options | 58 | 4 |
| x-xss-protect | 56 | 4 |
| content-security-policy | 25 | 2 |

表5内のHTTPレスポンスヘッダはユーザ側のセキュリティを考慮して設定されている。フィッシングサイトの短期間の運用、ユーザの情報を盗む様なサイトである点からこれらの設定が正規サイトに比べて設定されていないのではないかと考えた。表5より、明らかに各項目が正規サイトに比べてフィッシングサイトの方が設定されていない事が分かる。これによりフィッシングサイトはセキュリティに関するHTTPレスポンスヘッダの設定が行われているサイトが定量的に少ない事が確認出来る。

正規サイトで上記のレスポンスヘッダは比較的新しいサイトの方が設定されている割合が少し多かった。年数毎の区間で対象のレスポンスヘッダが何件設定されているかを表6に示す。

運用期間が2年以下と4年以下のサイトの結果は”content-security-policy”以外ほぼ設定されている事が確認出来る。

表6 年数ごとの正規サイトの設定件数

| | 2年以下 | 4年以下 | 6年以下 | 7年以上 |
|-------------------------|------|------|------|-------|
| x-frame-options | 5/7 | 6/6 | 9/10 | 50/77 |
| x-content-type-options | 5/7 | 5/6 | 5/10 | 43/77 |
| x-xss-protect | 6/7 | 5/6 | 6/10 | 39/77 |
| content-security-policy | 0/7 | 2/6 | 3/10 | 20/77 |

2.5 トップレベルドメインに関するデータ

実際に特徴量として使用する事は無かったが、正規サイトとフィッシングサイトのTLD(トップレベルドメイン)も特徴として使用する目的で収集を行なった。各TLDが100件中の件数を求めたものを表7に示す。

表7 各TLD毎の件数

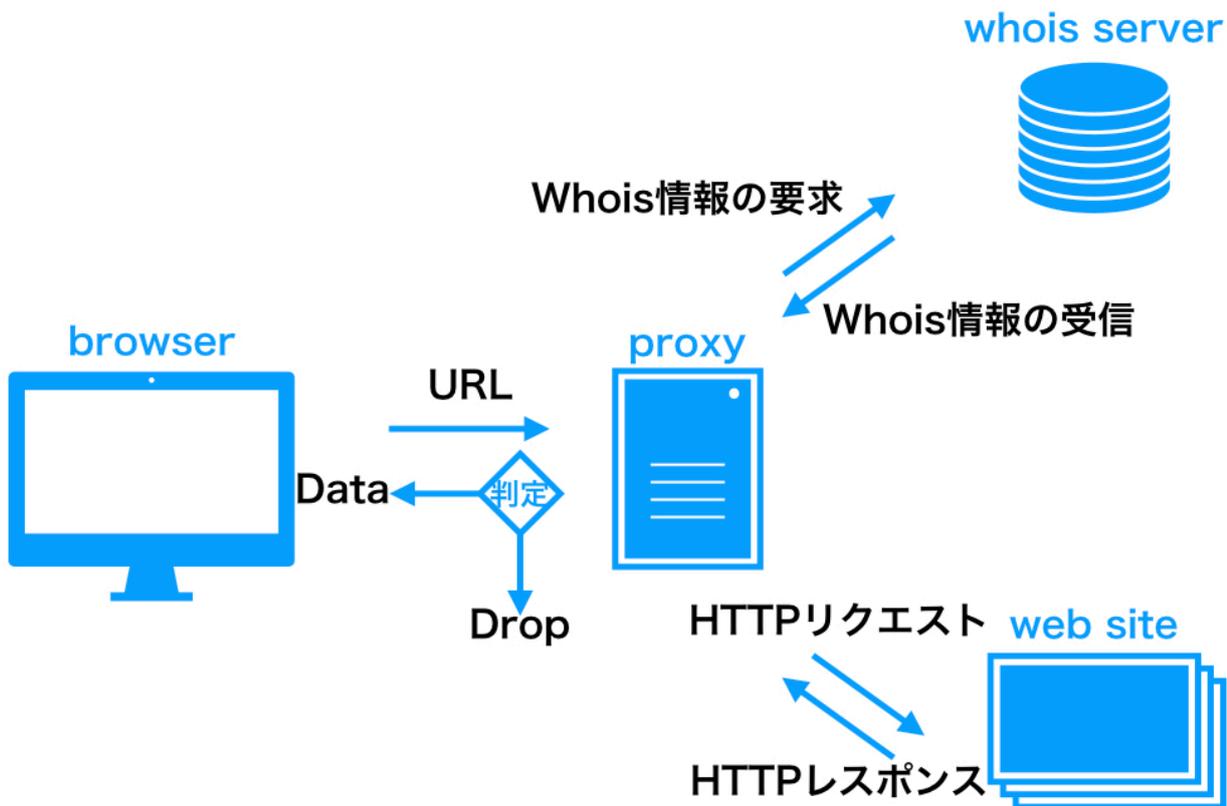
| TLD | 正規サイト | フィッシングサイト |
|---------|-------|-----------|
| com | 58 | 74 |
| jp | 29 | 1 |
| net | 4 | 4 |
| me | 2 | 0 |
| org | 1 | 2 |
| cc | 1 | 1 |
| is | 1 | 0 |
| life | 1 | 0 |
| live | 1 | 0 |
| mu | 1 | 0 |
| tv | 1 | 0 |
| xyz | 0 | 5 |
| info | 0 | 3 |
| am | 0 | 1 |
| cf | 0 | 1 |
| co | 0 | 1 |
| in | 0 | 1 |
| ir | 0 | 1 |
| support | 0 | 1 |
| top | 0 | 1 |
| uk | 0 | 1 |
| us | 0 | 2 |

データを収集する前では、フィッシングサイトでは”.xyz”等のドメインが多く使用されていると予測していたのだが、実際には正規サイトとフィッシングサイトどちらも”.com”が最も多く利用されており、正規サイトも様々なドメインを使用している点から特徴量として利用する事は難しいと判断した。

3. フィッシングサイト検出支援システム

求めた特徴を用いて、対象のサイトが安全であるかを判定するプログラムをpythonで作成した。プログラムの構成図を図3に示す。本システムはプロキシとして動作し、フィッシングサイトを検出した時にアクセスをブロックする。

図3 システム構成図



本システムでサイトが安全か危険かを判断する式を以下に示す。式は先行研究[5]を参考にし、表1の拡張ごとに与えた重みwに対して、特徴量をカテゴリ化して割り当てた定数 d_i を掛けた値の合計値が閾値 θ を超えた場合にサイトを安全と判断した。

$$\sum_{i=0}^7 w_i * d_i > \theta$$

上記の式で利用した各特徴毎の重みと特徴をカテゴリ化して割り当てた定数を表9, 表10, 表11, 表12, 表13に示す.

表9 特徴量ごとの重み

| 特徴量 | 重み |
|-------------------------|------|
| country | 0.4 |
| domain_interval | 0.45 |
| domain_lifetime | 0.45 |
| content-security-policy | 0.33 |
| x-frame-options | 0.33 |
| x-content-type-options | 0.33 |
| x-xss-protection | 0.33 |

表10 "country"の点数

| | 点数 |
|-------|----|
| TRUE | 50 |
| FALSE | 0 |

表11 "domain interval"の点数

| 期間 | 点数 |
|-------|-----|
| 1週間以下 | 0 |
| 1ヵ月以下 | 10 |
| 2年以下 | 50 |
| 4年以下 | 70 |
| 6年以下 | 80 |
| 7年以上 | 100 |

表12 "domain lifetime"の点数

| 期間 | 点数 |
|-----|-----|
| 1年 | 0 |
| 2年 | 10 |
| 3年 | 30 |
| 4年 | 40 |
| 5年 | 40 |
| 6年 | 50 |
| 7年 | 50 |
| 8年 | 70 |
| 9年 | 70 |
| 10年 | 100 |

表13 HTTPレスポンスヘッダの点数

| | 点数 |
|-------|-----|
| TRUE | 100 |
| FALSE | 0 |

表14に判定プログラムで交差検証を行い、計200件中2割を学習データとし、残りの8割をテストデータとした時の結果を示す。FNは正規サイトを危険と判断した件数 TNはフィッシングサイトを安全と判断した件数を意味しており、それぞれ80件中何件なのかを示す。その結果によって求めた全体正解率は約95%となり、有用性があると言える。

表14 FN,TNの平均値

| | 1回目 | 2回目 | 3回目 | 4回目 | 5回目 | 平均値 |
|---------------------|-----|-----|-----|-----|-----|-----|
| FN(正規サイトを危険と判断) | 8 | 7 | 7 | 2 | 5 | 5.8 |
| TN(フィッシングサイトを安全と判断) | 2 | 1 | 2 | 5 | 1 | 2.2 |

4. おわりに

本研究を通じてフィッシングサイトの特徴として、ドメインの生存期間と運用期間が短い点、フィッシングサイトでしか使用されていないIPアドレスの国の割り当て、セキュリティーに関するHTTPレスポンスヘッダが設定されている率が低い事が分かった。特にセキュリティーに関するHTTPレスポンスヘッダを設定しているフィッシングサイトは数件しか存在せず、正規サイトと大幅に異なる特徴となった。また、IPアドレスの割り当て国も正規サイトでは種類が少なく、フィッシングサイトでのみ使用されている国が多数存在する事が判明した。

現在のフィッシングサイト検出システムでは、ドメインの生存期間と運用期間だけでサイトが安全か否かを判定すると新しいウェブサイトを危険と判定してしまう事を防ぐためにHTTPレスポンスヘッダの特徴を加えたが、未だドメインの期間に関する情報は安全性を求める上で重要な特徴となっており、特徴ごとの重みと、特徴量をカテゴリ化して割り当てた定数の値を決める上で困難な点であった。なので、現在の特徴量以外にもウェブサイトが持つコンテンツの内容などの他の特徴量となり得るものに焦点を当て、識別の質を高める事を行う。また、フィッシングサイトを作成するためのフィッシングサイト構築キットが販売されており、被害に拍車を掛けている。今後はこの様なフィッシングサイト構築キットの特徴や攻撃者が送信するフィッシングメールに関する特徴等を調べる事で識別の精度を高める事を目標とする。

参考文献

- [1] フィッシング(詐欺) ([https://ja.wikipedia.org/wiki/フィッシング_\(詐欺\)](https://ja.wikipedia.org/wiki/フィッシング_(詐欺)), 2018年12月参照)
- [2] フィッシング対策協議会, "フィッシングレポート2018"(https://www.antiphishing.jp/report/pdf/phishing_report_2018.pdf, 2018年12月参照)
- [3] トレンドマイクロセキュリティブログ,"「クラウド時代の認証情報」を狙いフィッシング詐欺が急増、2018年上半期の脅威動向を分析"(<https://blog.trendmicro.co.jp/archives/19461>, 2018年12月参照)
- [4] phishtank(<https://www.phishtank.com>, 2018年12月参照)
- [5] 中村元彦,寺田真敏,千葉雄司,土井範久,"プロキシを利用した HTTPリクエスト解析によるフィッシングサイト検出システムの提案" 情報処理学会論文誌 Vol.48 No.10,2007.
- [6] トレンドマイクロセキュリティブログ"サイバー攻撃に利用されやすいオランダのホスティングサービス"(<https://blog.trendmicro.co.jp/archives/13279>,2018年12月参照)
- [7] MDN web docs HTTPヘッダー(<https://developer.mozilla.org/ja/docs/Web/HTTP/Headers>, 2018年12月参照)

謝辞

本研究に際して,様々なご指導をいただきました菊池浩明教授に深く感謝します.最後に,菊池研究室の皆様へ感謝の意を表すると共に,謝辞にかえさせていただきます.

付録A 実際に収集したウェブサイトのデータ

表A.1に収集したフィッシングサイトの各特徴をまとめたデータの一部を示す。

表A.1 フィッシングサイトのデータ

| No | fr | country | domain_interval | domain_lifetime | x-frame-options | x-content-type-options | x-xss-protection | content-security-policy |
|----|------|---------|-----------------|-----------------|-----------------|------------------------|------------------|-------------------------|
| 1 | Fake | US | 1M | 1 | False | False | False | False |
| 2 | Fake | DE | 2Y | 1 | False | False | False | False |
| 3 | Fake | BG | 1W | 1 | False | False | False | False |
| 4 | Fake | US | 2Y | 2 | False | False | False | False |
| 5 | Fake | AU | 2Y | 1 | False | False | False | False |
| 6 | Fake | US | 1M | 1 | False | False | False | False |
| 7 | Fake | IR | 1W | 1 | False | False | False | False |
| 8 | Fake | KR | 1W | 1 | False | False | False | False |
| 9 | Fake | US | 1W | 1 | False | False | False | False |
| 10 | Fake | US | 1W | 1 | False | False | False | False |
| 11 | Fake | CA | 2Y | 1 | False | False | False | False |
| 12 | Fake | US | 2Y | 1 | False | False | False | False |
| 13 | Fake | US | 2Y | 1 | False | False | False | False |
| 14 | Fake | NL | 1M | 1 | False | False | False | False |
| 15 | Fake | Unknown | 2Y | 1 | False | False | False | False |
| 16 | Fake | RO | 2Y | 1 | False | False | False | False |
| 17 | Fake | US | 1M | 1 | False | False | False | False |
| 18 | Fake | RU | 1M | 1 | False | False | False | False |
| 19 | Fake | TW | 1M | 1 | False | False | False | False |
| 20 | Fake | UK | 1W | 1 | False | False | False | False |

表A.2に収集した正規サイトの各特徴をまとめたデータの一部を示す。

表A.2 正規サイトのデータ

| No | fr | country | domain_interval | domain_lifetime | x-frame-options | x-content-type-options | x-xss-protection | content-security-policy |
|----|------|---------|-----------------|-----------------|-----------------|------------------------|------------------|-------------------------|
| 1 | Real | US | 7Y | 10 | True | True | True | True |
| 2 | Real | US | 6Y | 10 | True | True | True | True |
| 3 | Real | US | 7Y | 10 | True | True | True | False |
| 4 | Real | IE | 7Y | 10 | True | True | False | True |
| 5 | Real | US | 7Y | 10 | True | True | True | True |
| 6 | Real | JP | 7Y | 10 | True | True | True | False |
| 7 | Real | JP | 7Y | 10 | True | False | True | False |
| 8 | Real | US | 7Y | 10 | True | True | True | True |
| 9 | Real | US | 7Y | 10 | True | False | True | False |
| 10 | Real | US | 7Y | 10 | True | True | True | True |
| 11 | Real | US | 7Y | 10 | True | True | True | False |
| 12 | Real | US | 7Y | 10 | False | True | True | False |
| 13 | Real | US | 7Y | 10 | True | True | False | True |
| 14 | Real | US | 7Y | 10 | False | True | False | True |
| 15 | Real | US | 7Y | 10 | True | True | True | True |
| 16 | Real | US | 7Y | 10 | True | False | True | False |
| 17 | Real | US | 7Y | 10 | True | True | True | False |
| 18 | Real | US | 7Y | 10 | True | False | False | True |
| 19 | Real | IE | 7Y | 10 | True | True | True | False |
| 20 | Real | US | 7Y | 10 | True | True | False | False |

付録B FFRIデータセットのマルウェアのAPI情報による識別実験

B.1 はじめに

B.1.1 研究背景

近年,マルウェアによる被害は増加の一途を辿っている.マルウェアの種類も様々な物が存在し,被害を防ぐために不明瞭なファイルはダウンロードしない等の注意は重要ではあるが,ユーザが気付かない間にマルウェアに感染している事例は日々起こっている.

その様な被害を防ぐためにアンチウイルスソフトはマルウェアの種類を捉え,スキャン対象であるファイルがマルウェアであるか否かの判定の足がかりにしている.そこでマルウェアの特徴を知る事が出来ればさらなる感染防止に貢献出来ると思ひ,本研究ではセキュリティ・リサーチ会社のFFRIが自社で収集したマルウェアの検体の動的解析ログであるFFRIのデータセットを用いてマルウェアの特徴を探る.

本研究の目的は,マルウェアが用いるAPIの情報がマルウェアの識別において有効なものであるかを確認すると共に,マルウェアが使用するAPIにおいて特徴を掴む事である.

B.1.2 研究方法

本研究では,FFRIのデータセットのデータをサポートベクターマシン(以下SVM)を用いて行っている.SVMとは教師あり機会学習を用いるパターン認識モデルの1つであり,線形入力素子を利用して2クラスのパターン識別器を構成する手法である.SVMにFFRIのデータセットのAPIに関する情報を掛けることで,マルウェアの識別を行う.

B.2 収集したデータについて

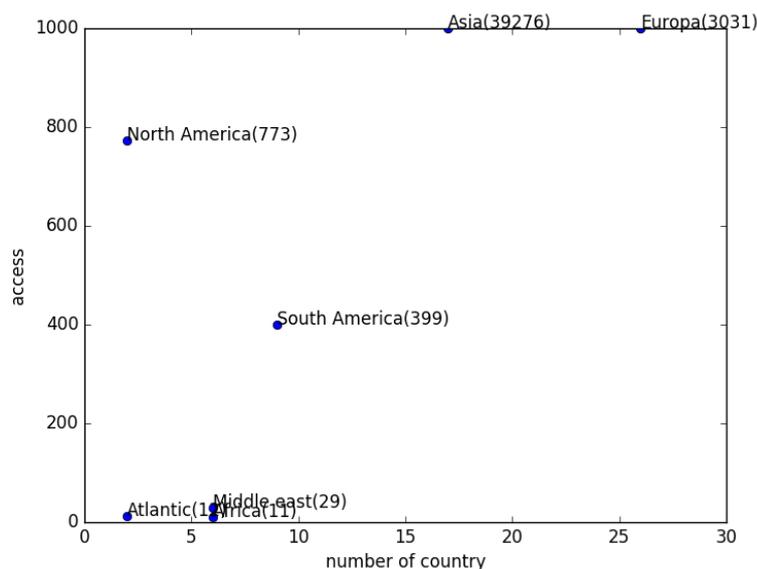
B.2.1 概要

実際にFFRIのデータセットを分析する前に理解を深めるため,SSHハニーポットcowrieを運用し,1ヶ月分のログを収集,分析した.cowrieで集めたログの地域/国毎に関するデータ, パスワードに関するデータ, ユーザ名に関するデータをpythonを用いて求め,グラフして可視化した.

B.2.2 エリア毎のデータ

エリア毎に飛んできたアクセスを分類し,その結果をプロットした散布図を図B.1に示した.アジアからのアクセスが最も多く,逆にアフリカからのアクセス数は少ないという結果だった.

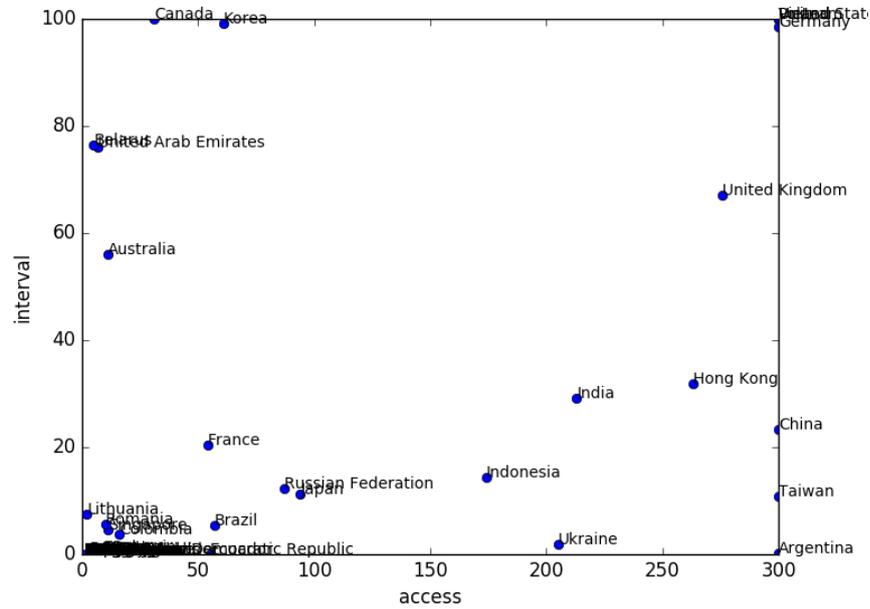
図B.1 エリア毎のデータの散布図



B.2.3 国毎のデータ

エリア毎のデータだと少し大雑把なので,国ごとのデータを求めた結果をプロットした散布図を図B.2に示す.結果としてはベトナムが最も多く(19896件),続いて台湾(17310件),中国(1157件)となった.

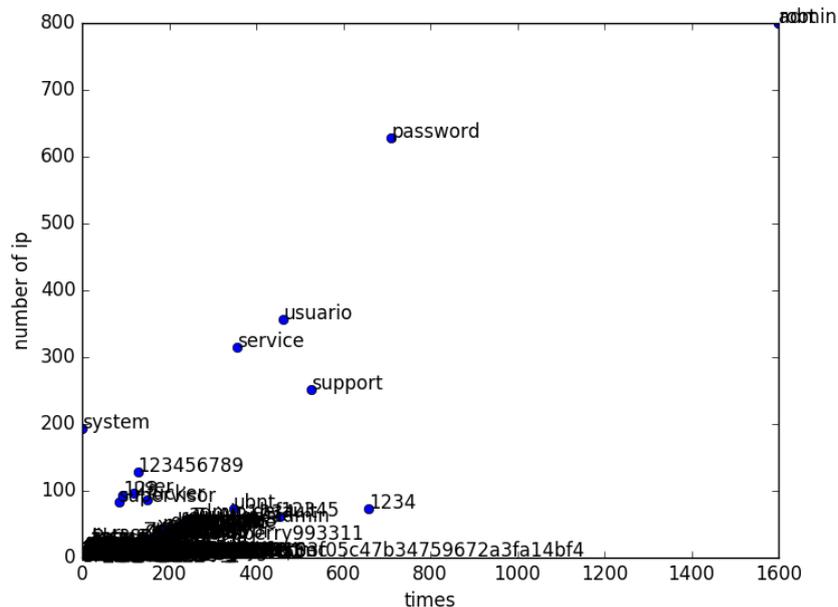
図B.2 国毎のデータの散布図



B.2.4 パスワードのデータ

どの位の端末があるパスワードを何回入力しているのかを知り,どのパスワードが狙われやすいのかを求めた結果をプロットした散布図を図B.3に示す.結果としてはrootが最も多く(13099回),次点でadmin(7331回/1194個),他にも”password”, ”usuario”, ”support”, ”service”が多かった.

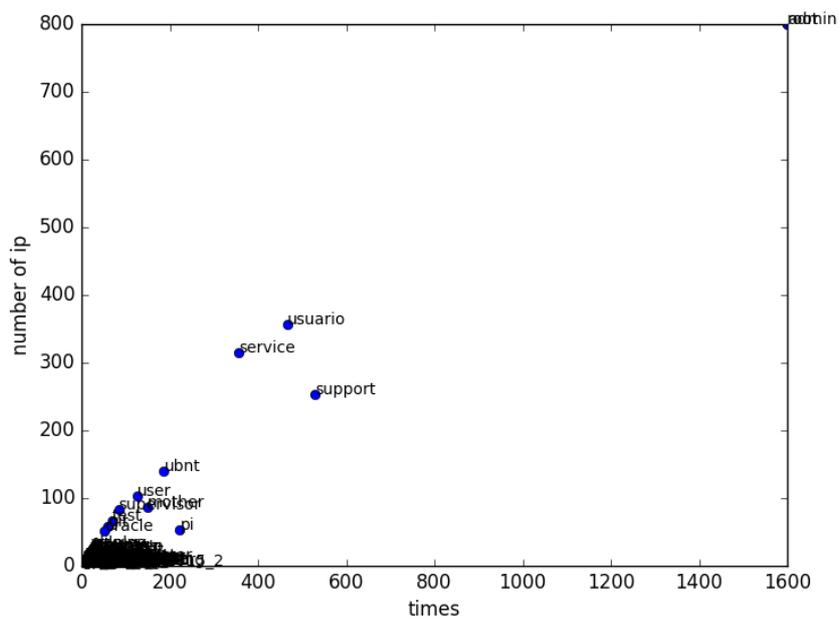
図B.3 パスワード毎のデータの散布図



B.2.5 ユーザ名のデータ

パスワードと同じ様にどのユーザ名が狙われやすいのかを求めた結果をプロットした散布図を図B.4に示す。結果としてパスワードと同じくrootが最も多く(17795回),次点でadmin(10570回)だった。他にも”usuario”, “support”, “service”, “ubnt”, “user” などが多かった。

図B.4 ユーザ名のデータ



B.3 FFRIのデータについて

B.3.1 概要

FFRI DataSet 2017は一つの検体を動的解析した時の様々な情報が記録されている.具体的なファイルの内容は表A.1にまとめた.本研究では検体が何処にデータを送信しているのか(network項目)とAPIの情報(behavior項目)を使用した.また,今回使用する検体ごとに何個あるのかを表A.2にまとめた.

表A.1 FFRI Dataset 項目

| 項目 | 内容 |
|------------|---------------------|
| info | 解析の開始,終了時刻,id等 |
| signatures | ユーザ定義シグニチャとの照合結果 |
| virustotal | VirusTotalから得られる情報 |
| static | 検体のファイル情報 |
| dropped | 検体の実行時に生成したファイル |
| behavior | 検体実行時のAPIログ |
| traget | 解析対象検体のファイル情報 |
| debug | 検体解析時のcuckooのデバッグログ |
| strings | 検体中に含まれる文字列情報 |
| network | 検体の実行時に行った通信の情報 |

表A.2 検体ごとの個数

| 項目名 | 説明 | 検体数 |
|--------|--------------------------------------|-----|
| TROJ | トロイの木馬型不正プログラム | 572 |
| RANSOM | 身代金要求不正プログラム | 385 |
| BKDR | 主にバックドア活動を行うTROJ | 225 |
| TSPY | 情報漏洩に繋がる行動を行うTROJ | 124 |
| PE | PE形式の実行ファイルに感染するファイル感染型ウイルス | 105 |
| Mal | ジェネリック検索で検出された不正な疑いのあるファイル | 44 |
| Worm | 主にワーム活動を行う不正プログラム | 40 |
| PUA | ユーザのセキュリティ,プライバシーに予期しない影響を与える可能性がある物 | 16 |

B.3.2 送信先の国について

送信先の国を求めた結果を表A.3にまとめた。上位3カ国はフランス、ベルギー、韓国となっている。結果を見てみると、以外な事に全体的にヨーロッパ諸国の国が多く、上位10カ国中半分がヨーロッパ及びその周辺だった。

表A.3 送信先の国上位10カ国

| 国名 | 回数 |
|-------------------|--------|
| France | 463908 |
| Belgium | 127200 |
| Republic of Korea | 83087 |
| United States | 78257 |
| Norway | 30422 |
| Germany | 30354 |
| India | 29207 |
| Greece | 27904 |
| Russia | 25197 |
| Singapore | 23539 |

B.3.3 API

昨今、マルウェアを検知する技術として機会学習等が注目されており本研究においてもマルウェアが使用しているAPIを特徴量としてサポートベクターマシン(以下SVM)を用いて、検体を識別出来るかを調べた。マルウェアのラベル名はVirusTotal項目内のTrendMicroの結果を用いており、TrendMicroによって判断出来た検体数は1742件である。

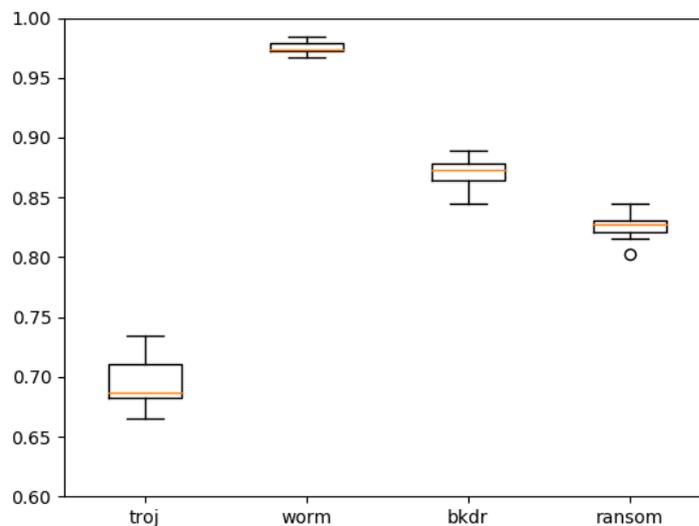
はじめに、どのAPIが全体的に使用されているかを求め、上位10種類を表A.4にまとめた。この事から全検体がIAT(Import Address table)に対するフックや書き換えのためにこれらの関数を使用している可能性が高いという事が推測できる。そしてSVMの結果を箱ひげ図で表している図A.1を見てみるとTROJは7割ほど、TORJ以外は正答率が8~9割ほどで良い結果を出している。これにより使用しているAPIは識別するのに使用する特徴量として良い結果を出したという事が分かった。

表A.4 使用されているAPI上位10種

| API名 | 回数 |
|------------|------|
| NtClose | 1742 |
| LdrLoadDll | 1742 |

| API名 | 回数 |
|-------------------------|------|
| NtAllocateVirtualMemory | 1742 |
| LdrGetProcedureAddress | 1742 |
| NtFreeVirtualMemory | 1739 |
| RegCloseKey | 1738 |
| NtCreateFile | 1737 |
| NtOpenKey | 1736 |
| LdrGetDllHandle | 1730 |
| NtQueryValueKey | 1698 |

図A.1 SVMの結果



B.4 まとめ

マルウェアのデータ送信先国を求め、何処にデータが多く送信されているかを求めたところ、ヨーロッパ諸国が多い事が分かった。また、マルウェアが使用しているAPIによる分類を行い、その有用性を実証した。FFRIデータセットにはAPI以外の情報もあるので、今後はそれらの情報を使用してマルウェアの分類をする事に取り組んで行きたいと思う。

参考文献

- [1] 鈴木貴之,宮保憲治 “データマイニング技術を活用したマルウェア分類法の検討”, 情報処理学会第76回全国大会
- [2] FFRI Dataset 2017のご紹介 (http://www.iwsec.org/mws/2017/20170606/FFRI_Dataset_2017.pdf)
- [3] TrendMicro - ウイルス名の接頭辞について (<https://success.trendmicro.com/jp/solution/1306448>)
Wikipedia - サポートベクターマシン (<https://ja.wikipedia.org/wiki/%E3%82%B5%E3%83%9D%E3%83%BC%E3%83%88%E3%83%99%E3%82%AF%E3%82%BF%E3%83%BC%E3%83%9E%E3%82%B7%E3%83%B3>)
- [4] %E3%82%B5%E3%83%9D%E3%83%BC%E3%83%88%E3%83%99%E3%82%AF%E3%82%BF%E3%83%BC%E3%83%9E%E3%82%B7%E3%83%B3