

研究指導 中澤 真 准教授

# 過去の受信メールに基づく 訓練用標的型攻撃メールの自動生成に関する一考察

伊藤 裕史

## 1. はじめに

近年、インターネットの普及によって人々の生活は便利で豊かなものになったが、そのことを悪用して企業や官公庁にサイバー攻撃を仕掛けるといった事例が増加している[1][2][3]。そのうえ、最近では中小企業や個人もターゲットとなっているため、法人だけではなく個人単位でも対策をする必要がある[4][5]。しかし、対策をしていてもサイバー攻撃による被害は増加しており、その影響も大きくなっている。

なかでも標的型攻撃メールは特に問題になっている。従来のコンピュータウイルスは不特定多数を対象にばらまく形であり、ウイルスが添付されているメールの文章は誰が受け取るものも同じものであったため、受信者もメールの怪しさに気づきやすかった。これに対し標的型攻撃メールは受信者に合わせて内容を変化させて送信するため非常にだまされやすいといえる。攻撃への対策は受信者に対して訓練用の標的型攻撃メールを送り、攻撃を疑似体験させる方法が一般的である[6]。訓練の実施にあたって1つ問題がある。それは訓練用の標的型攻撃メールの文面にテンプレートを利用することが多い点である。これではメールの内容を変化させる標的型攻撃メールの訓練にならない。そのため一人ひとりに適したメールを用いて訓練する必要がある。

そこで本研究では過去の受信メールの名詞、動詞、形容詞の単語の頻度から、開封のしやすさに影響する特徴的な単語を抽出し、これらの単語を用いてその利用者がだまされやすい傾向のメールを選び出す方法を示す。この結果を応用することで、利用者一人ひとりにパーソナライズされた訓練メールの生成につなげることができるようになる。

## 2. 近年のサイバー攻撃

### 2.1 インターネットを悪用したサイバー攻撃

パソコンやスマートフォンの爆発的な普及により、誰もが簡単にインターネットを利用するようになった。しかし、それを悪用してマルウェアを用いたDDos攻撃[7][8]や特定の組織や団体に送られる標的型攻撃メール[1][2]などのサイバー攻撃が後を絶たない。標的型攻撃メールとは実在する発信元に偽装し、関係者を装ったメールにすることで添付ファイルやURLを開かせウイルスに感染させる攻撃である。その標的型攻撃メールの開封率は、2012年の22.5%から2016年の9.2%と低下を続けているが[9]、それでもなお9人に1人は添付ファイルや

URLを開いてしまっており、依然として感染リスクの高い攻撃であることは間違いない。実際に2016年には大手旅行会社JTBにおいて取引先になりすましたメールの添付ファイルを開いてしまいウイルスに感染し顧客情報が流出した[2]。さらに富士山大学では非常勤講師が標的にされ、添付ファイルを開いてしまったため個人情報や研究データが外部に送信された事例もある[10]。

### 2.2 標的型攻撃メール対策

標的型攻撃メールの対策としては、訓練用の標的型攻撃メールを利用者に送り、攻撃を疑似体験させる方法が一般的である。しかし訓練メールの文面はテンプレートのみになっているためバリエーションが限られ、回数をこなすと訓練メールと容易に気づかれるようになり意味が薄れてしまう恐れがある。しかし手作業でメール文を作成するのは過大な労力を必要とするため現実的ではない。

この問題に対して岩田ら[11]では自動で継続的に訓練を行うメールクライアントシステムを提案している。しかし、メールを解析し自動生成する機能については、具体的な解析方法までには踏み込んでいなかった。

そこで本研究では訓練用標的型攻撃メールを自動生成するための第一歩として、過去の受信メールの中から最も訓練メールにふさわしいもの、すなわち添付ファイルやURLを開きやすいメールを選び出す方法を提案する。具体的には、添付ファイルやURLを開きやすいメールとそうでないメールでは、件名や本文に現れる単語の出現傾向に違いがあると考え、特徴的な単語の出現頻度に基づき二つのメール群を見分ける評価値を算出する方法を示す。

## 3. 過去の受信メールの分析

### 3.1 受信メールの取得

今回は自身のGmailに保存されている2016年4月5日～2016年12月12日のメール1010件を取得した<sup>1</sup>。取得したメールはmbox形式であるため、MozillaThunderbird<sup>2</sup>を使いすべてのメールデータをテキスト形式に変換した。

分析の目的は添付ファイルやURLを開きやすいメールと、開きにくいメールの傾向を明らかにすることであるため、メールに添付ファイルが無いあるいはURLのリンクが本文に含まれていないメールは分析の対象とはしない。そこで、先の取得したメールの中からこの条件を満たさないものを除外し、最終的に

<sup>1</sup> Google のデータツールへアクセスし作成したアーカイブをダウンロードした

<sup>2</sup> <https://www.mozilla.org/ja/thunderbird/>

460件のメールを対象に分析を行う。なお、対象メールのうち、開いた可能性があるメールは122件、開かなかったメールは338件であった。

### 3.2 RMeCabを用いたメール分析

開封・未開封メールの品詞の単語の特徴を明らかにするために件名と本文に対して形態素解析を行った。解析をする際に用いたソフトはRMeCab[12][13][14]である。まず添付・URLの開封の有無を基準にそれぞれの品詞の割合を図1に示す。この結果から開封メールと未開封メールの品詞の割合に特徴はなく、品詞の構成割合だけでは開封しやういメールとそうでないメールを見極めることができないといえる。

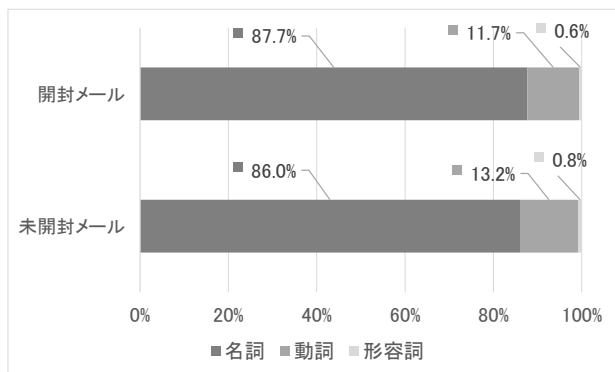


図1:メール別の品詞割合

次に、件名や本文に現れる単語の出現傾向が、開封メールと未開封メールで差異があるのではないかと考え、単語の出現頻度を分析した。分析にあたっては名詞、動詞、形容詞の品詞ごとに出現率の高い単語を開封メール、未開封メールでそれぞれ取り出した。ここで出現率は、各単語の出現頻度をすべての単語の出現頻度総数で除した値としている。

まず名詞についての結果を表1、表2に示す。名詞では細分類で記号、数字類などに属する意味がない単語を除いて出現率の上位20個の単語を抽出した。表1と表2を比べると、どちらの表にも登場する高い出現率の単語がある。例えば本学のアドレスのドメインである「jc.u-aizu.ac.jp」や学科名の一部である「情報」などの単語が挙げられる。開封メールにも未開封メールにも共通しているこれらの単語は、それぞれの特有の単語ではない、つまり特徴語ではないと考えるのが妥当である。そこで本研究では、ある単語の出現頻度の順位が、開封メールと未開封メールでその差が30以上あれば、出現に偏りがある単語として特徴語として定義した。表1、表2にはその順位差も記し、順位差が基準より大きいものを網掛けによる強調表示にした。

表1から開封メールの特徴語としては、受信者の知り合いの名前や、所属するコミュニティ特有のキーワード「課題」「提出」などが現れた。一方、表2の結果からは未開封メールの特徴語として「事務」「お知らせ」などメールマガジンや事務連絡に

関係する単語に多い傾向が示された。

表1:開封メールの出現頻度上位の名詞

開封順位	未開封順位	順位差	単語	出現頻度	出現率
1	1	0	jc.u-aizu.ac.jp	657	6.36%
2	3	1	情報	179	1.73%
3	10	7	中澤	139	1.35%
4	74	70	課題	135	1.31%
5	20	15	経営	127	1.23%
6	7	1	ゼミ	115	1.11%
7	37	30	提出	108	1.05%
8	8	0	お願い	106	1.03%
9	34	25	コス	105	1.02%
10	399	389	新聞	94	0.91%
11	2	-9	学生	84	0.81%
12	160	148	報告	80	0.77%
13	375	362	下田	80	0.77%
14	259	245	遠藤	77	0.75%
15	201	186	美香子	75	0.73%
16	198	182	添付	71	0.69%
17	451	434	甫	71	0.69%
18	568	550	昴	71	0.69%
19	5	-14	https	71	0.69%
20	449	429	佳	69	0.67%

表2:未開封メールの出現頻度上位の名詞

未開封順位	開封順位	順位差	単語	出現頻度	出現率
1	1	0	jc.u-aizu.ac.jp	1149	1.10%
2	11	-9	学生	1090	1.04%
3	2	1	情報	955	0.91%
4	37	-33	月	858	0.82%
5	19	-14	https	800	0.76%
6	31	-25	メル	629	0.60%
7	6	1	ゼミ	623	0.60%
8	8	0	お願い	612	0.59%
9	40	-31	企業	612	0.59%
10	3	7	中澤	594	0.57%
11	62	-51	job.rikunabi.com	552	0.53%
12	233	-221	事務	452	0.43%
13	24	-11	伊藤	439	0.42%
14	101	-87	リクナビ	423	0.40%
15	109	-94	参加	420	0.40%
16	33	-17	学科	418	0.40%
17	108	-91	お知らせ	415	0.40%
18	48	-30	連絡	414	0.40%
19	50	-31	時間	405	0.39%
20	5	15	経営	398	0.38%

名詞と同様に動詞と形容詞についても出現率上位の単語を取り出し、開封メールと未開封メールの特徴語の抽出を試みた。この結果を表3～6に示す。なお、表6にある”該当なし”は開封メールにその単語が出現しなかったことを意味する。このため、この単語は未開封メールの特徴語と判断している。

表3:開封メールの出現頻度上位の動詞

開封順位	未開封順位	順位差	単語	出現頻度	出現率
1	1	0	する	657	39.32%
2	2	0	ある	95	5.69%
3	3	0	なる	84	5.03%
4	4	0	できる	49	2.93%
5	5	0	行う	35	2.09%
6	6	0	思う	33	1.97%
7	12	5	書く	21	1.26%
8	14	6	送る	21	1.26%
9	13	4	読む	19	1.14%
10	119	109	よむ	19	1.14%



表10:件名に重み付けした評価値による上位10件のメール

No.	件名	本文	課題	提出	新聞	報告	添付	月	企業	job	事務	リクナビ	参加	お知らせ	連絡	時間	評価値
1	必ず読む★2年生:cc	0	0	0	19	0	-1	0	0	0	0	0	0	0	-1	0	17
2	Re:【重要】経営情報	14	0	4	0	1	-3	0	0	0	0	0	0	0	0	0	16
3	【重要】夏休みの過ごし方	13	0	3	0	1	-2	0	0	0	0	0	0	0	0	0	15
4	夏休みの過ごし方	7	5	3	0	0	0	0	0	0	0	0	0	0	0	0	15
5	【課題】経営情報	8	7	0	0	0	-2	0	0	-1	0	0	0	0	0	0	12
6	第7回新聞中澤先生	4	5	3	0	0	0	0	0	0	0	0	0	0	0	0	12
7	第7回新聞中澤先生	4	3	4	0	0	0	0	0	0	0	0	0	0	0	0	11
8	第4回新聞中澤先生	4	3	3	0	1	0	0	0	0	0	0	0	0	0	0	11
9	第2回新聞中澤先生	3	4	3	0	1	-1	0	0	0	0	0	0	0	0	0	10
10	第9回新聞中澤先生	3	3	3	0	1	0	0	0	0	0	0	0	0	0	0	10

実際に図2と図3のメールを見ると、図2の評価値が上位のメールには所属するコミュニティに関係する内容が書いてあり確かに開きやすいメールだった。また評価値が下位のメールにはメールマガジンなどの誰にでも関係するような内容になっており、開封することは絶対にはないものであった。

<p><b>必ず読む★法律で定められた…</b></p> <p>2年生: (cc:1年生) 進路活動の「報告」をしていない 2年生が大変多い ようです。 以下を読み、 直ちに対応…</p>	<p><b>Re:【重要】夏休み課題について</b></p> <p>経営情報コース 1年生各位 青木です。 昨日、送りました 「新聞課題」の 要項に誤りがあり ましたので 下記のとおり…</p>	<p><b>【重要】夏休み課題について</b></p> <p>経営情報コース 1年生各位 教務厚生委員の 青木孝弘です。 8/5～9/29 までの 夏休み期間中 全ゼミ共通に取り 組む3つの…</p>
--	---	--

図2:評価値上位のメール

<p><b>【東京、大阪、福岡など…】全…</b></p> <p>伊藤裕史さん こんばんは、リクナビ編集部です。 日本全国で開催する【交通費あり】の インターン…</p>	<p><b>【伊藤さんへ】インターンシップ…</b></p> <p>12月まで開催中！ インターンシップLIVE いよいよ本格的に 寒くなってきましたが、 いかがお過ごし でしょうか。風邪を…</p>	<p><b>【福島県で開催】人気のインタ…</b></p> <p>伊藤裕史さん リクナビ編集部です。 本日は、伊藤さんがお住まいの 「福島県」で開催 予定の人気インターン シップ…</p>
---	--	--

図3:評価値下位のメール

表11:開封・未開封メールの平均評価値

開封メールの平均評価値	未開封メールの平均評価値
2.20	-7.63

#### 4. 研究のまとめ

本研究では、過去の受信メールから開封・未開封メールを品詞ごとの単語頻度で分析し、開封しやすいメールの特徴を明らかにした。この特徴語に基づき開封のしやすさの尺度となる評価値の算出方法を示したことにより、評価値の高いメールを選ぶことで開封しやすいメール、すなわち訓練用標的型攻撃メールのベースとなるものを選び出すことができたようになった。これにより、テンプレートに頼らず利用者一人ひとりに特化したメールを自動生成へとつなげることが可能になる。

今回は一人のメールによる分析であったため、分析結果の信頼性について課題がある。今後は複数の利用者のメールを対象にした分析に取り組み、この手法の効果を検証する必要がある。さらに提案した手法に基づき、選択された開封しやすいメールをアレンジして実際の訓練メールの自動生成につな

げる必要がある。

#### 参考文献

- [1] 経団連事務局コンピュータマルウェア感染, 2016/11/15, <http://www.keidanren.or.jp/announce/2016/1115.html>
- [2] 不正アクセスによる個人譲歩流出の可能性について-現状報告と再発防止策-, 2016/8/24, <https://www.jtbcorp.jp/jp/160824.html>
- [3] サイバー攻撃ランサムウェア日立の感染原因 独頭微鏡装置か, 2017/7/4, <https://mainichi.jp/articles/20170704/ddm/008/020/060000c>
- [4] 【重要】pixivの一部アカウントに対する「なりすましログイン」の報告とパスワード変更のお願い, 2016/12/2, <https://www.pixiv.net/info.php?id=3897>
- [5] 都立動物園・水族園のホームページへの不正アクセスについて(続報), 2016/7/8, <http://www.metro.tokyo.jp/INET/OSHIRASE/2016/07/20q78600.html>
- [6] 標的型攻撃メール訓練の目的と利用～効果を上げる方法～, <https://www.ipa.go.jp/files/000053336.pdf>
- [7] JICAなどのサイトが閲覧しづらい状態-またアノニマスのDDos攻撃か, 2016/2/19, <http://www.atmarkit.co.jp/ait/articles/1602/19/news078.html>
- [8] 毎秒1テラビットという史上空前のDDos攻撃が発生, 攻撃元はハッキングされた14万5000台ものウェブカメラ, 2016/9/29, <http://gigazine.net/news/20160929-record-breaking-ddos/>
- [9] サイバーセキュリティ傾向分析レポート2017 “～気軽なIT利用～”, 2017/7/26, <https://www.nri.com/jp/event/mediaforum/2017/pdf/forum257.pdf>
- [10] 富山大学水素同位体科学研究センターに対する標的型サイバー攻撃について(概要), 2016/6/14, <https://www.u-toyama.ac.jp/news/2016/doc/1011.pdf>
- [11] 岩田一希, 中村嘉隆, 高橋修, “標的型メール攻撃対策のための自動訓練メールクライアントシステム”, 情報処理学会第78回全国大会, 6V-06, No3, pp.561-562, 2016, <http://naka.jc.u-aizu.ac.jp/IPSJ2016/data/pdf/6V-06.pdf>
- [12] RとLinuxと…, <http://rmecab.jp/wiki/index.php?FrontPage>
- [13] 小林雄一郎, Rによるやさしいテキストマイニング, オーム社, 2017,
- [14] 石田基弘, Rによるテキストマイニング入門, 森北出版, 2017,
- [15] 標的型攻撃への対策, 2013, [http://www.soumu.go.jp/main\\_sosiki/joho\\_tsusin/security/business/staff/05.html](http://www.soumu.go.jp/main_sosiki/joho_tsusin/security/business/staff/05.html)
- [16] 情報セキュリティ10大脅威, 2017, 2017/3/30, <https://www.ipa.go.jp/files/000058504.pdf>
- [17] IPA情報処理推進機構, 情報セキュリティ, 2017/5/30, <https://www.ipa.go.jp/security/vuln/10threats2017.html>