

企業プレスリリースからのサイバーインシデント情報の自動収集と分析

池上 和輝† 山田 道洋† 菊池 浩明† 乾 孝治†

明治大学総合数理学部†

1 はじめに

Japan Network Security Association(JNSA)は、インターネットニュースやプレスリリースからインシデント情報を収集し発生企業・被害人数などの情報を人手で分類して、毎年報告を行っている [1]. しかし、我々が朝日新聞の記事データベース「聞蔵 II」 [2] を調査したところ、JNSA の報告にはない 2015 年のインシデントが 134 件報道されていた。情報収集に利用するメディアなどにより、報道されるインシデントに偏りがあることを示している。

そこで、本研究ではサイバーインシデントを、メディアなどによる偏りなく網羅的に自動で収集して分類することを目的とする。そのために、企業ウェブサイトからインシデント情報を自動で収集・分類するシステムを開発する。本システムにより、収集したインシデント情報の統計情報と分類精度の評価、各種セキュリティマネージメントの効果を検証する応用が可能であることを示す。

2 クローラー・自動分類システム

2.1 クローラー開発

システムの全体構成を図 1 に示す。1. 企業ウェブサイトの URL 与える。2. ウェブサイトの html とそのページ内のリンクを収集する。3. 収集したコンテンツをテキストに変換する。4. 特定のキーワードを含むテキストを保存する。

キーワードには、[3] の調査結果に基づき情報漏洩に関するプレスリリースに含まれるであろう単語、「お詫び」、「漏えい」、「漏洩」を使用した。2 ではリンク先ページについても任意の回数繰り返し、URL を収集した各企業について 1-4 を繰り返し行う。

日付、被害人数は正規表現により抽出する。漏洩原因は、次の手順で推定する。1. 特定の単語の出現頻度とその単語を含む文書の出現頻度から定めた tf-idf 値を用いて、各漏洩原因の文章の特徴語を抽出する。2. 1 で抽出した各漏洩原因の特徴語から成る 49 次元の特徴語ベクトルを作成する。3. 未知原因のプレスリリースを参照し、2 で求めた特徴語の TF 値を求める。4. プレスリリースから抽出した特徴語ベクトルと各漏洩原因ごとの

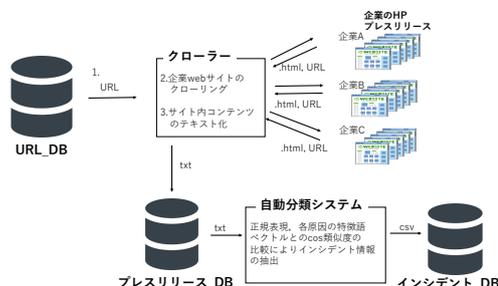


図 1 システム構成図

表 1 クローラーのインシデント収集実験結果

期間	企業数	取得記事数	インシデント記事数	割合
2004/10/1 - 2018/11/2	537	17,957	191	(0.01)

表 2 推定結果

	日付	被害人数	漏洩原因	日付&規模&原因
適合率	0.882	0.792	0.719	0.505
	157/178	141/178	128/178	90/178

特徴語ベクトルの cos 類似度を求める。最も高かった漏洩原因を出力する。

2.2 調査対象企業

株式会社東洋経済新報社は、上場企業全社および主要未上場企業に調査票を送付し、その回答から社会的責任投資 CSR データベース [4] を作成している。本研究では各種セキュリティマネージメントの効果を検証するため、CSR データベース内の 537 企業を調査対象とした。537 企業のうち ISMS を導入しているのは 169 社 (31%), CIO を設置しているのは 150 社 (28%) である。

2.3 実験結果

表 1 に実験結果を示す。2018 年 11 月、2.2 節の企業ウェブサイトの本システムを用いてインシデント収集・分類を行なった。取得記事のうちインシデントに関する記事は 191 件であり、全体の 1% であった。また、191 件にはインシデント詳細ページとインシデント一覧ページの重複があったため、一意なインシデントは 179 件である。

表 2 にインシデントの日付、被害人数、漏洩原因の推定精度を示す。適合率は抽出できた全インシデントのうち、正しく要素を抽出できたインシデントの割合を表している。日付、人数などは全て 7 割を超えるが、全てを組み合わせた時の適合率は 50% であった。

Analysis and Development of Cyber Incident Information Crawler from Public Press Release Statements

†Kazuki Ikegami, Hiroaki Kikuchi, Department of Frontier Media Science, School of Interdisciplinary Mathematical Sciences, Meiji University, Kikuchi Laboratory

†Michihiro Yamada Graduate School of Advanced Mathematical Sciences, Meiji University

†Koji Inui, Department of Mathematical Sciences Based on Modeling and Analysis, School of Interdisciplinary Mathematical Sciences, Meiji University

表3 調査結果と JNSA との比較

	JNSA	本調査	共通
企業数	65	34	23
インシデント数	251	141	80

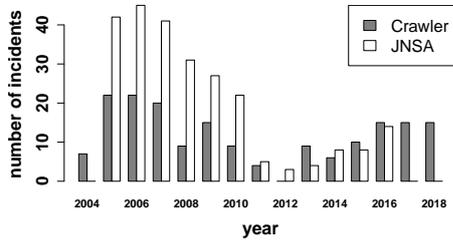


図2 インシデント数経年変化

インシデントリリース(入力)	取得項目(出力)	
	抽出結果	JNSA(正解)
株式会社ディー・エヌ・エー 2016/04/01 DeNAが運営している「Mobage」 において、第三者が「Mobage」 ユーザに成りすまし、不正にログ インしたと思われる事象が判明し た。不正ログインの確認されたID 件数最大104,847件	企業名 業種 日付 被害人数 漏洩原因 インシデント 内容要約 参照記事のURL 社会的責任度 漏洩情報区分 漏洩経路 事後対応	ディー・エヌ・エー 情報通信業 情報通信業・サービス その他 2016/04/01 2016/04/01 104847 104847 不正アクセス 不正アクセス 有 未定義 有 個人情報 インターネット 普通

図3 プレスリリースからの取得項目の違い

3 取得したインシデント情報の評価

3.1 JNSA との比較

表3にインシデントを確認できた企業数とインシデント数を示す。比較対象として、JNSA データセットを使用する。公平に比較するために、表3では本調査結果を2005年から2016年の34企業、141インシデントに限定する。

インシデント数の年による変化を図2に示す。JNSAのインシデントは本調査の調査対象であった537企業に限定している。2005年から2012年にかけてのインシデント数は、JNSAよりも少ないが、これはインシデントについてのリリースがHPに存在しないことが原因である。

図3に、インシデントの例とJNSA データセットと自動分類システムにより抽出した項目を示す。インシデントリリースを入力として、抽出できた取得項目を出力する。正解の出力はJNSAのデータセットにある項目とした。検出結果には、正しい検出と誤検出があるが、この例では全て正しく検出できている。インシデント内容要約は未実装である。

3.2 考察

インシデントの自動分類は適合率50%以下であった。これは、各要素の適合率で最も低かった漏洩原因の誤推定が原因だと考えられる。表4に判定の結果を示す。紛失・置忘れの誤推定のうち75%が管理ミスと判定された。

JNSA データセットとの比較の結果、本調査の537企業におけるインシデント数はJNSAの約半分であった。しかし、収集できていないインシデントの約5.5割はHP上に存在せず、約2.5割は期間が古いため確認できなかった。2013年、2015年、2016年のように近年のインシデントについてはJNSAよりも多く収集できていた。

表4 各漏洩原因の識別数

	紛失・置忘れ	管理ミス	盗難	誤操作	不正アクセス	ワーム等	その他
紛失・置忘れ	53	9	1	0	0	0	2
管理ミス	5	11	1	0	0	0	0
盗難	3	2	29	0	0	0	4
誤操作	0	0	0	12	2	0	1
不正アクセス	0	0	0	1	10	1	2
ワーム等	0	0	0	0	1	5	1
その他	0	1	0	0	4	2	8

4 セキュリティマネジメントの効果

本稿では、提案データセットの効果を確認するため、本調査によるデータセットとCSRデータセットの突合を行い、ロジスティック回帰を行なった。本調査と[5]による結果を表5に示す。正の係数はインシデントの生起確率が増加することを示す。本調査で正の係数となったISMSやCIO、内部監査が[5]では負の係数になるなどの違いがあった。

5 おわりに

本研究では、537企業のプレスリリースについてクローリングを行い、171件のインシデントを収集し、JNSA、「聞蔵II」によるデータセットとの比較を行った。クローリングでは取得記事のうちインシデントは1割であった。収集したインシデントとJNSAデータセットの比較結果、61件の独自のインシデントが見つかった。インシデントの日付、被害人数、漏洩原因は70%以上の適合率で正しい分類ができたが、全要素を正しく分類する例は50%以下であった。

今後は、インシデントの識別精度を高め、対象の企業数を増やし汎用性の高いDBを提供することを課題とする。

参考文献

- [1] 情報セキュリティインシデント調査報告書 (JNSA データセット)
- [2] 聞蔵II, 朝日新聞社, オンライン記事データベース
- [3] 池上和輝, 菊池浩明, インシデント調査に基づく漏えい原因のデータマイニング, 情報処理学会第80回大会, 2W-06, vol.3, pp.543-544, 2018.
- [4] 東洋経済データサービス CSR データ
- [5] 山田道洋, 池上和輝, 菊池浩明, 乾考治, 経営マネジメント状況による情報漏洩インシデント削減効果の評価 (2), 情報処理学会コンピュータセキュリティシンポジウム (CSS2018) 2018, pp.376-384, 2018.

表5 本調査によるデータセットを用いたロジスティック回帰

	(Intercept)	log(従業員数)	ISMS	CIO	内部監査	外部監査	内部窓口	外部窓口
本調査	-23.260	0.399	1.200	0.000	2.429	-0.959	-1.717	-0.969
[5]	-8.300	0.948	-0.217	-1.097	-2.070	0.117	-0.050	-0.685