

# 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<sup>1</sup> 山田 道洋<sup>2</sup> 菊池 浩明<sup>1</sup> 乾 孝治<sup>3</sup>

1:明治大学総合数理学部先端メディアサイエンス学科

2: 明治大学大学院先端数理科学研究科

3: 明治大学総合数理学部現象数理学科

# 研究背景

- 不正アクセスや内部犯行などによる情報漏洩の増加

- 企業の被害

- データの損失・破損
- 事業妨害
- 損害賠償
- 広報活動

- セキュリティマネジメントが必要

- ISMSを導入するとインシデントを20%削減することが示されている [1]



[1] 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)", Computer Security Symposium 2018, pp.376-384, 2018.

# 先行研究・問題点

## 先行研究

- メディアから手動で収集[2]

JNSA	朝日新聞[聞蔵Ⅱ]	共通
788	279	145

## 問題点

1. 新聞やネット記事を元にするると、ニュース性の高い特異なインシデントが多く取り上げられ網羅的に収集できない可能性がある。
2. 人手を使った収集は、コストと時間がかかる。
  - 1企業の1年間のリリースを確認するのに平均6分
  - JNSAは手作業により収集

[2] 池上和輝, 菊池浩明, "インシデント調査に基づく漏洩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.

# 研究目的・解決方法

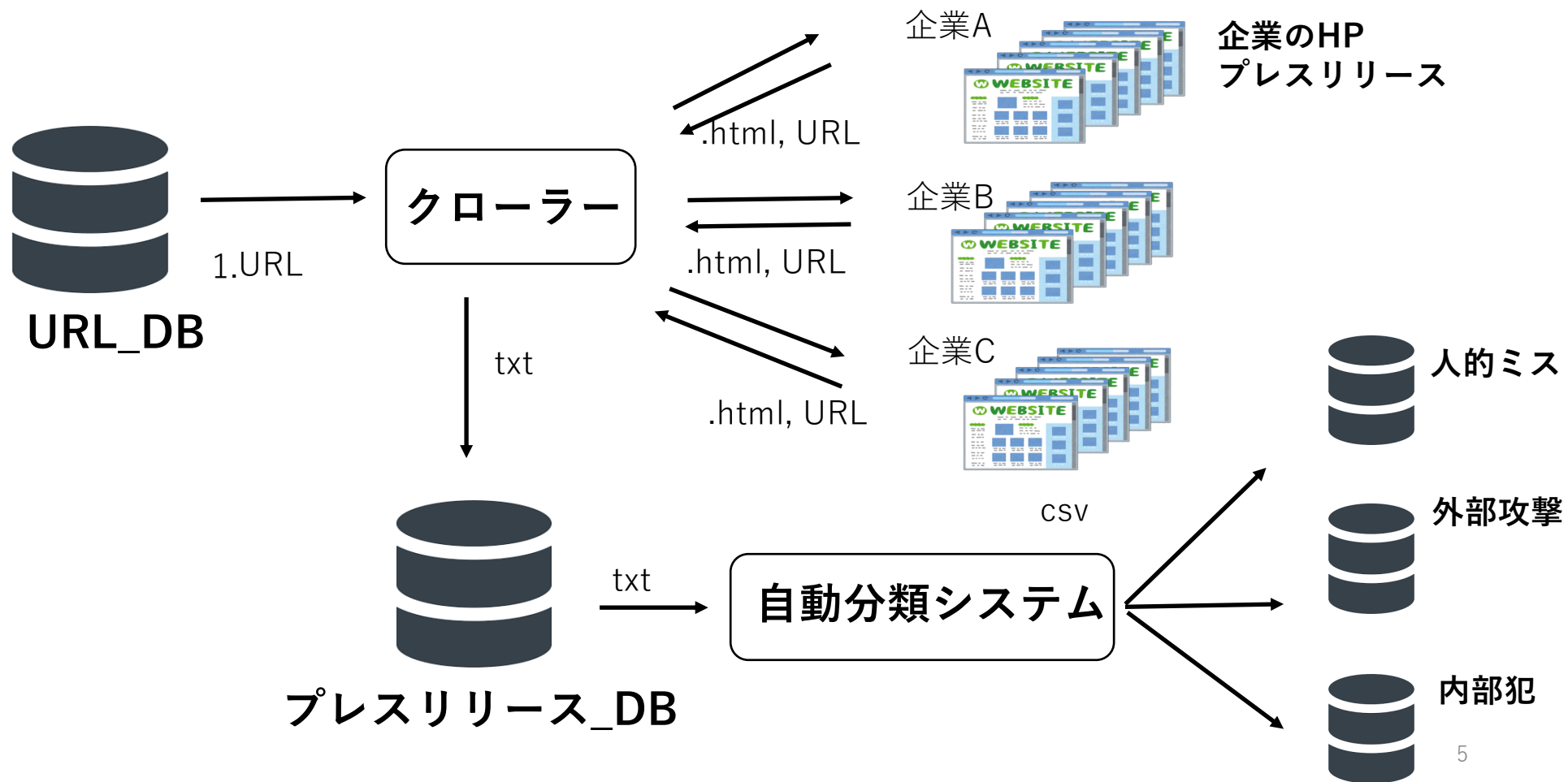
- 研究目的

- メディアなどの偏りなく、網羅的にインシデントを自動収集・分類すること

- 解決方法

1. 独自クローラーによる企業サイトのプレスリリース収集
2. 自然言語処理と49次元特徴語ベクトルのcos類似度による漏洩原因の自動分類システムの開発

# システムの全体構成図



# TF・IDF値と49次元ベクトル例

TF値 : 文書中の単語出現頻度

IDF値 : 文書集合の中で特定の単語を含む文書の割合の逆数

不正アクセス

	文書に含まれる単語	TF	IDF	TF-IDF
特徴語	パスワード	0.006	2.696	0.018
	ウェブサイト	0.003	3.572	0.012
	不正	0.004	2.696	0.012
不要	サーバー	0.001	3.235	0.006
	アカウント	0.003	2.783	0.008
	対策	0.002	2.879	0.008

# 自動分類システム(例)

## 特徴語の出現頻度比較

### 記事(入力)

2019/2/2  
A社の所有するパソコンに不正アクセスがあったことを確認した。  
約3000件の個人情報が出た可能性がある。

正規表現により抽出

特徴語	入力
紛失	0
不正	1
パソコン	1
委託	0
メール	0

cos類似度

人的ミス	不正アクセス	内部犯行
3	0	0
0	2	0
1	1	4
1	1	3
2	0	0

0.33

0.87\*

0.45

出力：日付：2019/2/2，規模：3000，原因：不正アクセス

# JNSAデータセットとの項目比較

インシデントリリース(入力)

取得項目(出力)

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

	抽出結果	正解(JNSA)
企業名	ディー・エヌ・エー	ディー・エヌ・エー
業種	情報通信・サービス その他	情報通信業
日付	2016/04/01	2016/04/01
被害人数	104,847	104,847
漏洩原因	不正アクセス	不正アクセス
インシデント 内容要約	未定義	有
参照記事のURL		
社会的責任度		普通
漏洩情報区分		個人情報
漏洩経路		インターネット
事後対応		普通



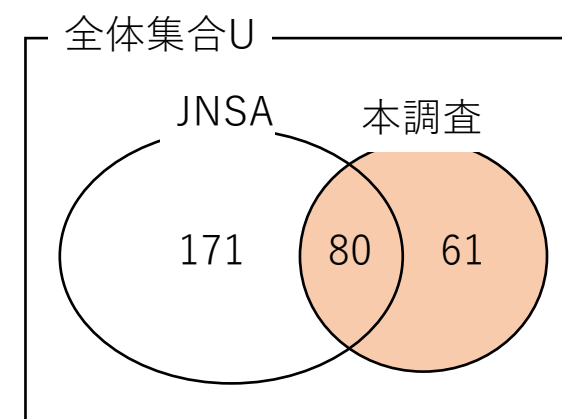
# 収集結果・JNSAとの比較

## ・ 収集結果

企業数	期間	取得記事数	インシデント数	割合
537	2004/10/1~ 2018/11/2(4年 間)	17,957	191	(0.01)

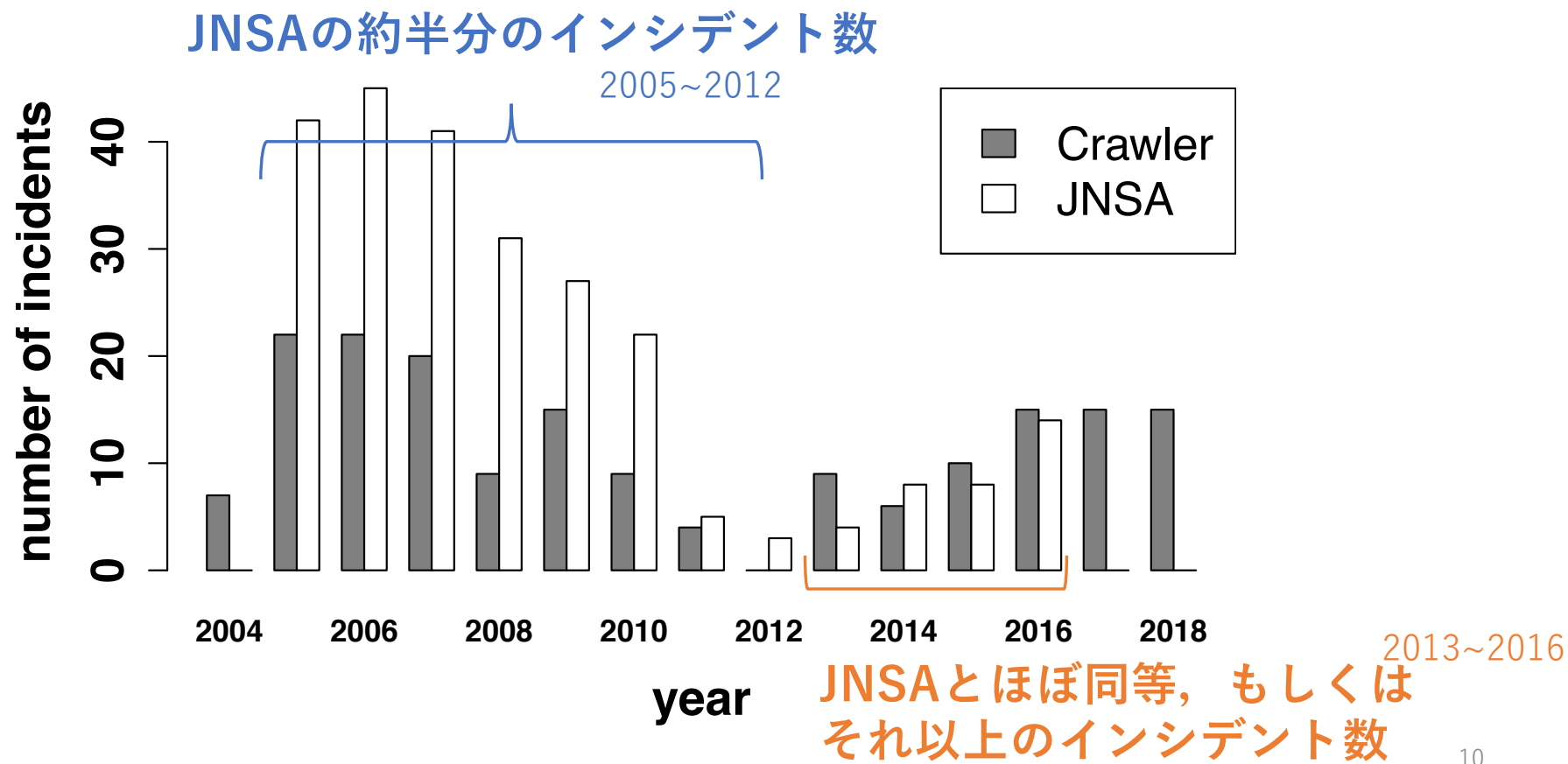
## ・ 比較

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



全インシデントにおけるクローラーの収集割合45%

# 経年変化



# インシデント例

サイバーエージェント

2010/01/01

不正アクセスにより450件のID、パスワードが流出した可能性がある

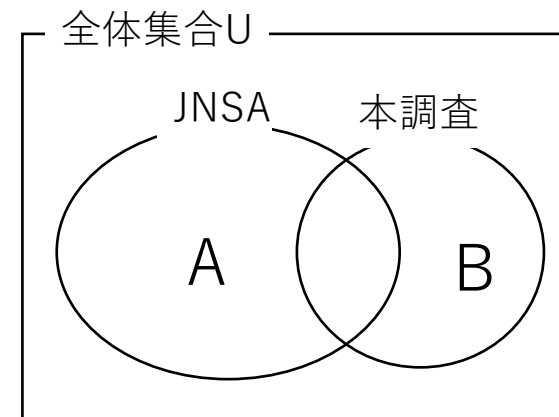
A

東京ガス

2006/12/8

お客様のガス料金収納を担当する事業所において、お客様情報の記載された3,463件の領収書を紛失した。

B



# 自動分類の結果

$$\text{推定正解率} = \frac{\text{正しく要素を抽出できたインシデント}}{\text{収集した全インシデント}}$$

	日付	被害人数	漏洩原因	日付&規模&原因
推定正解率	0.882	0.792	0.719	0.505

各要素はそれぞれ**7割以上**の推定正解率で分類できた  
すべての要素を合わせると約5割

# 応用例：マネジメント効果(ロジステック回帰)

- 交絡因子（業種，企業規模，年）を排除し，マネジメントの効果を測る

- ある企業のインシデントの生起確率  $p = \frac{1}{1+e^{-z}}$

$$z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{x_m} X_m$$

- $\beta_i$ ：係数， $X_i$ ：発生要因（例：従業員数、ISMS、外部監査）
- 上記の式を書き直すと以下になる
    - $\log \frac{p}{1-p} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{x_m} X_m$
    - 例：  $\log \frac{p}{1-p} = 0.33 + 1.21 * M_1 + 0.44 * M_2 - 0.61 * M_3$

# 応用例：マネジメント効果(ロジステック回帰)

- $\log \frac{P}{1-P} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{x_m} X_m$

	(Intercept)	従業員	ISMS	CIO	外部監査
Estimate( $\beta$ )	-23.26	0.399	1.222	0.0002	-0.959

ISMSを導入していると、していない企業よりインシデントが生じやすい

## まとめ

- クローラーにより191のインシデント記事を収集した
- 2013年以降ではJNSAと同等のインシデントを収集できた
- 要素別の分類では全て7割を超えた推定正解率で抽出した
- ISMSを導入していると、していない企業よりインシデントを報告しやすい