

## 匿名加工情報取扱事業者を調査するクローラーシステムの開発

金子 侑紀 †    小野 敦樹 †    伊藤 聡志 †    菊池 浩明 †    服部 充洋 ‡    飯田 泰興 ‡  
 藤田 真浩 ‡    山中 忠和 ‡  
 明治大学総合数理学部 †    三菱電機株式会社 ‡

### 1 はじめに

2017年5月30日に全面施行された改正個人情報保護法によって、中小企業をはじめとする全ての事業者が個人情報保護法の対象となった。また、匿名加工情報という新たな情報の類型が定義されたことにより、一定の条件の下で、本人の同意がなくても第三者提供や目的外利用が可能となった。しかし、匿名加工情報を第三者提供をする際に、匿名加工情報取扱事業者は提供する情報に含まれる個人に関する情報の項目の公表が義務付けられているが、届け出は義務ではないため、現在個人情報保護委員会は匿名加工情報公表ページから、個人に関する情報の項目の公表を約300社分手作業で取得を行っている。

そこで本研究では、匿名加工情報公表ページのクローリングを行い、提供項目とその手法について、APIを用いて自動取得を試みる。しかし、本研究で使用したAPIは、1度の動作で検索結果上位100件のデータのみしか取得できない制限がある為、検索対象の業種を絞り込むシードキーワードの導入をし、APIとシードキーワードの組み合わせによって効果的に取得した。

### 2 手動での匿名加工情報公表サイト調査

自動クローラーシステムの開発にあたり、手動により、匿名加工情報公表サイトと個人に関する情報の項目を収集した。

#### 2.1 調査結果と考察

2019年5月から2019年8月までに公表されていた匿名加工情報公表サイト数は表1の通り。

表1 手動での匿名加工情報公表サイト調査結果

期間	匿名加工情報公表サイト数
2019/05-2019/08	308

本調査により取得した308企業を、日本標準産業分類を基とした26業種に分類をした。匿名加工情報公表サ

Development of crawler system for anonymously processed information handling business operators

†Yuki Kaneko, Atsuki Ono, Satoshi Ito, Hiroaki Kikuchi, School of Interdisciplinary Mathematical Science, Meiji University.

‡Mitsuhiro Hattori, Yasuoki Iida, Masahiro Fujita, Tadakazu Yamanaka, Mitsubishi Electric Corporation

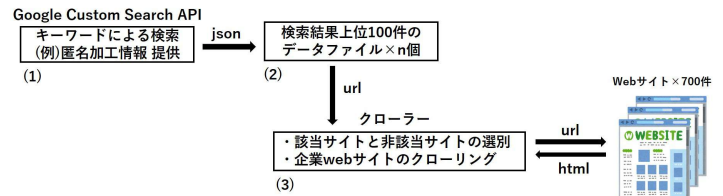


図1 システム構成図

イト企業業種別数を表2に示す。上位5業種に、医療業(病院), 健康保険組合, 小売業(薬局)が含まれていることから、匿名加工情報公表データには医療データが多く含まれていることが分かる。

表2 匿名加工情報公表サイト企業業種別数

医療業(病院)	58	職業紹介, 労働者派遣業	4
情報通信業	45	年金相談センター	4
健康保険組合	34	不動産業	4
小売業(薬局)	29	一般社団法人	3
製造業	21	運輸業	3
保険業	19	卸売業	3
情報サービス業	17	健康保険協会	3
サービス業	14	公益社団法人	2
金融業	14	特別民間法人	2
小売業	12	建設業	1
医療業(製薬)	5	信用格付け機関	1
社会保険, 社会福祉, 介護事業	5	弁護士会	1
教育, 学習支援業	4	保険労務士法人	1

### 3 自動クローラーシステムの開発

#### 3.1 システム構成

システムの全体構成を図1に示す。(1)Google Custom Search APIを使用し、キーワードによる検索を行う。(2)検索結果上位100件のサイトタイトル及びURLをjsonファイルとして取得する。(3)サイトタイトルに特定のキーワードを含むサイトのhtmlを保存し、テキスト化を行う。キーワードには手動調査結果に基づき、匿名加工情報公表サイトタイトルに含まれるであろう単語、「匿名加工情報」、「個人情報保護」、「プライバシーポリシー」をシードキーワードとして選定した。

表3 システム構成図(1)における検索に使用したシードキーワード

業種	共通キーワード	シードキーワード
病院	匿名加工情報 作成 提供 提供方法 -法律事務所 -個人情報保護委員会	病院
薬局	匿名加工情報 作成 提供 提供方法 -法律事務所 -個人情報保護委員会	調剤
健康保険関連	匿名加工情報 作成 提供 提供方法 -法律事務所 -個人情報保護委員会	健康保険組合 or 健康保険協会
生命保険	匿名加工情報 作成 提供 提供方法 -法律事務所 -個人情報保護委員会	生命保険
銀行	匿名加工情報 作成 提供 提供方法 -法律事務所 -個人情報保護委員会	銀行
年金関連	匿名加工情報 作成 提供 提供方法 -法律事務所 -個人情報保護委員会	年金
その他	匿名加工情報 作成 提供 提供方法 -法律事務所 -個人情報保護委員会	-病院 -健康保険組合 -健康保険協会 -銀行 -年金 -生命保険

## 4 自動クローラーシステムの精度評価

### 4.1 実験方法

1. 2019年11月18日, 病院, 薬局, 健康保険関連, 生命保険, 銀行, 年金, その他の7分類に検索ターゲットを分け, 計700件のデータを取得する(図1(2)より, 本実験の場合  $n = 7$ ). 検索に使用したキーワードを表3に示す.

2. 図1(3)による, 700件の取得データの, 匿名加工情報公表サイト, 非匿名加工情報サイトへの自動選別システムの精度を, 自動クローラーが収集したデータのうち, 重複して収集をしたサイト, 正解でないサイトを手作業で調査し, 評価する. 2において収集したデータと2章の実験において手作業で収集したデータとの差異を調査し, 考察する.

### 4.2 実験結果

自動クローラーにより, 計321社の匿名加工情報公表サイトの収集ができた. 結果を表4に示す. 321社のデータを2章の手動クローラーで収集したデータと比較をしたところ, 新たに210社の匿名加工情報公表サイトを取得することができ, 111社のサイトが重複していた.

### 4.3 考察

計7分野のうち, 「病院」, 「薬局」, 「健康保険関連」分野では手作業取得データ数を上回るデータ数を本自動クローラーにより取得することができた. しかし, 「その他」の分野において手動クローラーでは169社のデータを取得できたのに対し, 自動クローラーでは75社のみしか取得できなかった. それは「その他」に分類されて

表4 自動クローラーによる匿名加工情報公表サイト取得結果

業種・団体	重複データ	新データ	合計	手作業収集データ
病院	28	48	76	58
薬局	15	58	73	29
健康保険関連	8	78	86	37
生命保険	3	0	3	5
銀行	2	2	4	4
年金関連	0	4	4	6
その他	55	20	75	169

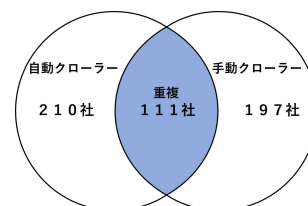


図2 自動クローラー手動クローラー収集データの差異

いる企業を共通のターゲットとしての検索が困難であること, また本自動クローラーの検索システムとして使用している Google 社の CustomSearchAPI の仕様により, 検索結果上位100サイトのみばかりしか取得出来ないことも大きな要因であることが考えられる.

## 5 おわりに

本研究では, 4章実験結果より, 自動クローラーシステムにより計321社の匿名加工情報公表サイトを収集することができ, 手動クローラーによる収集数を上回る結果を示した. また, 手動クローラーでは収集することのできなかつた新たな210社のデータを収集できた. 検索ターゲット細分化や検索ワードの選定精度をより向上させる事で, 有用性を向上する予定である.

今後は, 本研究で得た匿名加工情報公表サイトのデータを基に, 検索ターゲット細分化や検索キーワードの選定精度向上の為の分析を用いての新たな自動クローラーシステムの開発を試す必要がある.

## 参考文献

- [1] 濱田, 荒井, 小栗, 菊池, 黒政, 中川, 西山, 波多野, 村上, 山岡, 山田, 渡辺' 匿名加工再識別コンテンツの設計 履歴データの一般化, 再識別', PWS Cup 2018, pp935 - 940, 2018.
- [2] 小林祐貴' 一般化匿名加工された購買履歴データの顧客・商品の RFM 分析', 明治大学菊池研究室 2018 年度卒業論文, 2018.
- [3] GoogleCustomSearchAPI(<https://developers.google.com/custom-search/v1/overview?hl=ja>), 2019年11月参照