

匿名加工情報公表サイト調査 (2) 加工対象情報の抽出

金子 侑紀†

明治大学総合数理学部 先端メディアサイエンス学科 菊池研究室†

表 1 匿名加工情報公表サイト抽出結果

企業名	個人に関する情報の項目	提供手法
株式会社日本医薬総合研究所	年齢(生年)、性別、処方せん情報、調剤情報、各種アンケート回答	電子メール、CD-ROM、USB等の外部記憶媒体、HTTPS
セキ薬局	氏名、生年月日、被保険者記号番号、公費受給者番号、医師氏名、処方日、調剤日、性別、生年、処方・調剤履歴	セキュリティが確立された転送方式
イオン銀行	性別、年代、申込手段、当行普通預金口座の有無、現在の借入の有無、契約から初回借入までの経過日数、現在の返済実績等	パスワードで保護し、CD-ROMで手交

1 はじめに

個人情報取扱事業者が個人情報データベースを匿名加工し、作成した匿名加工情報を第三者へと提供する場合には、あらかじめ、提供する匿名加工情報に含まれる、個人に関する情報の項目を公表するとともに、提供先に対し、匿名加工情報である旨を明示しなければならない(改正個人情報保護法第 37 条)。

しかし、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は義務付けられているが、届出や申請の必要はない。匿名加工情報の全貌を知るためには各社の公表ページから、個人に関する情報の項目を手作業で収集する他はない。

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、5 パターンの正規表現を用いて提供項目とその手法についての自動取得を試みる。前者のクローリングについては [3] で報告し、本稿では、後者の自動取得について述べる。本論文では匿名加工取扱事業者のサイトを匿名加工情報公表サイトと定義する。

2 抽出システムの開発

抽出システムの構成を図 1 に示す。クローラーが取得した匿名加工情報公表サイトの HTML ファイルから提供する項目と手法を抽出する。抽出した結果は匿名加工情報公表企業及び項目一覧 DB として N 行 \times 4 列(企業名、個人に関する情報の項目、情報の提供手法、URL)の形式で表で管理する。 N は企業数である。

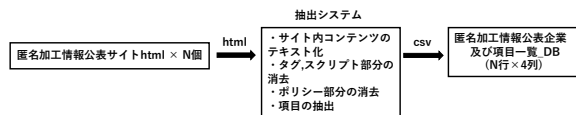


図 1 システム構成図

```
...<h2 class="mMH01">匿名加工情報の提供について</h2>
<p class="mText">当組合では、保健事業や疫学調査等のために匿名加工情報を継続的に作成し、電子的な通信手段もしくはDVD等の物理媒体を用いてレセプト分析業者に提供します。作成及び提供する匿名加工情報に含まれる情報の項目は、性別、生年月、医療保険の資格情報(加入時期、脱退時期、本人・家族区分等)、診療報酬明細書の受診履歴、健診の受診履歴です。なお、個人を特定できる情報は含まれておりません。</p>...
```

入力

出力

図 2 抽出システム処理例

抽出システムの処理の例を図 2 に示す。入力文章にタグの処理とヘッダ、フッタなどの対象外の部分を取り除く処理を施した後、情報の項目、提供方法について正規表現による抽出をする。情報の項目を抽出する正規表現パターンは 5 種類用意する。正規表現の例を表 3 に示す。

3 抽出システムの評価

3.1 実験方法

2019 年 11 月 18 日に、クローラーを動作させて収集した 321 件のファイルに対して、本システムによる抽出を試みる。

表 2 抽出目標項目

項目	例
個人に関する情報の項目	性別、年代
提供方法	電子メールによる提供

†Yuki Kaneko, Department of Frontier Media Science, School of Interdisciplinary Mathematical Science, Meiji University, Kikuchi Laboratory.

3.2 実験結果

本システムにより、情報の項目と提供方法の抽出結果を表4に示す。取得したデータが過剰である例を図3に示す。

表3 正規表現の例

正規表現	件数
(次のとおり 下記 以下)+([\s\S]+)\n([\s\S]+ 提供 [の])*方法	44
DPC+([\s\S]+)\n([\s\S]+ 提供 [の])*方法	14
情報 [の]*項目 ([\s\S]+)\n([\s\S]+ 提供 [の])*方法	4
項目は ([\s\S]+) です	3
([\s\S]+)(上記項目 提供 [の])*方法 第三者に提供)	60

表4 匿名加工情報公表サイト抽出結果

適切に取得できた	取得したデータが過剰	取得できなかった	計
125	27	167	319

...サイト内検索(中略)取り組み
匿名加工情報の作成及び第三者提供について
匿名加工情報の取り組み
匿名加工情報に関する問合せ窓口
匿名加工情報の取り組み
DPC制度の導入の影響評価及び今後のDPC制度の見直しを図る目的で、厚生労働省が収集し管理する情報となるデータ（DPCデータ）を作成しております。また、審査支払機関への請求のため診療に係る費用を診療報酬明細書（レセプト）として作成しております。

DPCデータは、診療録からの情報および診療報酬明細書からの情報で構成されており、レセプトデータは、医療機関情報・保険者情報・診療行為情報・医薬品情報・特定器材情報等から構成されております。

DPCデータ並びにレセプトデータを活用することで、医療の質向上および病院経営の改善に役立てる事が可能になるため、匿名加工後のデータを第三者へ提供しております。第三者提供するこれらのデータは氏名、住所、電話番号は含みません。なお、地域傾向や受診年齢層等を分析するため必要があるため、郵便番号（上3桁のみ）、生年月日(生年月及び入院時年齢に変換を行い100歳以上は100歳に一括り)、各種保険証に関する情報については保険者番号（健康保険事業の各運営主体を指す番号）のみを含みます。

当院は上述の通り、診療情報から匿名加工情報を作成（毎月継続）し、第三者に提供しております。
匿名加工情報の提供の方法

図3 取得したデータが過剰である例

3.3 考察

取得できなかった公表サイトの理由を表5に示す。

クローラー側から渡されたファイルはすべてHTML形式だったが、PDFが多くあった。そこで、PDF形式に変更して処理を行ったところ、75件全ファイルのデータが空であり、取得できていなかった。

PDF以外のファイルについて、文字コードが間違っ保存されているものが多く存在した。PythonのrequestsライブラリはHTTPレスポンスヘッダのcontent-typeが指定されていないファイルはdefaultのISO8859-1形式で保存されるため、読み込めないすべてのファイルに対して強制的に変換を行っている。25件中2件がサーバーエラー、23件がISO8859-1形式への変換後も問題が解消されなかったファイルである。文字コードが原因

のものは適切な変換手法が見つければ項目を抽出できる可能性がある。

取得したデータが過剰であるファイルは、ヘッダやフッタ部分が特徴的で削除しきれなかったものや、匿名加工情報についての記載がプライバシーポリシーの一部として記載されていた。

抽出対象が含まれていないファイルとは法令で定める基準に従う旨が記載されており、個人に関する情報の項目や提供方法が明記されていないファイルである。

表5 取得に失敗した公表サイトの内訳

理由	件数	例
PDFの破損	75	抽出不可
サーバーエラー・ファイル破損	25	500等、ページ取得に失敗
提供していないと明記しているサイト	5	現在、匿名加工情報は取扱っておりません。
公表サイトでない	3	匿名加工情報作成ツールの記事
項目と方法は別ファイルに記載	11	弊社が作成する匿名加工情報に含まれる個人に関する情報の項目、弊社が第三者に提供する匿名加工情報に含まれる個人に関する情報の項目および提供の方法については、こちらをご覧ください。
抽出対象が含まれていない	42	匿名加工情報を作成した場合には、匿名加工情報に含まれる個人に関する情報の項目を公表いたします。匿名加工情報を第三者提供する場合には、提供する匿名加工情報に含まれる個人に関する情報の項目および提供の方法について公表するとともに、提供先に、提供される情報が匿名加工情報である旨を明示いたします。
抽出失敗	8	抽出不可

4 おわりに

本研究では匿名加工情報取扱事業者の公示ページに含まれる情報の自動抽出を実現した。

参考文献

- [1] 濱田, 他, “匿名加工再識別コンテキストの設計 履歴データの一般化, 再識別”, CSS 2018, pp935-940, 2018 情報処理学会 2018.
- [2] 匿名加工情報 - 個人情報保護委員会, (<https://www.ppc.go.jp/personalinfo/tokumeikakouInfo/>, 2019年12月11日参照).
- [3] 小野敦樹, “匿名加工情報公表サイト調査 (1) 自動クローラーシステムの開発”, 明治大学菊池研究室 2019年度卒業論文, 2019