

# $Pk$ -匿名化データの有用性評価実験

正木 彰伍<sup>1,a)</sup> 五十嵐 大<sup>1</sup> 菊池 亮<sup>1</sup> 齋藤 恆和<sup>1</sup> 千田 浩司<sup>1</sup> 廣田 啓一<sup>1</sup>

**概要:** パーソナルデータの安全な利活用には、データの安全性と有用性の両立が必要である。安全性については代表的な  $k$ -匿名性や、 $k$ -匿名性を確率的な指標に拡張した  $Pk$ -匿名性が提案されるなど、匿名化技術が広く研究されている。一方で、匿名化データの有用性についての議論は未だ限定的である。特に、レコード数などのデータの特徴と有用性の関係性を明らかにすることは、実用上非常に有益である。しかし、これまで行われてきた、ウェブ上で公開されている実データなどを用いた実験では、用いるデータの特徴が限定的になり、議論が困難となっていた。そこで本稿では、多くのデータを包含する一般的な模擬データモデルを利用した評価法を提案し、この模擬データに、 $Pk$ -匿名化を適用した実験を行う。さらに実験結果から、有用性と模擬データモデルのパラメーターの関係について調べ、特定のパラメーターから有用性を予測できることがわかった。

## 1. はじめに

ICTの発達に伴い、サービス事業者には、ユーザーの購買履歴や位置情報といったデータが豊富に蓄積されるようになってきた。昨今、そういったパーソナルデータを二次利用し、サービス開発や商品販売戦略などに活用したいというニーズが高まっている。一方で、ユーザー個人に紐づくパーソナルデータの利活用にはプライバシー侵害のリスクがあり、慎重な取り扱いが求められる。対応策の一つとして、データを安全な形に加工する匿名化技術が広く研究されている。

匿名化技術の中でも代表的な指標として、 $k$ -匿名性 [1], [2] が知られている。 $k$ -匿名性は、1つのデータ中に“同じ属性の値の組み合わせを持つレコードが  $k$  個以上存在する”という直観的でわかりやすいプライバシー保護指標である。一方、データへのノイズ付加といった確率的な操作を伴う匿名化には適用できなかった。これに対し、五十嵐ら [3] は、 $k$ -匿名性を確率的な指標に拡張し、“ある人のデータを  $1/k$  以上の確率で当てることができない”ことを保証する  $Pk$ -匿名性を提案した。 $Pk$ -匿名性を満たす匿名化方法 ( $Pk$ -匿名化) としては、カテゴリ属性に対する維持置換攪乱 [4]、数値属性に対するラプラスノイズ加算 [5]、有界ノイズ加算 [6]、データ依存型維持置換攪乱 [7] および属性間の相関を考慮した攪乱再構築法 [8] がある。

$Pk$ -匿名化は、確率的なノイズ付加を伴う加工であるた

め、匿名化されたパーソナルデータは元と異なるデータになってしまう。匿名化データを利活用するうえで、 $Pk$ -匿名化によってどれほど変化するか、有用性はどれほど保たれるかは明らかにする必要がある。特に、データの特徴から匿名化データの有用性の関係が明らかになれば、 $Pk$ -匿名化を施さずとも有用性が予測でき、実用上非常に有益である。これまで、Adult Data Set [9] に代表されるウェブで公開されているいくつかのデータを使った、実験的な有用性評価が行われてきた [10], [11]。しかし、ウェブ上の公開データや購入可能なデータを用いるだけでは、利用可能なデータの特徴が限定的になり、データの特徴と匿名化データの有用性の関係を議論することが困難となっていた。

### 1.1 本稿の成果

本稿では、多くのデータを包含する一般的な模擬データモデルを利用した評価法を提案し、この模擬データに、 $Pk$ -匿名化を適用した実験を行う。実験結果から、匿名化データの有用性はデータ分布にあまり依存せず、特定のパラメーターから有用性を予測できることがわかった。

## 2. 準備

### 2.1 テーブルとクロス集計

本稿で扱うテーブルとは以下のようなものである。テーブルは情報提供者から提供されたデータの集合である。各情報提供者からのデータはテーブル上では1行に表現され、これをレコードと呼ぶ。各レコードはあらかじめ定められた項目に対する値から成り立っており、この項目を属性といい、取りうる値の範囲を属性の値域という。表1にテ

<sup>1</sup> 日本電信電話株式会社 NTT セキュアプラットフォーム研究所  
180-8585 東京都武蔵野市緑町 3-9-11

a) masaki.shogo@lab.ntt.co.jp

ブルの例を示す。“性別”及び“年代”が属性であり，各行“女性，10代”等がレコードである。

表1 テーブルの例

ID	年代	性別
1	10代	女性
2	40代	男性
3	20代	男性
4	30代	女性
⋮	⋮	⋮

クロス集計とは，テーブルの複数の属性に着目し，その全ての属性の値が等しいレコードの数を表すものである。このクロス集計は，アンケート集計で使われるほかに数多くのデータマイニング手法の中で使用される基本的かつ重要な統計量である。一般的にクロス集計は行ベクトルとして扱われる。その大きさは各属性の組み合わせの総数，つまり各属性の値域の積となり，ベクトルの要素の和はレコードの総数と一致する。表2が表1におけるテーブルでのクロス集計である。ベクトルの第2要素である“{10代，女性}”はテーブルの中にある10代かつ女性である人数の総和であり，右辺のベクトルの第二要素はそれが45人であることを示している。

表2 クロス集計の例

$$\begin{aligned} & \{10代, 男性\}, \{10代, 女性\}, \\ & \{20代, 男性\}, \{20代, 女性\}, \dots \\ & = (61, 45, 73, 50, \dots) \end{aligned}$$

## 2.2 攪乱再構築法によるPk-匿名化

攪乱再構築法とは，データを改変することによりプライバシーを保護する攪乱と，攪乱データから元データを推定する再構築という2つのアルゴリズムからなる匿名化手法である。

### 2.2.1 維持置換攪乱

具体的な攪乱方法として，属性値を一定の確率で維持し，それ以外では取りうる値の中でランダムに置換する維持置換攪乱がある[4]。ある属性*i*の値*v<sub>i</sub>*が，*v'<sub>i</sub>*に遷移する確率を，サイズ*M<sub>i</sub> × M<sub>i</sub>*の行列(*A<sub>i</sub>*)<sub>*v<sub>i</sub>, v'<sub>i</sub>*</sub>を用いて，

$$(A_i)_{v_i, v'_i} = \begin{cases} \rho_i + \frac{1-\rho_i}{M_i} & \text{for } v_i = v'_i \\ \frac{1-\rho_i}{M_i} & \text{for } v_i \neq v'_i \end{cases} \quad (1)$$

と表す。ここで， $\rho_i$ は属性*i*の維持確率と呼ばれ $0 \leq \rho_i \leq 1$ であり，*M<sub>i</sub>*は属性*i*の値域である。属性が*n*個あった場合，行列  $A = A_1 \otimes A_2 \otimes \dots \otimes A_n$  を遷移確率行列と呼ぶ。ここで， $\otimes$ はクロネッカー積である。

この維持置換攪乱は，維持確率を適切に設定することで，Pk-匿名性を満たすことができる。 $\mathcal{R}$ をテーブルのレコード集合とした場合，パラメーター間で

$$k = 1 + (|\mathcal{R}| - 1) \prod_{i=1}^n \left[ \frac{1 - \rho_i}{1 + (M_i - 1)\rho_i} \right]^2 \quad (2)$$

の関係が成立する必要がある[3]。

### 2.2.2 再構築

具体的な再構築方法としては，反復ベイズ推定法が知られている[12]。これは，攪乱テーブルのクロス集計に対し，遷移確率行列*A*を用いてベイズ推定を繰り返し行うことで，元テーブルのクロス集計を推定するものである。

反復ベイズ推定法のアルゴリズムをAlgorithm 1に示す。ここで，*j*は反復回数を表し，ベクトル  $\vec{v}_1, \vec{v}_2$  に対し

#### Algorithm 1 ベイズ推定法による再構築アルゴリズム

---

```

Input:  $A, \vec{y}$ 
Output:  $\vec{x}_j$ 
 $\vec{x}_0 := \vec{y}$ 
 $j := 0$ 
while  $|\vec{x}_{j+1} - \vec{x}_j|_{L_1} \geq \epsilon$  do
     $\vec{x}_{j+1} := \vec{x}_j * (A(\vec{y}/\vec{x}_j A))^t$ 
     $j := j + 1$ 
end while
    
```

---

て  $\vec{v}_1 * \vec{v}_2$  と  $\vec{v}_1 / \vec{v}_2$  はベクトルの成分ごとの積と商を表す。 $|\vec{x}_{j+1} - \vec{x}_j|_{L_1}$  は  $\vec{x}_{j+1}, \vec{x}_j$  の  $L_1$  距離と呼ばれ，各要素の差の絶対値の総和であり，本稿では総レコード数で割ったものを用いる。 $\epsilon$ はあらかじめ定める収束半径である。遷移確率行列*A*と攪乱テーブルのクロス集計 $\vec{y}$ が入力であり，出力として推定した元テーブルのクロス集計 $\vec{x}_j$ を得る。

## 3. 提案する評価法

### 3.1 概要

ウェブ上の公開データや購入可能なデータを用いるだけでは，利用可能なデータの特徴が限定的になり，網羅的に有用性評価実験を行うことは困難になる。そこで，多様なデータを用いて実験を行うために，データの特徴づけるパラメーターを設定し，値を動かして多数の模擬データを作成し，実験に用いる。

手順としては以下の通りである。

- (1) データの特徴づけるパラメーターの設定
- (2) パラメーターの値を変えながら，模擬データを作成
- (3) 作成した模擬データに対してPk-匿名化を実行
- (4) 匿名化前後のデータを比較し，有用性を評価

ポイントとなるのは，データの特徴の1つである属性値の分布のモデル化およびパラメーター化である(3.3.3を参照)。

### 3.2 1属性データと属性の連結化を用いたPk-匿名化

本稿では，属性が1つの模擬データを作成し，実験に用いる。この属性が1つのみの模擬データにPk-匿名化を適用することは，属性が複数あるデータに属性の連結化を用

いた  $Pk$ -匿名化を適用することと等しい。

属性の連結化とは、複数の属性の値を一括りにして、1つの属性の値として扱うことである。表3に、表1の“年代”と“性別”を連結化した例を示す。当該表の属性は、“年

表3 表1の属性を連結化したテーブル

ID	年代・性別
1	10代・女性
2	40代・男性
3	20代・男性
4	30代・女性
⋮	⋮

代・性別”の1つとなり、IDが1のレコードの属性値は、“10代・女性”となる。

属性の連結化によるメリットは、 $Pk$ -匿名化されたデータの有用性向上が期待できる点にある。 $Pk$ -匿名化は攪乱を伴うため、属性を連結せずに  $Pk$ -匿名化を行うと、元データには存在しない属性値の組合せが、匿名化データに現れることがある。属性を連結することでこれを防ぐことができ、 $Pk$ -匿名化されたデータの有用性向上、特に連結された属性間の相関がより良く保存されることが期待できる。

### 3.3 パラメーター値の設定

模擬の1属性データを作成するにあたって、1属性データの特徴づけるパラメーターとして、次の3つを考える。

- レコード数
- 値域
- 値の分布

#### 3.3.1 レコード数 $N_{\text{rec}}$

現実的な値として、1,000 (=  $10^3$ ) から 10,000,000 (=  $10^7$ ) の範囲を考える。具体的には、 $10^3$ ,  $2 \times 10^3$ ,  $5 \times 10^3$ ,  $10^4$ ,  $2 \times 10^4$ ,  $5 \times 10^4$ , ...,  $10^7$  の13個の値を取る。

#### 3.3.2 値域 $V_{\text{att}}$

値域とは、属性値の種類の数である。レコード数が増えれば取りうる値域が増えるため、レコード数に対する比をパラメーターとする。本稿では、具体的には、 $V_{\text{att}}/N_{\text{rec}} = 10^{-6}, 2 \times 10^{-6}, 10^{-5}, 2 \times 10^{-5}, \dots, 10^{-1}, 2 \times 10^{-1}$  の値を取る。ただし、実験の冗長さを防ぐため、利活用の価値が比較的低いと考えられる  $V_{\text{att}}$  が10に満たない模擬データは実験から除外する。

#### 3.3.3 値の分布に関するパラメーター $a$

本稿では、1つのカテゴリ属性を持つ模擬データを考える。数値属性と違いカテゴリ属性の分布は、値に順序が無いいため、正規分布等のモデルで単純に特徴づけることは難しい。そこで、度数が大きい順に、値を並び替え順位をつける。そして、度数を順位の関数として表し、その関数にパラメーターを用いることとする。度数を  $y$ 、順位を  $x$  とした時に、

$$y(x) \propto x^{-a} \quad (3)$$

と近似するモデルを用いる。このようなべき乗分布は、自然界、人間界における出現頻度分布に多く見られることが知られている。ここで、 $x, y$  は正数であり、特に  $1 \leq x \leq V_{\text{att}}$  である。 $a$  がパラメーターである。ただし、

$$N_{\text{rec}} = \sum_{x=1}^{V_{\text{att}}} y(x) \quad (4)$$

という関係式が成り立つ。パラメーター  $a$  の値として、0.05, 0.1, 0.2, 0.5, 1, 2 の6個の値を取る。 $a$  が小さいほど分布は平坦で、大きいほど高順位と低順位の度数の差は大きい。

3つのパラメーター値の組合せにより、計602通りの模擬データを作成し、実験に用いる。

### 3.4 模擬データの作成

式(3)を確率関数とし、それに従う乱数を逆関数法を用いて生成し、模擬データを作成する。値域が与えられているため、各値で1レコードを生成させ、すべての値が最低でも1度は現れるようにする(つまり、 $y \geq 1$ )。その後、乱数生成により残りのレコードを生成する。

図1に作成した模擬データの例を示す。レコード数が

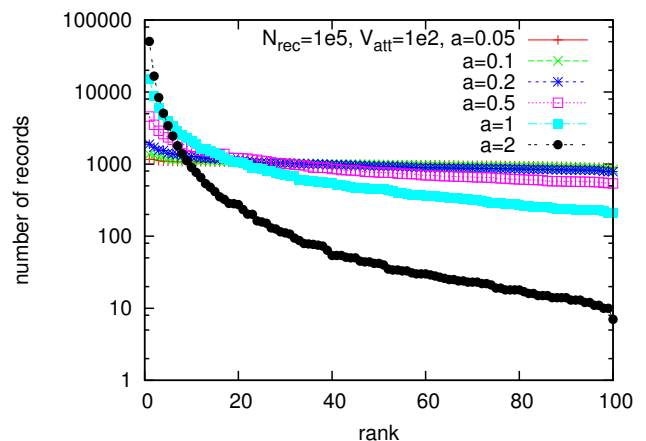


図1 レコード数が100,000、値域が100の模擬データ。横軸が度数の順位、縦軸が度数。分布のパラメーター  $a$  を6種類取っている。

100,000、値域が100の場合である。横軸が度数の順位、縦軸が度数となっている。分布のパラメーター  $a$  を6種類取っている。 $a$  が小さいほど、平坦な分布となっていることがわかる。

## 4. 実験結果

本稿では、匿名化前後の模擬データ間の  $L_1$  距離を有用性の指標として用いる。1つのパラメーター値のセットに対して、模擬データ作成と  $Pk$ -匿名化の実行を10回繰り返

し、それぞれの試行で  $L_1$  距離を測る。以降の図では、10個の  $L_1$  距離の平均値をそのパラメータ値セットの  $L_1$  距離とし、標準偏差を誤差として示す。

以下、 $L_1$  距離が、模擬データを特徴づけるパラメータを変化させた時、どのように影響を受けるか見ていく。匿名性指標  $k$  の値は  $k = 2$  で固定してある。

#### 4.1 値域 $V_{att}$ による影響

図2に、 $N_{rec} = 10^3, 10^4, 10^5, 10^6, 10^7$  の場合の、匿名化前後のデータの  $L_1$  距離を示す。縦軸が  $L_1$  距離である

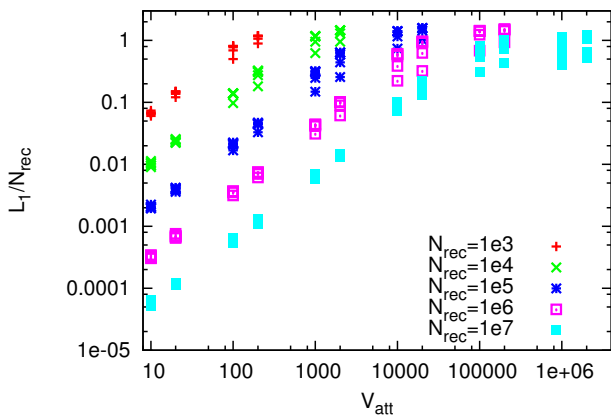


図2  $N_{rec} = 10^3, 10^4, 10^5, 10^6, 10^7$  の場合の、匿名化前後のデータの  $L_1$  距離。縦軸は、レコード数  $N_{rec}$  で正規化した  $L_1$  距離。横軸は、値域  $V_{att}$ 。見やすさのため、誤差は示さない。

が、レコード数  $N_{rec}$  で正規化している。横軸は値域  $V_{att}$  である。レコード数を固定すると、値域が小さいほど  $L_1/N_{rec}$  は小さくなる。これはあるレコード数に対して、値域が小さいほど維持確率が大きくなるためである。 $N_{rec}, V_{att}$  が与えられた時に、わずかにバラつきがある。これは、分布パラメータ  $a$  の違いによるものである。

#### 4.2 分布のパラメータ $a$ による影響

図3には、分布のパラメータ  $a$  の影響を示す。レコード数が  $N_{rec} = 10^6$ 、レコード数と値域の比が  $V_{att}/N_{rec} = 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}$  の結果を示している。分布のパラメータ  $a$  の影響は比較的小さいことが見て取れる。ただし、値域が大きい場合は、 $a$  が大きいと  $L_1$  距離が下がる傾向にある。

### 5. 考察

#### 5.1 $L_1$ 距離とパラメータの相関係数

データを特徴づけるパラメータと  $L_1$  距離との相関関係を評価するために、Spearman の相関係数を求めた。結果は、表4の通りである。 $L_1$  距離は、値域  $V_{att}$  と強い相関があることがわかる。レコード数  $N_{rec}$  とは弱い相関が

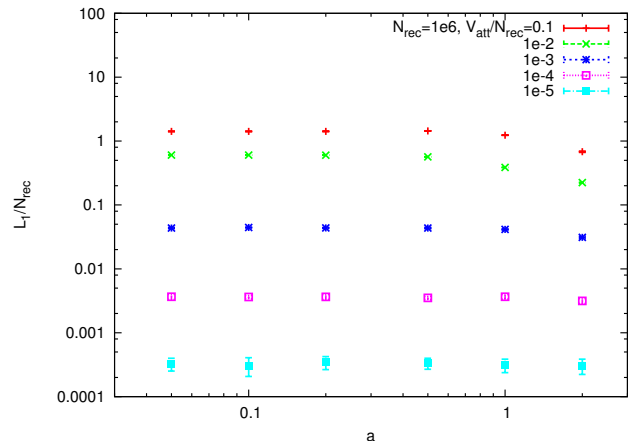


図3  $N_{rec} = 10^6$  の場合の、匿名化前後のデータの  $L_1$  距離。縦軸は、レコード数  $N_{rec}$  で正規化した  $L_1$  距離。横軸は、分布のパラメータ  $a$ 。レコード数と値域の比  $V_{att}/N_{rec}$  は、 $10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}$  の場合を示す。

表4  $L_1$  距離とデータを特徴づけるパラメータの相関係数。値域  $V_{att}$  と強い相関がある。

パラメータ	相関係数
レコード数 $N_{rec}$	-0.25
値域 $V_{att}$	0.75
分布のパラメータ $a$	-0.067

あり、分布のパラメータ  $a$  は、ほとんど相関が無い。

#### 5.2 $L_1$ 距離とパラメータの関係

図3と表4で見たように  $a$  と  $L_1$  距離との相関関係は非常に弱い。そこで、 $L_1$  距離に影響を及ぼすのは、レコード数  $N_{rec}$  と値域  $V_{att}$  の2つだと考えられる。そのためこの2つの変数を使って  $L_1$  距離をフィッティングして予測可能にすることを考える。ただし、分布パラメータ  $a$  による影響は全く無いわけではないので、これを誤差として評価する。具体的な評価法としては、 $a$  の値がいくつであっても、 $L_1/N_{rec}$  が  $X \pm Y$  の範囲に収まるような  $X$  と  $Y$  を求める。例えば、図2中のある  $N_{rec}, V_{att}$  のペアに対して

$$L_1/N_{rec} = 0.5 \pm 0.1 \quad (a = a_0) \quad (5)$$

$$L_1/N_{rec} = 0.8 \pm 0.2 \quad (a = a_1) \quad (6)$$

$$L_1/N_{rec} = 1.0 \pm 0.2 \quad (a = a_2) \quad (7)$$

という  $a$  が異なる3点があったとする。 $a$  の影響を誤差として捉えるために、 $1.0 + 0.2 = 1.2$  を上限、 $0.5 - 0.1 = 0.4$  を下限、中央値  $X$  を  $(1.2 + 0.4)/2 = 0.8$ 、誤差  $Y$  を  $(1.2 - 0.4)/2 = 0.4$  とし、3点の情報を

$$L_1/N_{rec} = 0.8 \pm 0.4 \quad (8)$$

と1点にする。こうすることで、 $a = a_0, a_1, a_2$  の誤差範囲をすべて包含する、誤差を最も大きく見積もった誤差範囲を決めることができる。この操作を行った時に、図2は、

図4となる。

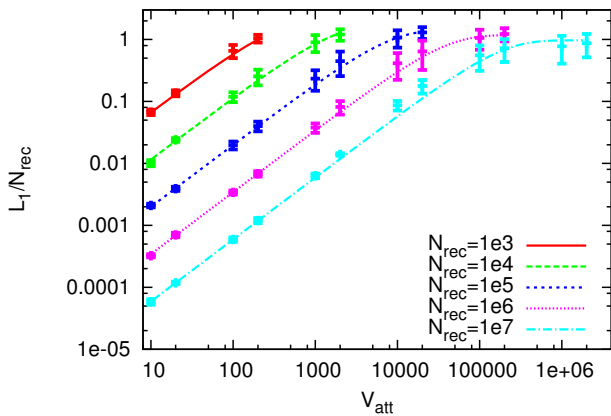


図4 分布のパラメーター  $a$  による変化を誤差としてとらえた場合の図2。線は、フィッティング関数である。

分布のパラメーター  $a$  による影響を誤差として扱った場合、 $L_1$  距離に影響を及ぼすのはレコード数  $N_{rec}$  と値域  $V_{att}$  の2つとなる。ここで、実験で得られた  $L_1/N_{rec}$  の値を  $N_{rec}$  と  $V_{att}$  の2パラメーターの関数でフィットする。フィッティング関数として次のモデルを用いる。

$$L_1/N_{rec} = A(N_{rec}) \times [1 - \exp\{-B(N_{rec}) \times V_{att}\}] \quad (9)$$

$$A(N_{rec}) = \alpha \times N_{rec}^\beta \quad (10)$$

$$B(N_{rec}) = \gamma \times N_{rec}^\delta \quad (11)$$

あるレコード数  $N_{rec}$ 、匿名性指標  $k$  に対して値域  $V_{att}$  が増加すると、維持確率は減少し、一様ランダムなランダム化に近づいていくため、 $L_1$  距離が一定値に漸近することが予想される。そこで、式(9)のように一定値に漸近する関数を用いた。

最小二乗法によるフィッティングの結果、

$$\alpha = 3.69 \pm 0.30, \beta = -0.0830 \pm 0.0066 \quad (12)$$

$$\gamma = 0.358 \pm 0.036, \delta = -0.681 \pm 0.008 \quad (13)$$

の値を得た。図4中の線は、このフィッティング関数を示しており、実験結果と良い一致を見せることがわかる。

このフィッティング関数を用いることで、実際に処理をせずとも、 $Pk$ -匿名化を施したデータの有用性の予想ができ、実用上便利である。

## 6. まとめ

本稿では、1属性のみ持つ模擬データを用いて、匿名化データの有用性検証実験を行った。模擬データ作成にあたって、レコード数、値域、データ分布といったデータの特徴付け量をパラメーター化し、値を変えることで多数の模擬データを作成し、実験に用いた。特にカテゴリ属性のデータ分布については、度数の高い順に属性値の順位を

つけ、度数を順位の関数として表し、その関数にパラメーターを用いてモデル化した。本稿で取り扱った範囲では、匿名化データの有用性指標(本稿では  $L_1$  距離)に対するデータ分布のパラメーターの影響は、他の2つのパラメーターに比べて小さいという結果を得た。この結果に基づき、レコード数と値域によって、 $L_1$  距離を近似的に表すフィッティング関数を求めた。得られたフィッティング関数を用いることで、実際に匿名化を行わなくとも、データの特徴付ける量によって、匿名化データの有用性を予測できる。

## 7. 今後の課題

### 7.1 連結化によるプライバシーリスク

本稿の有用性検証実験では、1属性のみ持つデータを取り扱ったため、 $Pk$ -匿名化の中でも、属性の連結化を用いた  $Pk$ -匿名化を使用したことに対応している。属性の連結化によって、匿名化データの有用性向上や連結された属性間の相関関係がより良く維持されることが期待できる。

一方で、連結化によるプライバシーリスクが発生する場合も考えられる。 $Pk$ -匿名性は、“ある人のレコードがどれか、確率  $1/k$  以下でしか当てられない”ことを保証するが、この匿名性だけで防ぐことができないプライバシー侵害からデータを守ることができないことも考えられる。1つの例として属性の値域といったパーソナルデータのメタ情報にはプライバシー情報は含まれていないという前提に立っているため、メタ情報に対して何も保証していない。しかし、連結化した途端にメタ情報にプライバシー情報が含まれ、そこからプライバシー情報が流出する可能性がある\*1。したがって、メタ情報にプライバシー情報が含まれないように連結化をすることが必要となる。本稿では、パーソナルデータ中の全属性を連結化することを想定したが、一部の属性のみを連結化するなどの方法が考えられる。

### 7.2 有用性検証実験の今後

模擬データ作成において、カテゴリ属性のデータ分布をべき乗型でモデル化した。本稿で扱った範囲では、データ分布と匿名化データの有用性の間の相関は非常に弱いという結果を得た。しかし、全ての実データの分布がべき乗型になる保証は無い。例えば、偏りのあるデータではステップ関数のような分布を持つことも考えられる。今後、べき乗型以外の現実的な分布の形を用いた実験を行い、匿名化データの有用性への影響をより網羅的に調べていく必要がある。

本稿では、1属性データに  $Pk$ -匿名化、あるいは複数属性を持つデータに属性の連結化を用いた  $Pk$ -匿名化を適用した場合の有用性評価実験を行った。実際は、複数属性デー

\*1 ただし、そのような場合であっても  $Pk$ -匿名性が破られているわけではないことに注意。

タに対して、連結化を用いない他の  $Pk$ -匿名化<sup>\*2</sup>を適用する場合も当然考えられる。したがって、模擬1属性データだけでなく、模擬複数属性データを用いた有用性評価実験を行うことも今後の課題の1つである。

## 参考文献

- [1] Pierangela Samarati & Latanya Sweeney, “Generalizing Data to Provide Anonymity when Disclosing Information”, PODS, 1998.
- [2] Latanya Sweeney, “k-Anonymity: A Model for Protecting Privacy”, International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, 2002.
- [3] 五十嵐大, 千田浩司, 高橋克己, “k-匿名性の確率的指標への拡張とその適用例”, CSS, 2009.
- [4] Rakesh Agrawal, Ramakrishnan Srikant & Dilys Thomas, “Privacy Preserving OLAP”, SIGMOD, 2005.
- [5] 五十嵐大, 千田浩司, 高橋克己, “数値属性における, k-匿名性を満たすランダム化手法”, CSS, 2011.
- [6] 五十嵐大, 長谷川聡, 納竜也, 菊池亮, 千田浩司, “数値属性に適用可能な, ランダム化により k-匿名性を保証するプライバシー保護クロス集計”, CSS, 2012.
- [7] 菊池亮, 五十嵐大, 千田浩司, 濱田浩気, “データ分布依存処理によって高い有用性を実現する確率的 k-匿名化”, SCIS, 2013.
- [8] 齋藤恆和, 五十嵐大, 菊池亮, 廣田啓一, 正木彰伍, “属性間の相関を考慮した攪乱再構築法の提案”, CSS, 2014.
- [9] UCI repository of machine learning databases, 1998, <http://archive.ics.uci.edu/ml/datasets/Adult>.
- [10] 千田浩司, 菊池亮, 濱田浩気, 五十嵐大, 廣田啓一, 富士仁, 高橋克己, “動的集合匿名化データの有用性評価と開示リスクに関する一考察”, CSS, 2013.
- [11] 正木彰伍, 廣田啓一, 齋藤恆和, 菊池亮, 濱田浩気, 五十嵐大, 千田浩司, “データ分布依存処理を用いた確率的 k-匿名化の有用性の検証”, SCIS, 2014.
- [12] Rakesh Agrawal & Ramakrishnan Srikant, “Privacy-Preserving Data Mining”, SIGMOD, 2000.

---

<sup>\*2</sup> 連結化を用いていないが, データ分布依存型攪乱 [7] も有用性維持が期待できる  $Pk$ -匿名化の1つである。