ユークリッド距離を用いた再識別手法と PWSCup2015の匿名加工データを用いた評価

伊藤聡志 菊池浩明

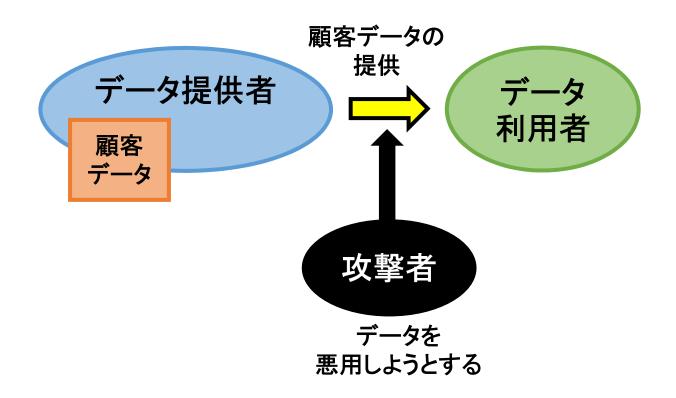
明治大学 総合数理学部 先端メディアサイエンス学科

発表概要

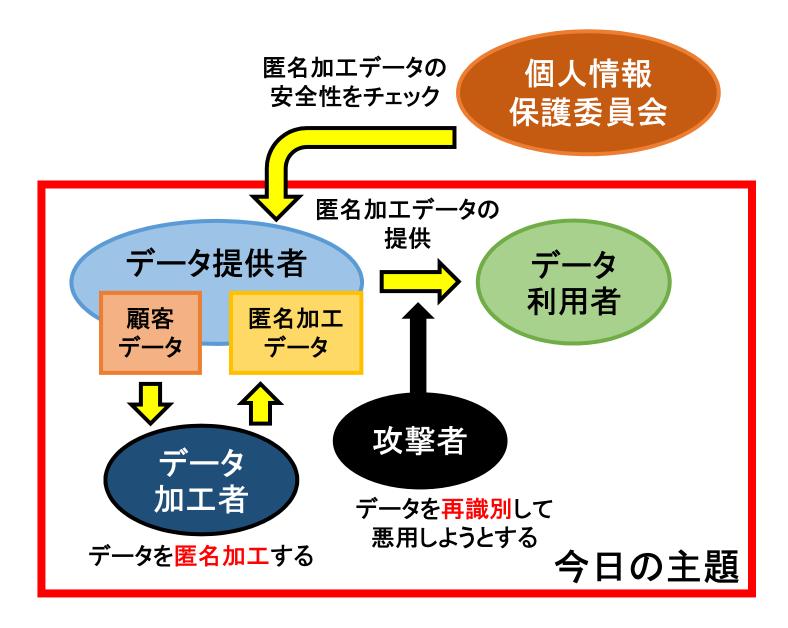
・新たな再識別手法の提案

• PWSCUP2015に提出された匿名加工データの解析

匿名加工とは?



匿名加工とは?



保護法改正とPWSCUP

2015年9月に個人情報保護法が 改正され、「**匿名加工情報**」が 定義された.

2015年10月に第1回「PWSCUP 匿名加工・再識別コンテスト」が 長崎で開催された. 第2回は秋田で開催される.



問題点と研究目的

- 1. 既存再識別手法には問題点がある.
- →既存手法の弱点を改善した新たな再識別手法を提案し、 PWSCUP2015の本戦に提出された匿名加工データを用いて 既存手法との比較を行う
- 2. PWSCUP2015の本戦に提出された匿名加工データがどのような 手法で加工されているか不明である.
- →本戦に提出された匿名加工データの解析を 小規模データを用いて行う

問題点1:既存再識別手法

「PWSCUP 2015 匿名加工・再識別コンテスト」で匿名加工 データの安全性評価に用いられた4つの再識別手法

identify.rand : ランダムに再識別を行う

identify.sa : ある1つのSAを用いて再識別を行う

identify.sort : SAの合計値をソートして再識別を行う

identify.sa21 : 特定のSAを用いて再識別を行う

問題点

再識別に用いる属性の数が少ないため、加工に弱い

問題点2:提出された匿名加エデータ

PWSCUP2015本戦には様々な企業や大学から13チームが参加し、24の匿名加エデータが提出された.

しかし、どのデータがどのような手法で匿名加工されているかは不明である。

今回は5チーム(自チーム含む)が提出した12個のデータを 研究に用いる.

データ名	作成チーム	成績
D_1,D_2	T_A (明治大学)	
D_3,D_4	T_B	2位
D_5,D_6	T_{C}	
D_7,D_8,D_9	T_D	1位
D_{10}, D_{11}, D_{12}	T_E	3位

提案手法 identify.euc

identify.euc

匿名加工データのレコードと同じQIのベクトルを持つレコードの中からSAのユークリッド距離 $D(\boldsymbol{a},\boldsymbol{b})=\sqrt{\sum_{i\in s}^n(b_i-a_i)^2}$ で再識別を行う.

	元データ					匿名	加エラ	データ			
	QI1	QI2	QI3	SA1	SA2	1111	QI1	QI2	QI3	SA1	SA2
0	2	1	1	100	100	14.142	2	1	1	110	90
X	2	1	1	200	400	322.8	2	1	1	220	390
	1	1	2	300	200	322.0	1	1	2	280	210
	1	1	2	400	500		1	1	2	390	520

既存手法との違い

既存手法 identify.sa

元データ 匿名加エデータ QI1 QI2 QI3 SA1 QI3 SA2 QI1 QI2 SA1 SA2 50 100 1 100 2 150 100 X 110 300 160 300 40 300 200 350 200 400 2 500 450 500

提案手法 identify.euc

元データ

QI1 $\mathbf{QI2}$ QI3 SA1 SA2 QI1 QI2 50 100 100 1 1 O 110 300 X 203.96 300 200 400 500 1

匿名加エデータ

QI3 SA1

150

160

350

450

1

SA2

100

300

200

500

提案手法 EUC1, EUC2

EUC1

元データ

QI1	QI2	QI3	SA1	SA2
2	1	1	100	100
2	1	1	200	400
1	1	2	300	200
1	1	2	400	500

匿名加工データ

QI1	QI2	QI3	SA1	SA2
2	1	1	110	90
2	1	1	220	390
1	1	1	280	210
1	1	1	390	520

EUC2

元データ

QI1	QI2	QI3	SA1	SA2
2	1	1	100	100
2	1	1	200	400
1	1	2	300	200
1	1	2	400	500

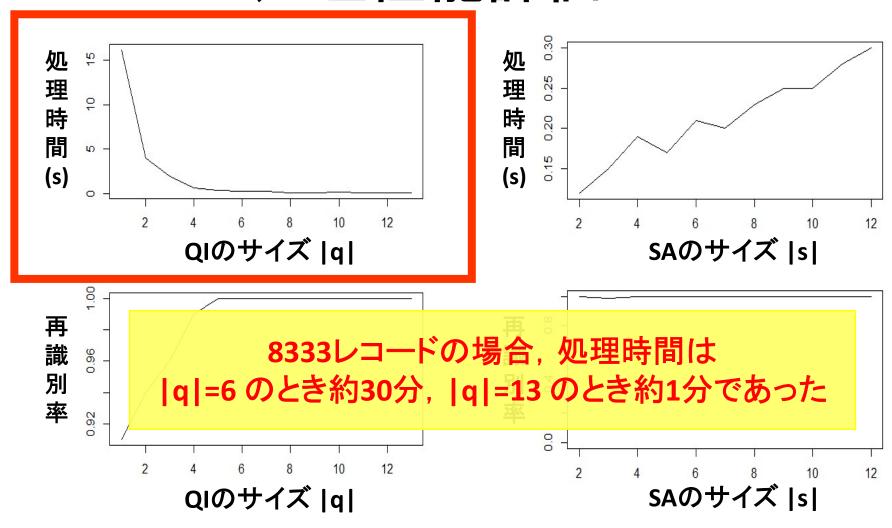
匿名加エデータ

QI1	QI2	QI3	SA1	SA2
2	1	1	110	90
2	1	1	220	390
1	1	1	280	210
1	1	1	390	520

実験結果1: 既存手法との比較(EUC1の精度)

			提案方式		
匿名加工 データ	Id-rand	Id-sa	ld-sort	Id-sa21	EUC1
D_1	0.0326	0.8238	*1.0000	0.1858	0.3010
D_2	0.6485	*0.6507	0.0012	0.0022	0.4780
D_3	0.1990	0.2412	*0.2482	0.0511	0.2070
D_4	0.1894	0.2401	*0.2526	0.0455	0.2110
D_5	0.0000	0.0223	0.0004	0.0002	*0.0743
D 1-5- 2	<u>, 0.0000</u>	C = 40.0223	0.0004	0.0002	*0.0743
東 週第	数と平均す	月 識別率	か最大で	かる014	*0.8762
Dg					*0.0011
D_9	0.0001	0.0002	0.0004	0.0000	*0.0024
D_{10}	0.0060	*0.0066	0.0001	0.0005	0.0043
D_{11}	*0.0180	0.0164	0.0001	0.0001	0.0080
D ₁₂	*0.0214	*0.0214	0.0004	0.0001	0.0080
平均	0.0931	0.1723	0.1261	0.0240	*0.1871
標準偏差	0.1741	0.2578	0.2681	0.0499	0 2426
最適数	2	3	3	0	5

EUC1の処理性能評価



処理性能評価は小規模データ(100レコード, 25属性)を用いて行った

実験2: 匿名加工データの解析

単一の手法を用いて加工した小規模匿名加工データ $D_A,...,D_H$ を用いて $D_1,...,D_{12}$ の加工手法を予測する.

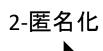
データ名	匿名加工手法
D_1	?
D_2	?
D_3	?
D_4	?
D_5	?
D_6	?
D_7	?
D_8	?
D_9	?
D_{10}	k-匿名化+SA平均化
	?
$\begin{array}{c c} D_{11} \\ \hline D_{12} \end{array}$?

データ名	匿名加工手法	加工対象
$D_{\mathbf{A}}$	k-匿名化	QI
$D_{ m B}$	SAノイズ付加	SA
D_{C}	山岡匿名化	ID
$D_{ m D}$	QI統一(対象外)	QI
$D_{ m E}$	QI統一(対象内)	QI
$D_{ m F}$	SA平均化	SA
D_{G}	QI内スワップ	SA
$D_{ m H}$	レコード削除	レコード

匿名加工手法の例

•k-匿名化

QI1	QI2	QI3	SA1	SA2
2	1	1	100	100
2	1	2	200	400
1	1	1	300	200
1	1	2	400	500



QI1	QI2	QI3	SA1	SA2
2	1	1	100	100
2	1	1	200	400
1	1	1	300	200
1	1	1	400	500

•SA平均化(ミクロアグリゲーション)

QI1	QI2	QI3	SA1	SA2
2	1	1	100	100
2	1	1	200	400
1	1	1	300	200
1	1	1	400	500



QI1	QI2	QI3	SA1	SA2
2	1	1	150	250
2	1	1	150	250
1	1	1	350	350
1	1	1	350	350

期待される効果

						匿名加	工手法	Ė		
		対象	k−匿名化	ノイズ 付加	YA	QI統一 (対象外)	QI統一 (対象内)	SA 平均化	QI内 スワップ	レコード 削除
	U1	SA	_	\(\)		1	-	1	_	×
	U2	QI,SA	×	Δ	_	1	×	1	_	×
有用	U3	QI	×	Δ	_	ı	×	_	_	×
性	U4	SA	_	Δ	_	I	_	×	×	×
1	U5	SA	_	Δ	×	-		×	×	×
		行数		_	_	_	_	_	_	×
			する				×	×	×	×
			の有				0	×	×	×
安	SA7	が対象	の攻	撃に対	付して	強い	Δ	×	×	×
安 全	E2	QI,SA		X	0		Δ	×	Δ	×
性	E3	SA	×	\(\)	0	×	×	0	0	×
	E4	SA	×	Δ	0	×	×	0	Δ	×
	EUC1	QI,SA	Δ	×	0	Δ	Δ	×	0	×

手法の組み合わせによる効果

	DA	Df	D10
加工手法	k−匿名化	SA平均化	k-匿名化+SA平均化
U1	1	1	_
U2	×	1	×
U3	×	ı	×
U4	1	×	×
U5	I	×	Δ
U6	1	1	_
S 1	0	×	0
S2	0	×	0
E1	Δ	×	Δ
E2	Δ	×	Δ
E3	×	0	0
E4	×	0	0
EUC1	Δ	×	0

実験2結果: 匿名加工データの予測結果

						匿	名加二	エデー	-タ				
		D1	D2	Dз	D4	D7	D ₅	D6	D8	D9	D10	D11	D12
	U1	_	_	_	_	_	_	_	_	_	_	_	_
有	U2	×	_	×	×	×	_	_	_	_	×	×	×
用用	U3	×	_	Δ	Δ	Δ	_	_	_	_	×	×	Δ
性	U4	_	Δ	_	Δ	Δ	Δ	Δ	Δ	Δ	×	×	×
11	U5	_	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ
	U6	_	_	_	_	_	_	_	_	_	_	_	_
	S1	×	×	Δ	Δ	Δ	×	×	×	×	0	0	Δ
	S2	×	×	Δ	Δ	0	Δ	0	0	0	0	0	0
安	E1	Δ	×	×	×	0	0	0	0	0	Δ	Δ	Δ
全	E2	×	×	×	×	Δ	Δ	Δ	0	0	Δ	Δ	Δ
性	E3	×	0	×	×	0	0	0	0	0	0	0	0
	E4	×	0	Δ	Δ	0	0	0	0	0	0	0	0
	EUC1	×	×	×	×	×	Δ	Δ	0	0	0	0	0
	DA	_	_	0	0	0	_	_	_	-	0	0	0
若	D в	_	-	1	1	-	-	_	_	1	_	_	_
名	Dc	_	-	-	-	_	0	0	0	0	_	_	_
加	D _D	_	-	-	1	0	0	0	_	1	_	_	_
エ	DE	0	_	_	1	_	_	_	0	0	_	_	_
手	DF	_	0	_	1	_		_	_	ı	0	0	0
法	DG	_	-	0	0	0	1	_	0	0	_	_	_
	Dн	_	_	_	_	_	_	_	_	_	_	_	_

実験2結果: 匿名加工データの予測結果

			匿名加工データ											
		D1	D2	Dз	D4	D7	D5	D6	D8	D9	D10	D11	D12	
	U1		_	_	_			_	_	_	_	_	_	
有	U2	×		×	×	×		_	_	_	×	×	×	
用用	U3	×				Δ		_	_	_	×	×	Δ	
性	U4		EUC	1がす	与効	Δ	Δ	Δ	Δ	Δ	×	×	×	
I I	U5					Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	
	U6		_	_	_			_	_	_	_	_	_	
	S1	×	×	Δ	Δ	Δ	×	×	×	×	0	0	Δ	
	S2	×	×	Δ	Δ	0	Δ		0	0	0	0	0	
安	E1	Δ	×	×	×	0			0,	0	Δ	Δ	Δ	
全	E2	×	×	×	×	Δ		グルー			Δ	Δ	Δ	
性	E3	×	0	×	×	0 [山岡	若名	化+他	自手法	0	0	0	
	E4	×	0	Δ	Δ	0	0	0			0	0	0	
	EUC1	×	×	×	×	×	Δ	Δ	0	0	0	0	0	
	DA		_	0	0	0		_	_	_	0	0	0	
匿	Dв		_	_	_			_	_	_		_	_	
名	Dc		_	_	_	_	0	0	0	0	ガ _ル	, —フ	3	
加	D _D		_	_	_	0	0	0	_	_				71.
エ	DE	0	_	_	_	_		_	0	K- 置	名化	+SA	平均·	1亿
手	DF		0	_	_	_		_	_	_	0	0	0	
法	DG	_	_	0	0	0		_	0	0	_	_	_	
	Dн	_	_	_	_	_	_	_	_	_	_	_	_	

実験2結果: 匿名加工データの予測結果

			匿名加工データ										
		D1	D2	Dз	D4	D7	D ₅	D6	D8	D9	D10	D11	D12
	U1	_	_	_	_	_	_	_		_	_	_	_
有	U2	×		×	×	×	_				×	×	×
用	U3	×	_	Δ	Δ	Δ	_				×	×	Δ
性	U4	_	Δ	_	Δ	Δ	Δ	Δ	Δ	Δ	×	×	×
	U5	_	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ
	U6	_	_	_	_	_	_				_	_	_
	S1	×	×	Δ	Δ	Δ	×	×	×	×	0	0	Δ
									0	0	0	0	0
安)8は	本戦	1位0)匿名	ら加コ	ロデー	-夕		0	0	Δ	Δ	Δ
	司匿名								0	0	Δ	Δ	Δ
									0	0	0	0	0
テーク	ヌ(ク)	ルー	フ2)7	が上	位に	多く見	しられ	た	0	0	0	0	0
									0	0	0	0	0
	DA	_	-	0	0	0	_	_			0	0	0
匿	D в	_	_	_	_	_	_				_	_	_
名	Dc	_	-	_	-	_	0	0	0	0	-	_	_
加	D _D	1	_	_	_	0	0	0	_	_	1	_	_
エ	DE	0	_	_	_	_	_	_	0	0	-	_	_
手	DF	1	0	_	_	_	_	_	_	-	0	0	0
法	D _G	_	_	0	0	0	_	-	0	0	_	_	_
	D н	_	_	_	_	_	_	_	_	_	_	_	_

まとめと今後の課題

まとめ

- ・既存手法の問題点を改善した、ユークリッド距離を用いる手法を提案した。12個の匿名加工データのうち、5個で最高の再識別率である。
- 提出された12個のデータで用いられていた匿名加工手 法を明らかにした 山岡匿名化と他手法を組み合わせ たデータが上位に多い

・ 今後の課題

- identify.eucのさらなる改善
- 新たな再識別, 匿名加工手法の開発

• 謝辞

・データを提供していただいた企業と大学の皆様に感謝いたします。

山岡匿名化について

山岡匿名化

データのIDを入れ替える匿名加工手法

元データ

ID	QI1	QI2	QI3	SA1	SA2
1	2	1	1	100	100
2	2	1	1	200	400
3	1	1	2	300	200
4	1	1	2	400	500

匿名加エデータ

ID	QI1	QI2	QI3	SA1	SA2
2	2	1	1	100	100
3	2	1	1	200	400
4	1	1	2	300	200
1	1	1	2	400	500