

明治大学総合数理学部

2021 年度

卒 業 研 究

マイクロアグリゲーションを用いた k 匿名化手法の提案と評価

学位請求者 先端メディアサイエンス学科

入 沢 響

目次

第 1 章	はじめに	2
第 2 章	ヘルスケアデータ	3
2.1	データ概要	3
2.2	利用データの匿名性とリスク	4
2.3	レコード削除式匿名化手法 [2]	5
第 3 章	匿名化手法の提案	6
3.1	手法の概要	6
3.2	RMSE による有用性の変化の調査	8
第 4 章	他匿名化手法との比較	10
4.1	評価手法	10
4.2	評価結果	12
4.3	考察	14
第 5 章	おわりに	14
第 6 章	謝辞	15
参考文献		15
付録 A	偽装 Wi-Fi アクセスポイントによる現在地情報のスプーフィング攻撃の調査	16
A.1	はじめに	16
A.2	位置情報スプーフィング	17
A.3	偽装 AP に利用した MAC アドレスについて	18
A.4	AP の分類について	18
A.5	実験	19
A.6	おわりに	25
参考文献		26

第 1 章

はじめに

昨今、ビックデータの利活用が企業・医療機関・金融機関など多様な場面で盛んになっている。なかでも健康診断データは病気の罹患を予測する有効な情報と考えられる。例えば、野田らは厚生労働省と総務省の許可を得て人口動態統計死亡票を目的外利用して、茨城県に住む 92,277 人の住民健診データを分析することにより、検査項目と死亡との関係を相対リスクなどを用いて明らかにした [1]。しかし、日本では 2017 年 5 月に改正個人情報保護法が施行され、要配慮個人情報を第三者に提供する際に、データに含まれる本人の同意をあらかじめとる (オプトイン) か、個人情報の第三者提供とならないようにデータを匿名加工情報とすることが必要となった。

そこで、池上らはヘルスケア企業が取得した 10 年間の 20 万人分の健康診断データと 28 万人分のレセプトデータのから成る匿名加工情報から傷病を予測するモデルを作成し、モデルの精度の変化から匿名加工の有用性を調査した [2]。池上らは、レコードを削除して k 未満のレコードの QI (Quasi-Identifier: 順識別子) が同じ値とならない匿名化手法を適応しても平均 F 値が最大 0.02 ほどしか変化しないことを示した。だが、 k を増加させた際に多くのレコードが削除されてしまうことが問題点として挙げられる。多すぎるレコードが削除されてしまうことでデータとしての有用性が損なわれ、疾患リスク等を求めた際に多くの影響が出てしまう恐れがある。

そのため、本稿ではレコードの削除をせずとも匿名性を満たす匿名加工手法を提案することを目的とする。提案手法は、マイクロアグリゲーションを段階的に用いることで k 匿名を満たす。池上らと同様のデータを使用し、QI を性別、年齢、身長として匿名化を行い、本提案の匿名化手法と池上らの使用した匿名化手法の有用性を調査する。有用性には、加工前後の傷病に対する OddsRatio (OR) と p 値を求めた際の値を利用する。

第 2 章

ヘルスケアデータ

2.1 データ概要

健康診断データには、各個人の体重や身長等の身体的特徴 21 属性と問診結果 28 属性の計 49 属性の健康診断結果が記録されている。一方、レセプトデータには、各個人に処方された医薬品の情報が記録された医薬品レセプトデータ (21 属性) と、各個人が診断された傷病の情報が記録された傷病レセプトデータ (15 属性) の 2 種類がある。健康診断データとレセプトデータには共通の仮 ID が振られている。本稿では、レセプトデータから、3 年以内に罹患した傷病を 1 とする傷病列を用いる。データサイズを表 2.1 に示す。

表 2.1 データサイズ

レコード数	健康診断データ列	傷病列	合計列数
203,521	53	1,428	1,481

2.2 利用データの匿名性とリスク

利用データの QI を本稿では性別、年齢、身長とする。次世代医療基盤法 [3] によると、性別と年齢は複数組み合わせることで個人の特定が可能な情報であるため、QI であるとされている。そのため、個人を特定されないように加工する必要がある。また、成人の身長は、同法によると、不変性が高いため静的属性に分類されている。そのため本稿では性別、年齢と同様に一般的に人を特定が可能な情報であるとし、QI とし匿名加工を行う。

表 2.2 は、身長を QI とした際と、性別、年齢、身長を QI とした際に一意に特定されてしまうレコード数を示す。例えば、身長を未加工だと少数第一位までの連続値であるが、身長のみで 16 レコードが一意に識別されてしまう。さらに、性別、年齢も組み合わせれば、6247 レコードが一意に識別される。また、身長を少数第一位を四捨五入し整数に加工しても、性別、年齢を組み合わせれば 425 レコードも特定される。

表 2.3 は、身長を QI とした際と、性別、年齢、身長を QI とした際の、QI を同一の値とするレコード数の平均値である。例えば、個人の性別、年齢、身長が特定されてしまった際には、平均 8 レコードほどまで絞られてしまう。

以上のことから、身長には個人を特定するリスクがあると考え、本稿では QI とし、匿名加工処理をする列の対象とする。

表 2.2 QI と身長の四捨五入の有無による一意なレコード数

	身長の加工なし	身長の四捨五入後
QI = { 身長 }	16	0
QI = { 性, 年齢, 身長 }	6247	425

表 2.3 QI が同一の平均レコード数

	身長の加工なし	身長の四捨五入後
QI = { 身長 }	386.1	3700.4
QI = { 性, 年齢, 身長 }	8.1	54.0

2.3 レコード削除式匿名化手法 [2]

QI が同一のレコード数が k 以下のレコードを削除することによって k 匿名を満たす手法である。この手法は、PWSCUP2020[6] のサンプル加工で提供された匿名化手法である。レコード削除式匿名化手法はロジックの理解が容易で実装が簡単な点や、QI 列を加工せずに k 匿名性を満たせる利点があるが、 k を増やした際にレコード数が減少してしまう難点がある。

レコード削除式匿名化手法を利用して 2.1 章のデータに k 匿名を行った際のレコード数の推移を図 2.1 に示す。 $k = 100$ まで加工した際には、25% ほどのレコードが削除されてしまっている。

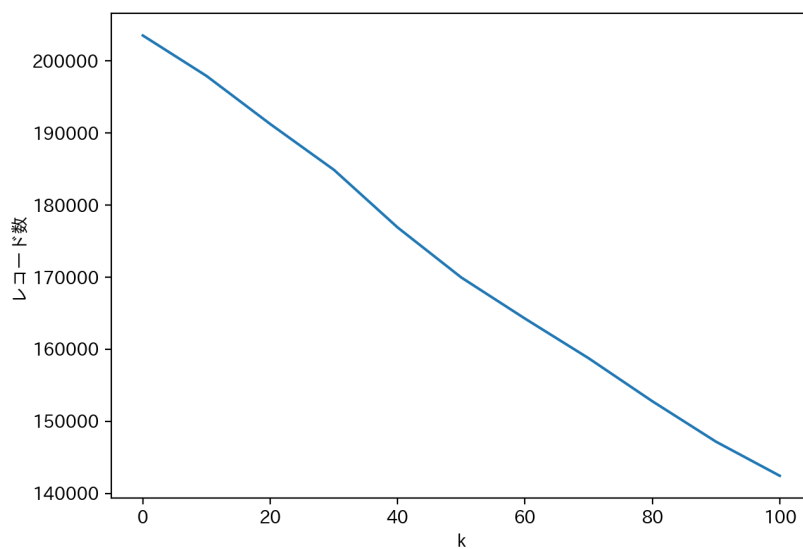


図 2.1 レコード削除匿名化によるレコード削除数

第3章

匿名化手法の提案

3.1 手法の概要

本稿では、QIを性別、年齢、身長とし、QIが同一の値のレコード数が k 以上になるような加工を k 匿名化と呼ぶ。また、身長は加工前に四捨五入で整数に丸め込む処理を予め行う。

提案手法は匿名加工手法は2段階的にマイクロアグリゲーションを使った k 匿名である。マイクロアグリゲーションとは、マイクロデータ(個別データ)を k 個のレコードを有する同質的なレコード群にグループ化した上で、そのレコードにおける個々の属性値を平均値等の代表値に置き換えることである[4]。

本加工は下記の加工①、加工②の順に2段階で加工を行うことで k 匿名を満たす。本加工を行う際のパラメータには k の他に自然数 C を用いる。本加工の数値例を図3.1に示す。

加工①では年齢列にマイクロアグリゲーションを用いた加工を行う。

1. 性別、年齢が等しいレコードを同じグループとしてグループ化する。
2. グループのレコード数が $C \times k$ 未満のグループは、性別が等しく年齢が近いグループと同じグループとしてまとめる。
3. 全てのグループのレコード数が $C \times k$ 以上になるまで2の処理を繰り返す。
4. グループごとに全レコードの年齢をグループの年齢列の平均値の小数点第一位を四捨五入した値に置換する。

加工②では年齢列にマイクロアグリゲーションを用いた加工を行う。

1. 性別、年齢が等しいレコードを同じグループとしてグループ化する。
2. グループのレコード数が $C \times k$ 未満のグループは、性別が等しく年齢が近いグループと同じグループとしてまとめる。
3. 全てのグループのレコード数が $C \times k$ 以上になるまで2の処理を繰り返す。
4. グループごとに全レコードの年齢をグループの年齢列の平均値の小数点第一位を四捨五入した値に置換する。

以上の処理の数値例を図3.1に示す。図3.1は $k=5$ 、 $C=2$ の際の処理であり、性別は全て同一のものとする。加工①によりレコード数が $C \times k$ 未満の20歳のグループが21歳のグループとまとめられ、20、21歳のレコードを含むグループの年齢が21歳に置換されている。その後、加工②によりレコード数が k 未満のグループが、身長の近いグループと同じグループにまとめられ、身長が代表値に置換されている。例えば、21

歳，身長 167cm のグループは 21 歳，身長 168cm とまとめられ，168cm に置換されている。
 本手法の特徴は以下の 2 点である。

- レコード数を保持できる。
- 加工後のデータの有用性に影響のある列を大きく変動しないように加工する。

身長(cm) \ 年齢(歳)	20	21	22
167	1*	2*	5
168	2*	2*	4*
169	2*	3*	3*
170	3*	3*	5
合計	8	10	17

身長(cm) \ 年齢(歳)	21	22
167	3*	5
168	4*	4*
169	5	3*
170	6	5
合計	18	17

身長(cm) \ 年齢(歳)	21	22
167	0	5
168	7	7
169	5	0
170	6	5
合計	18	17

図 3.1 提案匿名化手法の数値例 (*は k を満たしていないデータを示す)

3.2 RMSE による有用性の変化の調査

本手法を k を 10 から 100, C を 1 から 10 まで変化させて, 有用性を評価する. 加工前の列と加工後の列の RMSE(平均二乗偏差) で有用性を評価する. 加工前の値を y_i , 加工後の値を x_i とするとき, RMSE を

$$RMSE(y_i) = \sqrt{\frac{\sum_{i=1}^n (y_i - x_i)^2}{n}}$$

と定める. k と C の変化による身長 RMSE の変化を図 3.2 に, 年齢 RMSE の変化を図 3.3 に示す. 本加工では, 図 3.2 から身長は C の値が大きい方が, 加工後の身長の誤差が少なく, k の値を増やすほど, 加工前との身長の誤差がある傾向があることが分かる. C の値が大きいほど, $C \times k$ の値が大きくなり, 性, 年齢の総数が少ないグループの年齢が加工され, グループのレコード数が増えることで k を満たすデータが増え, 身長の加工の必要性が減るためである.

図 3.3 から年齢は C , k の値を増やすほど $C \times k$ の値が大きくなるため, 加工前との誤差があることが分かる.

安全性はグループの識別率で評価する. 識別率とは, レコードごとに QI が同一のレコード数で 1 を割った値とし, 特定される確率を表す. C , k ごとの安全性の推移を図 3.4 に示す. y 軸は識別率の平均値を示し, 値が低い方が安全性が高い. 結果は k の値の増加につれて識別率が低くなるため, 総和が低くなり安全性が高いことが分かる. だが, C による影響を強く受けなかった.

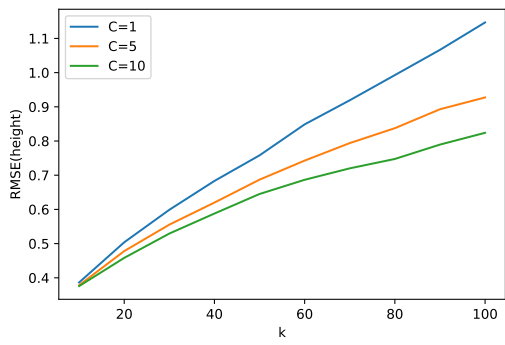


図 3.2 k についての身長 の RMSE

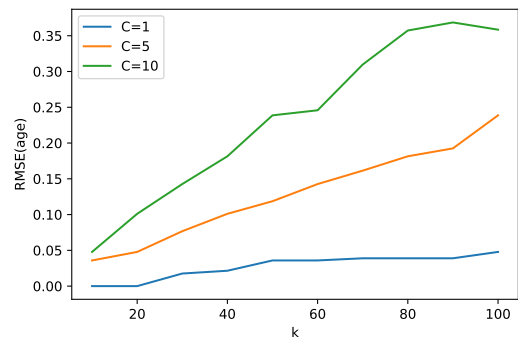


図 3.3 k についての年齢 の RMSE

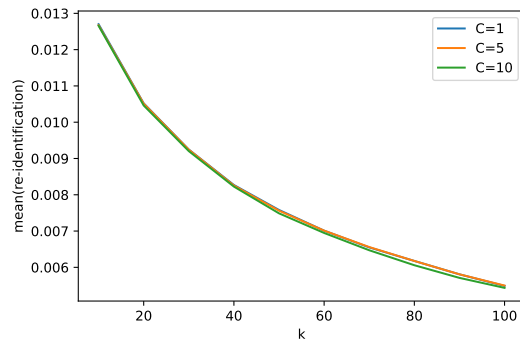


図 3.4 k についての識別率の平均値

第 4 章

他匿名化手法との比較

4.1 評価手法

本稿提案手法とレコード削除 k 匿名化手法を傷病ベクトルへの影響度で比較する。QI を含む健康診断データ 53 列を説明変数、傷病ベクトル 1428 列を一つずつ目的変数にし、ロジスティック回帰を行い、傷病に対する健康診断データの p 値と OR を算出した。本稿では、QI にしている性別、年齢、身長、傷病に対する p 値が 0.05 以下であり、罹患人数が 1000 人以上の傷病を、QI が傷病に対する影響があると定義する。QI が影響している傷病に対する p 値と OR を 4.1 に示す。

本稿提案手法とレコード削除 k 匿名化手法は、匿名加工後に、4.1 の QI の傷病についての OR と p 値を再び算出し、RMSE で有用性を評価する。RMSE が低いほど、他データに対する有用性を変化させておらず、匿名手法としては有用である。安全性に関しては、どちらも k 匿名を満たしているため同一である。

表 4.1 加工前の QI の傷病に対する OR と p 値

傷病情報		OR			p 値		
name	罹患者数	身長	性別	年齢	身長	性別	年齢
K04	30434	1.01E+00	6.51E-01	1.01E+00	1.65E-09	2.61E-72	6.25E-53
L30	29902	1.00E+00	4.96E-01	9.85E-01	2.39E-02	1.43E-185	9.15E-69
T88	22035	1.01E+00	7.46E-01	1.03E+00	4.62E-05	4.05E-27	3.56E-211
M54	21924	1.01E+00	7.43E-01	9.93E-01	3.74E-04	1.48E-27	1.35E-11
J01	16906	1.01E+00	5.57E-01	9.64E-01	5.03E-08	6.71E-79	1.20E-236
L85	15090	1.01E+00	4.84E-01	9.88E-01	1.25E-04	4.10E-110	4.61E-24
J03	13534	1.01E+00	7.20E-01	9.57E-01	2.43E-03	9.72E-22	1.59E-267
A49	13365	9.96E-01	8.14E-01	9.94E-01	3.06E-02	1.80E-09	5.60E-07
J40	13353	1.00E+00	6.93E-01	9.79E-01	1.21E-02	9.91E-27	3.51E-66
H16	13298	9.94E-01	5.06E-01	9.90E-01	2.71E-04	6.33E-87	1.52E-14
T14	12753	1.00E+00	6.84E-01	9.90E-01	2.35E-02	8.74E-28	1.07E-15
K03	12042	1.01E+00	5.70E-01	1.01E+00	2.55E-05	6.10E-56	4.29E-14
B35	11115	1.01E+00	8.96E-01	1.01E+00	1.06E-02	3.12E-03	3.20E-05
C18	10235	1.01E+00	9.16E-01	1.02E+00	1.05E-05	2.39E-02	3.38E-45
M51	9877	1.01E+00	7.85E-01	9.97E-01	8.72E-14	9.46E-10	3.85E-02
D50	9763	1.01E+00	2.91E-01	9.64E-01	6.68E-03	5.91E-176	7.85E-126
M75	8517	1.01E+00	6.88E-01	1.02E+00	4.26E-03	5.29E-19	1.34E-51
K12	8505	9.91E-01	7.30E-01	1.01E+00	3.45E-05	5.63E-14	1.54E-04
M79	8185	1.01E+00	6.62E-01	1.00E+00	1.12E-02	7.19E-22	2.99E-02
D25	7407	1.01E+00	2.08E-04	9.78E-01	2.45E-09	6.22E-49	6.80E-32
I49	6706	1.01E+00	6.80E-01	1.01E+00	2.12E-07	3.34E-16	8.36E-04
I50	6667	1.01E+00	8.16E-01	1.03E+00	2.48E-05	2.86E-05	2.38E-52
B07	6528	1.01E+00	6.17E-01	9.95E-01	1.41E-03	8.89E-24	4.12E-03
M17	5872	1.02E+00	3.96E-01	1.06E+00	1.93E-16	6.26E-75	7.58E-209
N76	5804	1.01E+00	4.90E-04	9.41E-01	5.16E-04	2.47E-64	4.67E-192
M13	5749	1.01E+00	5.83E-01	1.01E+00	5.46E-04	2.05E-26	6.92E-03
E28	5602	1.01E+00	4.80E-04	9.24E-01	3.27E-04	1.35E-64	1.06E-304
E03	5203	1.01E+00	2.99E-01	9.95E-01	1.89E-04	5.85E-109	1.17E-02
K07	4898	1.01E+00	5.87E-01	9.77E-01	1.21E-04	3.76E-22	2.75E-31
L70	4853	9.93E-01	4.51E-01	9.30E-01	8.22E-03	6.22E-44	2.31E-273
H90	4515	9.92E-01	6.99E-01	1.02E+00	4.58E-03	2.91E-10	1.49E-16
I63	4403	9.93E-01	8.79E-01	1.04E+00	9.51E-03	2.57E-02	5.35E-60
H43	4310	1.01E+00	5.34E-01	1.05E+00	2.36E-04	3.95E-27	3.38E-99
M06	4148	1.01E+00	3.49E-01	1.03E+00	7.06E-05	2.12E-68	8.40E-31
N64	3980	1.01E+00	3.48E-03	9.91E-01	1.03E-02	4.33E-138	9.52E-05
K80	3598	9.93E-01	5.38E-01	1.02E+00	3.36E-02	3.06E-22	1.62E-11
N63	3578	1.01E+00	3.03E-03	9.87E-01	2.25E-02	6.37E-112	7.55E-07
C61	3521	1.01E+00	9.66E+02	1.10E+00	1.76E-02	2.37E-22	2.54E-273
I67	3507	1.01E+00	5.66E-01	1.03E+00	2.39E-03	9.73E-19	4.84E-25
D37	3459	1.01E+00	7.88E-01	1.03E+00	1.01E-03	2.46E-04	1.26E-25
M10	3429	1.01E+00	3.72E+00	1.01E+00	2.57E-02	8.97E-55	8.46E-03
G43	3419	1.01E+00	3.35E-01	9.69E-01	6.70E-03	3.49E-59	1.15E-38
R52	3293	1.01E+00	6.13E-01	1.01E+00	5.11E-03	2.06E-13	9.58E-05
N95	3282	1.02E+00	5.99E-03	1.04E+00	2.05E-05	1.72E-200	1.05E-56
M65	3270	9.89E-01	6.24E-01	1.01E+00	1.06E-03	1.73E-12	2.16E-05
H11	3196	9.93E-01	6.54E-01	1.02E+00	2.78E-02	2.61E-10	2.28E-17
D38	3162	1.01E+00	8.47E-01	1.02E+00	5.99E-03	1.51E-02	4.01E-10
L81	3105	1.01E+00	1.35E-01	9.76E-01	3.67E-03	4.31E-144	2.11E-20
L72	3039	1.02E+00	5.64E-01	9.91E-01	2.51E-09	2.75E-16	6.59E-04
D68	2988	1.01E+00	5.78E-01	1.01E+00	1.94E-02	4.77E-15	3.42E-03
D39	2774	1.01E+00	4.09E-04	9.60E-01	4.23E-02	4.94E-28	2.07E-45
A56	2539	1.01E+00	3.16E-01	9.08E-01	2.50E-02	8.32E-47	1.19E-254
H65	2510	9.87E-01	7.03E-01	9.88E-01	3.76E-04	3.88E-06	1.65E-05
I80	2421	1.02E+00	4.64E-01	1.01E+00	7.78E-05	4.89E-23	1.69E-02
L84	2251	1.02E+00	4.22E-01	9.91E-01	5.12E-07	3.41E-26	1.96E-03
C67	2082	1.02E+00	6.62E-01	1.04E+00	1.28E-04	7.77E-07	1.55E-36
S93	2052	1.01E+00	5.03E-01	9.73E-01	1.37E-02	5.87E-16	1.29E-18
D22	1962	1.01E+00	2.63E-01	9.77E-01	2.95E-03	3.42E-50	3.42E-13
L82	1942	1.02E+00	4.91E-01	1.04E+00	3.00E-04	2.22E-16	6.26E-32
J15	1836	1.01E+00	5.90E-01	9.85E-01	1.80E-02	3.63E-09	2.83E-06
E22	1742	1.02E+00	3.97E-02	9.13E-01	1.73E-04	2.63E-142	1.80E-150
C78	1731	1.01E+00	4.47E-01	1.07E+00	1.55E-02	1.24E-18	3.20E-77
I25	1425	1.02E+00	7.96E-01	1.03E+00	8.21E-05	3.03E-02	5.12E-14
K57	1418	9.85E-01	1.36E+00	1.03E+00	1.87E-03	2.62E-03	1.27E-18
T81	1332	1.01E+00	4.73E-01	1.01E+00	4.52E-02	6.16E-13	2.68E-04
R22	1310	1.01E+00	6.05E-01	9.92E-01	1.93E-02	1.89E-06	4.60E-02
J38	1288	9.87E-01	7.09E-01	9.87E-01	1.01E-02	1.24E-03	9.25E-04
R80	1284	1.01E+00	6.45E-01	9.83E-01	1.55E-02	3.71E-05	5.52E-06
D44	1262	1.01E+00	2.18E-01	1.01E+00	1.20E-02	3.13E-43	1.76E-02
N97	1187	1.03E+00	7.74E-04	8.82E-01	1.33E-08	1.58E-23	4.24E-184
I34	1175	1.01E+00	6.37E-01	1.02E+00	1.89E-02	5.12E-05	1.35E-07
D41	1078	1.02E+00	7.64E-01	1.02E+00	2.25E-05	2.03E-02	4.78E-08
K06	1075	1.01E+00	4.86E-01	1.05E+00	2.77E-02	4.13E-10	1.67E-26
M67	1069	1.01E+00	4.22E-01	1.02E+00	4.80E-02	1.29E-13	5.63E-04
K52	1059	1.01E+00	6.17E-01	9.91E-01	2.15E-02	4.10E-05	3.58E-02

4.2 評価結果

評価結果の身長、年齢、性別の OR, p 値の RMSE の総和が最小の C を k を表 4.2 に示す. 表 4.2 の C を適応した際の k ごとの OR, p 値の RMSE の推移を図 4.1 から図 4.6 に示す. 図 4.1, 図 4.2, 図 4.3, 図 4.4 は, 身長と性別の OR, p 値の RMSE である. 身長, 性別のどちらも p 値, OR のどちらも k の値に関わらずレコード削除よりも低い RMSE を示した. 例えば, $k = 100$ で匿名化した際, レコード削除の身長の OR の RMSE が 9.2×10^{-3} であるのに対し, 提案手法では 1.2×10^{-3} と大きく差がついた. 加えて, その際の身長の p 値の RMSE もレコード削除式が 3.1×10^{-1} ほどであるのに対し, 提案手法は 3.2×10^{-2} ほどであった.

図 4.5, 図 4.6 は, 性別についての OR, p 値の RMSE である. 性別は他の QI と違い, 本手法では加工をおこなっていない. だが, 他の QI の値を加工することで性別の OR, p 値に誤差が生じた. その結果他の QI と違い, k の値により, レコード削除式よりも RMSE が高い値が示された. だが, レコード削除式と比較し, 提案手法の方が k が増えた際の増加分が少ないため, k を一定以上大きい値にした際には, 提案手法の方が RMSE の値が低くなることが予想される.

表 4.2 k ごとの最適 C

	k									
	10	20	30	40	50	60	70	80	90	100
最適 C	1	8	1	7	2	1	9	1	1	2

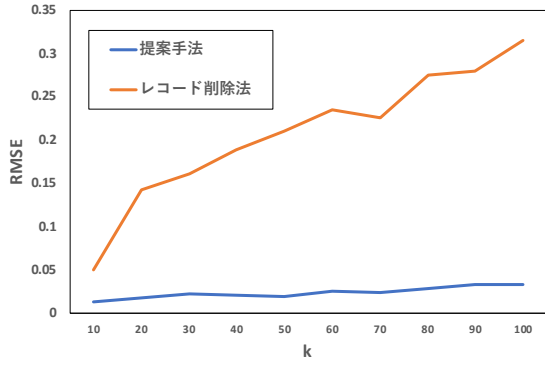


図 4.1 身長の p 値の RMSE の推移

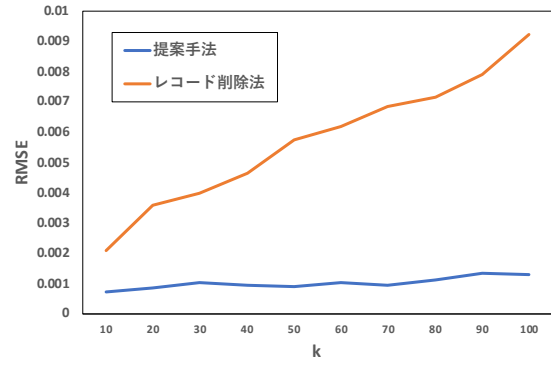


図 4.2 身長の OR の RMSE の推移

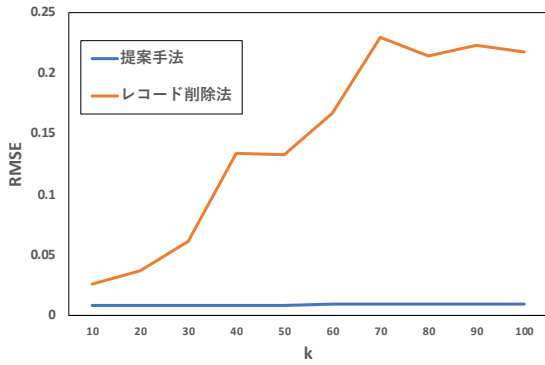


図 4.3 年齢の p 値の RMSE の推移

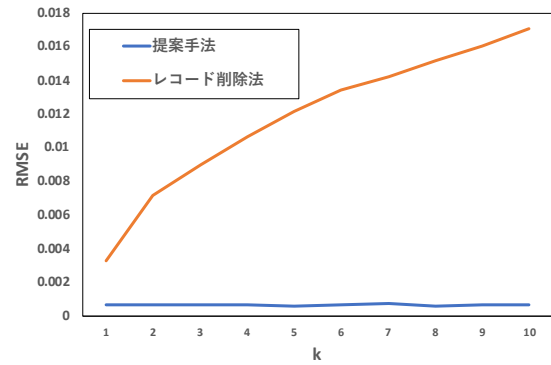


図 4.4 年齢の OR の RMSE の推移

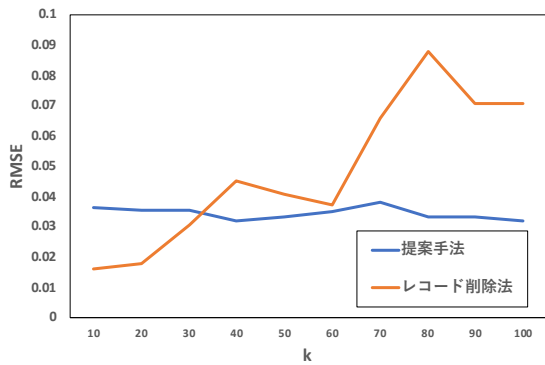


図 4.5 性別の p 値の RMSE の推移

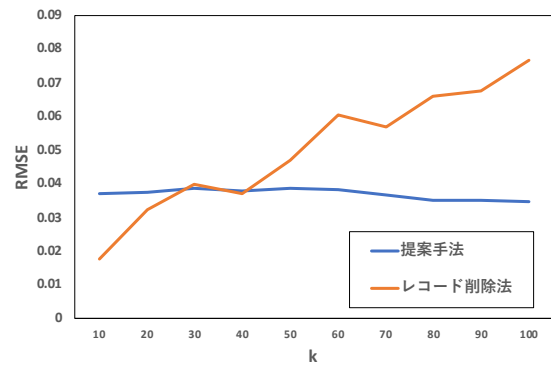


図 4.6 性別の OR の RMSE の推移

4.3 考察

結果から本稿提案手法には以下のような特徴があることが考えられる。

- 匿名性を保ちながら、レコード数を完全に保つことができる。
- 先行研究と比較し、傷病の予測の影響は低く、有用性を保つ。

第5章

おわりに

本研究では、あるヘルスケア企業が収集した 20 万人分の傷病、健康診断データを利用し、マイクロアグリゲーションを用いた匿名化手法を提案し評価を行った。その結果、先行研究と比較し有用性を下げない加工であり、 k が大きい際に特に有用な加工であると結論づけた。本研究には以下の 3 点の制約があり、それらを解消することを今後の課題とする。

- 比較対象である匿名化手法が 1 種類である
- 本稿で使ったデータに特化した匿名化手法であり、他のデータに適用するには、ロジックの調整が必要である
- 一度全通りの有用性を調査してから C を決めている

第6章

謝辞

本研究を行うにあたり、多くの方より御指導いただきました。

終始適切な助言を下さり、丁寧に指導してくださった、明治大学総合数理学部先端メディアサイエンス学科の菊池浩明教授に深く感謝申し上げます。

また、メンターの伊藤聡志さんには、多くの機会テーマや研究の進め方に関して助言をいただきました。心より感謝申し上げます。

最後に、共に励まし合い、切磋琢磨することで研究生活を有意義なものにしてくださった、菊池研究室の皆様に感謝の意を表するとともに、謝辞とさせていただきます。

参考文献

- [1] 野田博之, 磯博康, 西連地利己, 入江ふじこ, 深澤伸子, 鳥山佳則, 大田仁史, 能勢忠男, “住民健診 (基本健康検査) の結果に基づいた脳卒中・虚血性心疾患・全循環器疾患・がん・総死亡の予測”, 日本公衛誌, 53 卷 4 号, pp.265-277, 2006.
- [2] 池上, “匿名加工情報の応用 (2):各種傷病を予測する健康診断モデル”, コンピューターセキュリティシンポジウム (CSS2020), pp.26-29, 2020.
- [3] 水町雅子, “Q&A でわかる医療ビッグデータの法律と実務次世代医療基盤法・匿名加工医療情報の活用”, p.198, 2019
- [4] 伊藤伸介, 匿名化技法としてのマイクロアグリゲーションについて, 熊本学園大学, 経済論集, 15. 3:4, pp.197-232, 2009.
- [5] 伊藤聡志, 池上和輝, 菊池浩明, ”匿名加工情報の応用 (1):健康診断データとレセプトデータの分析とプライバシーリスク評価”, コンピューターセキュリティシンポジウム (CSS2020), pp.1222-1229, 2020.
- [6] PWS2020 実行委員会, PWSCup2020, ”<https://www.iwsec.org/pws/2020/cup20.html>”, 参照 2022 年 1 月 16 日

付録 A

偽装 Wi-Fi アクセスポイントによる現在地情報のスプーフィング攻撃の調査

A.1 はじめに

位置情報を利用したインターネットサービスが普及し、利用者に応じたサービスが多く存在する。

SNS やゲームなど様々な分野で利用される機会も多い。だが一方で、位置情報利用者に対し、偽装された位置情報を提示する位置スプーフィング攻撃（なりすまし攻撃）の脅威が指摘されている [1]。この攻撃を受けると誤った位置が推定され、船や自動車の自動運転などの位置情報を利用した社会インフラの混乱を引き起こす恐れがある。江藤らは文献 [2] において PC などの一部のデバイスが、周囲の Wi-Fi アクセスポイント (AP) の情報をホストサーバに送信し、位置情報を推定していることを利用し、そこに偽装された AP を複数台利用すれば、位置スプーフィング攻撃が実現できることを示している。しかし、位置情報推定に多用される GeolocationAPI のサーバへの通信は暗号化されていたため、スプーフィング攻撃の条件の詳細は不明であった。そこで、本稿では [2] の位置情報スプーフィング攻撃を再現し、次を明らかにする。

- Geolocation API に送信するパケットを観測、復号化し送信内容を明らかにすること
- スプーフィング攻撃が成立する条件について明らかにすること
- MAC アドレスを繰り返し送信することでサーバ内の DB を書き換える攻撃の効果を検証すること

A.2 位置情報スプーフィング

A.2.1 原理

GPS を搭載しないデバイスを対象とする。位置情報を推定する際、周囲の AP の情報を geolocation API サーバに送信し AP の登録された位置情報から推定する。そのため、攻撃対象のデバイスに、他の場所に設置されている AP に偽装した偽 AP を複数設置することで、別の場所を表示させることが可能である [2]。[2] のスプーフィング攻撃の手法の概要図を図 A.1 に示す。通常は、周囲の AP 登録されている位置情報から推定するが、偽 AP をある条件数以上設置すると、偽 AP の登録された位置情報に誤認する。だが、その際に必要な偽 AP の数の条件は明らかになっていない。

台湾の AP に偽装した偽 AP を、AP の少ない図 A.2 の場所に 5 台設置した際の、位置情報推定結果を図 A.3 に示す。AP の少ない場所で、別の場所の AP に偽装した AP を置くことで、実際にはいない場所に推定されている。

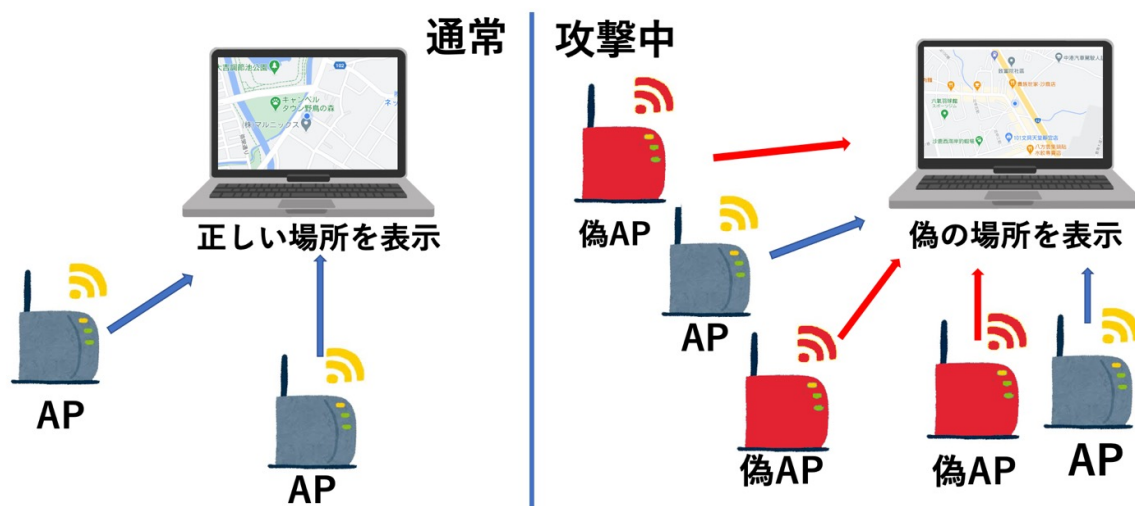


図 A.1 提案手法概要図



図 A.2 攻撃前の位置情報推定結果



図 A.3 攻撃後の位置情報推定結果 (台湾)

A.3 偽装 AP に利用した MAC アドレスについて

本稿で行うスプーフィング攻撃に利用する MAC アドレスについて述べる。2020 年 11 月 3 日, 明治大学中野キャンパス 1005 の菊池研究室にて 5 秒おきに 10 回 MacOS の `airport` コマンドを利用して観測した。平均信号強度が強い 5 つの AP の MAC アドレスを利用する。10 回のうち 1 度でも観測に失敗した AP の MAC アドレスは除く。利用した MAC アドレスを表 A.1 に示す。

表 A.1 偽装する AP に利用する MAC アドレス

観測場所	mac アドレス	平均信号強度
研究室	00:1a:eb:82:e0:**	-31
	98:f1:99:a4:7d:**	-39.3
	98:f1:99:8c:04:**	-37.9
	00:1a:eb:82:e0:**	-42.5

A.4 AP の分類について

本稿では, AP を以下の 3 種類に分類する。

- 静的 AP: サーバ内のデータベース内に存在すると予想される AP
- 動的 AP: サーバ内のデータベースにないとされる動的 AP
- 偽 AP: MAC アドレスを偽装した AP

本稿では, 静的 AP と動的 AP を観測場所で 3 日間連続観測可能かどうかにより判断する。3 日間連続で観測可能であれば, データベースに登録済みであると仮定し, 静的 AP と呼ぶ。

表 A.2 AP 分類表

状態	AP 種類		
	静的 AP	動的 AP	偽 AP
3 日間連続観測	可	不可	×
サーバ上の DB	有	無	有

A.5 実験

A.5.1 実験目的

本実験は、以下の2点を目的とする。

- スプーフィング攻撃が成立する条件を明らかにする
- 新規の攻撃手法の効果検証

A.5.2 実験環境

スプーフィング攻撃は、実験場所が攻撃対象場所となり、周囲のAPが少ない環境でなければ、攻撃不可能である。この条件を満たす本実験環境を表 A.3 に示す。

表 A.3 実験環境

実験場所	場所名称	緯度経度	平均 AP 数
実験場所 1	野鳥の森	35.914195,139.801881	5.2
実験場所 2	リユース駐車場	35.900391,139.819188	0.7
実験場所 3	自宅	35.89****,139.77****	19.7

A.5.3 事前調査

OS, デバイス, ブラウザの違いによるスプーフィング攻撃の可否

デバイス, OS, ブラウザごとに [2] のスプーフィング攻撃の可否を調査する. 偽 AP を 5 台利用し, 研究室で観測された MAC アドレスを利用し偽装する.

スプーフィング攻撃を複数回試行し, 一度でも成功すれば○, 正しい実験場所を示せば×とする. 結果を表 A.4 に示す. GPS を搭載した iPhone XS 以外のデバイスでは, OS, ブラウザに関わらず攻撃可能であった.

表 A.4 デバイス, OS, ブラウザごとの可否

デバイス	OS	chrome	IE or safari	firefox
iphoneXS	IOS 14.0.1	×	×	×
ThinkPad	Windows 10 1909	○	○	○
MacBookPro	macOS Catalina 10.15.7	○	○	○
vaio	Ubuntu 18.04	○	○	○
android	androidOS5.1.1	○	○	○

GeolocationAPI のサーバに送信するパケットの観測

Geolocation API に送信する暗号化されているパケットを, fiddler を利用し proxy サーバを経由させることにより, サーバへの送信内容を明らかにする.

Geolocation API は TLS 暗号化通信であることが必須であり, 暗号化されているためサーバに送る内容の傍受が困難である. fiddler を proxy サーバとして利用し, TLS トラフィックを復号することが可能である. 位置情報を推定したデバイスから Geolocation API のホストサーバへ送信するパケットを復号し, ホストに送信していた json ファイルを図 A.4 に示す. ホスト (www.googleapi.com) に周囲の AP の MAC アドレスと信号強度を json 形式で送信している. “macadress” は位置情報推定時の周囲の AP の MAC アドレスであり, “signalStrength” はその macadress の観測時の信号強度を表しており, “age” とは AP が検知されてからの時間 (ミリ秒) を表している.

```
{"wifiAccessPoints":  
  [{"age":0,"macAddress":"f2-a1-6d-d9-0c-***","signalStrength":-23},  
  {"age":0,"macAddress":"98-f1-99-a4-7d-***","signalStrength":-27},  
  {"age":0,"macAddress":"00-1a-eb-82-e0-***","signalStrength":-34},  
  {"age":0,"macAddress":"74-03-bd-e3-3e-***","signalStrength":-83},  
  {"age":0,"macAddress":"f8-b7-97-fc-5f-***","signalStrength":-83},  
  {"age":0,"macAddress":"fa-8f-ca-8f-f4-***","signalStrength":-83}]}
```

図 A.4 PI サーバに送信されるパケットに含まれる json データ

A.5.4 実験方法

実験 1：AP の数とスプーフィングの可否の調査

本実験では、静的 AP、動的 AP と偽 AP のそれぞれの数と [2] のスプーフィング攻撃の可否の関係を明らかにすることを目的とする。

研究室の AP に偽装した偽 AP を最大 5 台まで変動させながら、Google マップへアクセスし、現在地情報の推定を行う。その際送られた API 利用時のサーバへ送信するパケットを復号化し、送信していた AP を分類し、結果との関係性を考察する。実験場所 1、実験場所 2 において、偽 AP を増やす度に 2 度調査する。

実験 2：単一の MAC アドレスを持つ複数 AP によるスプーフィング攻撃

偽 AP の偽装する MAC アドレスを単一のものを使用した際の [2] のスプーフィング攻撃の可否を明らかにすることを目的とする。MAC アドレスは研究室で 10 回計測した際の信号強度の平均が一番高い 00:1a:eb:82:e0:**を使用する。偽 AP は 5 台利用する。実験は実験場所 1 で行う。

実験 3：架空の MAC アドレスを持つ偽 AP によるスプーフィング攻撃

偽 AP に存在しないと考えられる架空の MAC アドレスを設定した際のスプーフィング攻撃の可否について明らかにする。偽 AP は 5 台利用する。それぞれの MAC アドレスを表 A.5 の MAC アドレスに偽装し、スプーフィング攻撃を試みる。実験は実験場所 2 で行う。

表 A.5 偽 AP に設定する MAC アドレス

AP の記号	設定 MAC アドレス
A	aa:aa:aa:aa:aa:aa
B	aa:aa:aa:aa:aa:bb
C	aa:aa:aa:aa:bb:bb
D	aa:aa:aa:bb:bb:bb
E	aa:aa:bb:bb:bb:bb

実験 4：偽位置情報登録の調査

API のサーバに偽の位置情報を繰り返し送信することによる位置情報の攻撃の可否を調査する。

実験場所 1 で観測された静的 AP の MAC アドレスに偽装した偽 AP を AP の多い実験場所 3 に設置し、1 日に 10 回位置情報推定を行い、偽 AP の情報を送る。これにより、実験場所 1 の静的 AP の登録された位置情報が実験場所 3 に登録され、本来の実験場所 1 においては、位置情報が不能となると予想される。設置期間は 2020 年 12 月 1 日から 1 週間とし、2 日おきに攻撃対象である実験場所 1 で位置情報を推定した際の状態を確かめる。

A.5.5 実験結果

実験 1

スプーフィング攻撃後の位置情報推定の結果は、正しい位置、偽装された位置、特定不可能の3通りである。図 A.5, 図 A.6, 図 A.7 にそれぞれの推定結果を示す。各 AP の数と攻撃結果を表 A.6 に、各 AP の数の関係と攻撃結果の関係を表 A.7 に整理する。スプーフィング攻撃が成功したのは、静的 AP, 動的 AP の数を偽 AP の数が上回ったときのみである。また、表 A.6 の No2, 4 から動的 AP が偽 AP, 静的 AP の数を上回ったとき、位置情報が推定不可能になることが予想される。



図 A.5 実験場所表示時 (正規の場所の時)



図 A.6 攻撃成功時 (研究室に偽装)

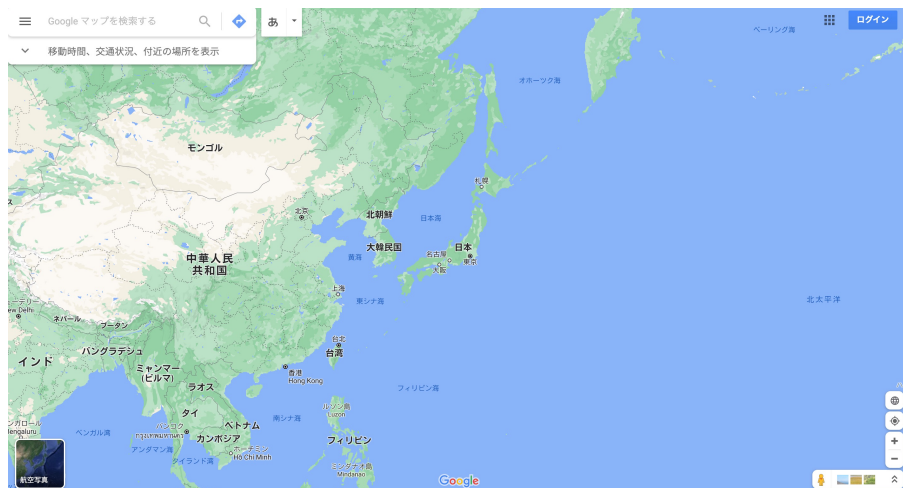


図 A.7 位置情報が特定不可能の時の表示画面

実験 2

常にスプーフィング攻撃が成功した。サーバに送信し利用される AP の MAC アドレスは種類ではなく、数に依存することが考えられる。

表 A.6 各 AP の数と結果

	No	偽 AP	静的 AP	動的 AP	結果	
実験場所 1	1	0	2	2	実験場所 1	正しい位置
	2	0	1	3	不明	特定不可能
	3	1	3	1	実験場所 1	正しい位置
	4	1	2	3	不明	特定不可能
	5	2	1	1	研究室	偽装された位置
	6	2	2	1	不明	特定不可能
	7	3	0	2	研究室	偽装された位置
	8	3	2	2	不明	特定不可能
	9	4	2	2	研究室	偽装された位置
	10	4	3	2	不明	特定不可能
	11	5	0	2	研究室	偽装された位置
	12	5	3	2	研究室	偽装された位置
実験場所 2	13	0	2	1	実験場所 2	正しい位置
	14	0	1	1	不明	特定不可能
	15	1	0	1	不明	特定不可能
	16	1	1	1	不明	特定不可能
	17	2	0	1	研究室	偽装された位置
	18	2	1	1	研究室	偽装された位置
	19	3	0	1	研究室	偽装された位置
	20	3	2	2	不明	特定不可能
	21	4	1	1	研究室	偽装された位置
	22	4	2	1	研究室	偽装された位置
	23	5	1	1	研究室	偽装された位置
	24	5	2	1	研究室	偽装された位置

表 A.7 各 AP の数の関係と攻撃結果回数

	偽装された位置 (%)	推定不明 (%)	正しい位置 (%)
偽 AP > 静的 AP, 動的 AP	12(80)	3 (20)	0 (0)
静的 AP ≥ 動的 AP, 偽 AP	0 (0)	2 (40)	3(60)
動的 AP ≥ 静的 AP, 偽 AP	0 (0)	4 (80)	1(20)
各 AP 同数	0 (0)	1(100)	0 (0)

実験 3

実験場所では位置情報が特定不可能となり、図 A.7 のように示された。存在しない MAC アドレスを利用した際は、動的 AP が多い際と同様で、位置情報が推定できない。

実験 4

結果を表 A.8 に示す。12/5 以降は、攻撃対象であった実験場所 1 では、図 A.7 のように位置情報が推定不可能であった。偽った位置へスプーフィングすることは不可能であったが、位置情報推定を不可能にする DoS 攻撃が可能であることが示された。API のサーバ内のデータベースにおいて、単一の MAC アドレスについて異なる 2 箇所の位置が登録され、その MAC アドレスによる位置情報推定が不可能となると考える。

表 A.8 偽 AP 設置日程と攻撃結果

調査日	正しい位置情報	推定不可能	正しい位置 (%)
2020/12/1(1 日目)	10	0	0 (0)
2020/12/3(3 日目)	4	6	3(60)
2020/12/5(5 日目)	0	10	1(20)
2020/12/7(7 日目)	0	10	0 (0)

A.6 おわりに

実験の結果に基づき偽 AP の数に対して静的 AP が多ければ本来の現在地を示し、動的 AP が多ければ位置情報を示さず、偽 AP が多ければスプーフィング攻撃が成功することが明らかになった。調査回数が少ないため、引き続き調査が必要である。

スプーフィング攻撃に用いる MAC アドレスは、MAC アドレスの種類に関わらず、サーバに登録されている AP の数が関係していることが考えられる。また、実験 4 の結果から、攻撃対象位置の AP に偽装した AP を 5 日以上別の場所に設置し続けると、DB が更新され、DoS 攻撃が成功することが明らかになった。

信号強度がどのように使用されているのかを明らかにすること、実験 4 の手法を利用して、AP の多い場所でも DoS 攻撃が可能か調査することを今後の課題とする。

参考文献

- [1] 片岡 義明, ”偽の電波で'GPS のなりすまし'攻撃 INTERNET Watch” (<https://internet.watch.impress.co.jp/docs/column/chizu3/1202619.html>, 2020 年 12 月参照).
- [2] 江藤 一樹, ”偽造 Wi-Fi アクセスポイントによる現在地情報のスプーフィング攻撃の脅威”, 2019 年度 明治大学卒業 論文.
- [3] googlemaps/google-maps-services-go(<https://github.com/googlemaps/google-maps-services-go/blob/master/geolocation.go>, 2020 年 12 月参照).
- [4] Fiddler を使って HTTPS トラフィックを確認する” (<https://idea.tostring.jp/?p=1077>, 2020 年 12 月参照).