

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

入沢響 †

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 †

1 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) が同じ値とされない匿名化手法を適用しても平均 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。

2 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 1 に示す。

表 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 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 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

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

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 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 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 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。

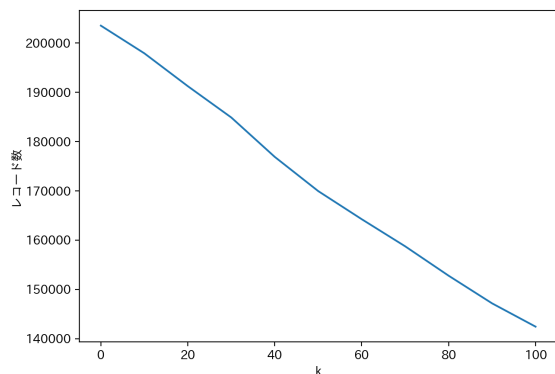


図 1 レコード削除匿名化によるレコード削除数

3 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 2 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では身長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、身長が等しいの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 k 未満のグループは、性別、年齢が等しく身長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 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身長をグループの身長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2に示す。図2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、身長に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。

加工前				加工①後				加工②後				
身長(cm)	20	21	22	身長(cm)	21	22	身長(cm)	21	22	身長(cm)	21	22
167	1*	2*	5	167	3*	5	167	0	5	168	7	7
168	2*	2*	4*	168	4*	4*	168	7	7	169	5	0
169	2*	3*	3*	169	5	3*	169	5	0	170	6	5
170	3*	3*	5	170	6	5	170	6	5			
合計	8	10	17	合計	18	17	合計	18	17			

図2 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 の変化による身長のRMSEの変化を図3に、年齢のRMSEの変化を図4に示す。本加工では、図3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4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5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。

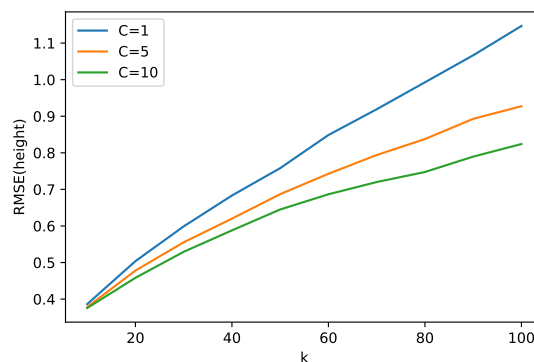


図3 k についての身長のRMSE

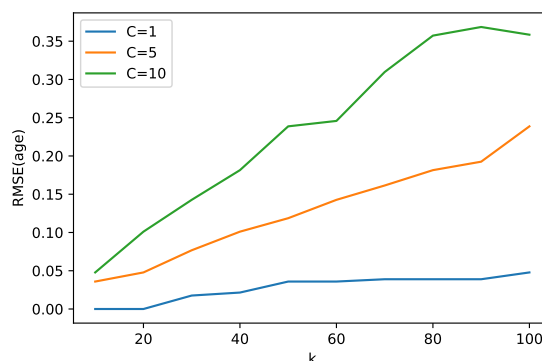


図4 k についての年齢のRMSE

4 既存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列を一列ずつ目的変数にし、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、身長の傷病に対する p 値が0.05以下であり、罹患人数が1000人以上の傷病を、 QI が傷病に対しての影響があると定義する。 QI が影響し

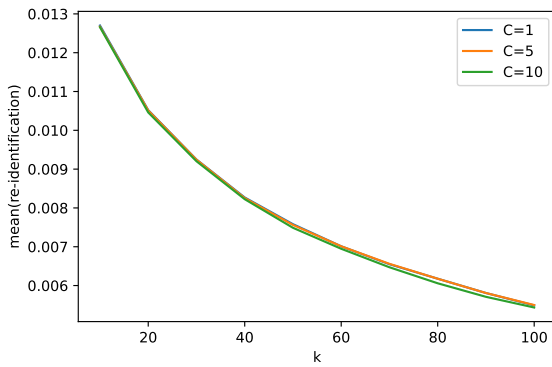


図5 kについての識別率の平均値

ている傷病に対しての p 値と OR を付録 A 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、付録 A 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。

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 に示す。表 4 の C を適応した際の k ごとの OR、p 値の RMSE の推移を図 6 から図 11 に示す。図 6、図 7、図 8、図 9 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 のどちらもが 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} と大きく差があった。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 10、図 11 は、性別についての OR、p 値の PMSE 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 k ごとの最適 C

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

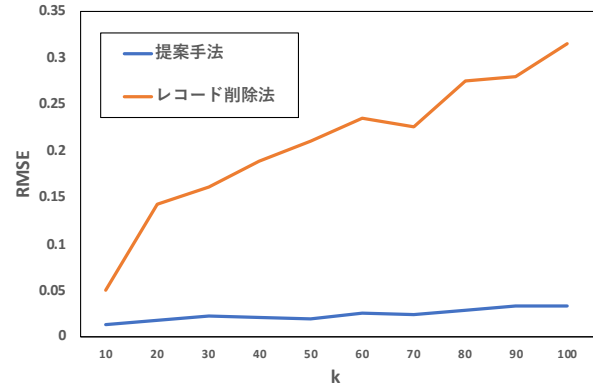


図6 身長の p 値の RMSE の推移

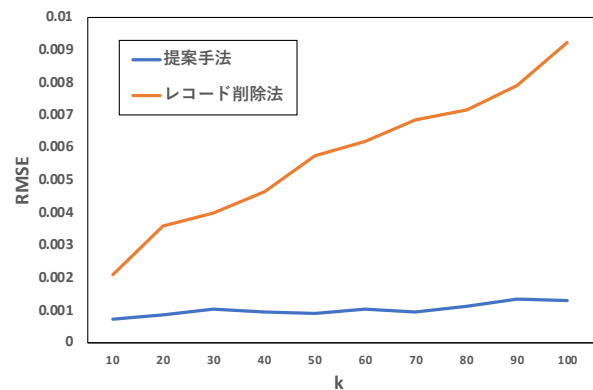


図7 身長の OR の RMSE の推移

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。

5 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

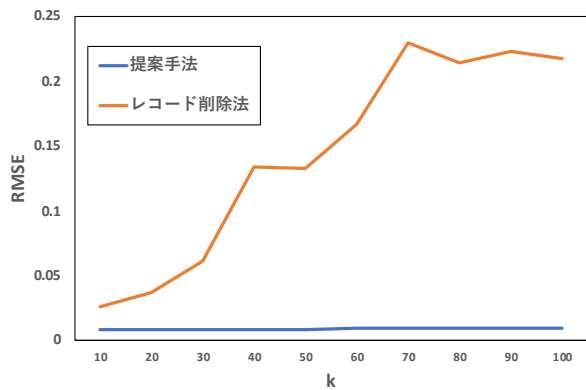


図 8 年齢の p 値の RMSE の推移

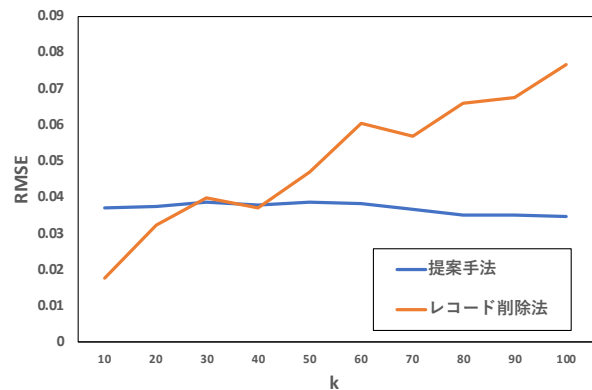


図 11 性別の OR の RMSE の推移

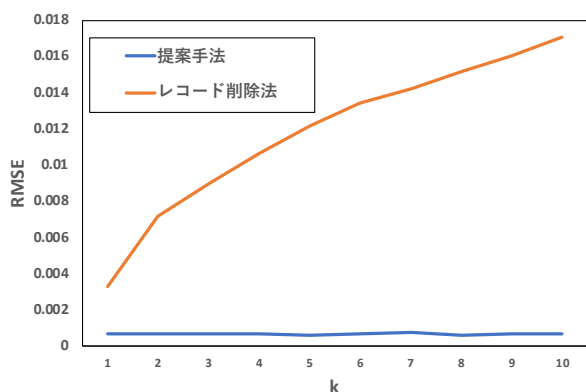


図 9 年齢の OR の RMSE の推移

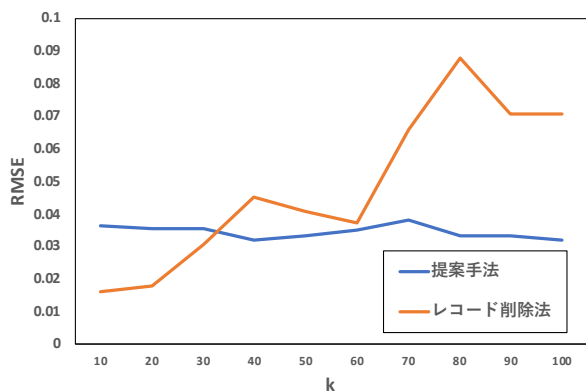


図 10 性別の p 値の RMSE の推移

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

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 巻 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 でわかる医療ビック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, PWS-Cup2020, “<https://www.iwsec.org/pws/2020/cup20.html>”, 参照 2022 年 1 月 16 日

付録 A 加工前の QI の傷病に対する OR と p 値

傷病情報		OR			p 値		
name	罹患人数	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