

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

菊池研4年 入沢 響

背景(ビックデータの活用)

- 世はビックデータの加工から新たな知見を得る時代！
- 集めたデータを**第三者提供**をすることでマーケティングできたり新たな知見を得たりしている。



他社との連携による活用



第三者提供例

匿名化とは

使用データ：性別、年齢、身長をもつデータ
悪用を試みる人：その3列で個人を特定できる人を想定する

性別	年齢	身長	病気記録
男	20	170	○○○
男	20	170	○×△
男	20	171	○○△
男	20	172	○○○
男	20	173	○△×
男	20	173	×△×
男	20	173	×△○

加工前データ

Aくんは
・性別男
・年齢20歳
・身長171センチ
だから上から3番目だ



悪用大好きくん

性別	年齢	身長	病気記録
男	20	170	○○○
男	20	170	○×△
男	20	171	○○△
男	20	171	○○○
男	20	173	○△×
男	20	173	×△×
男	20	173	×△○

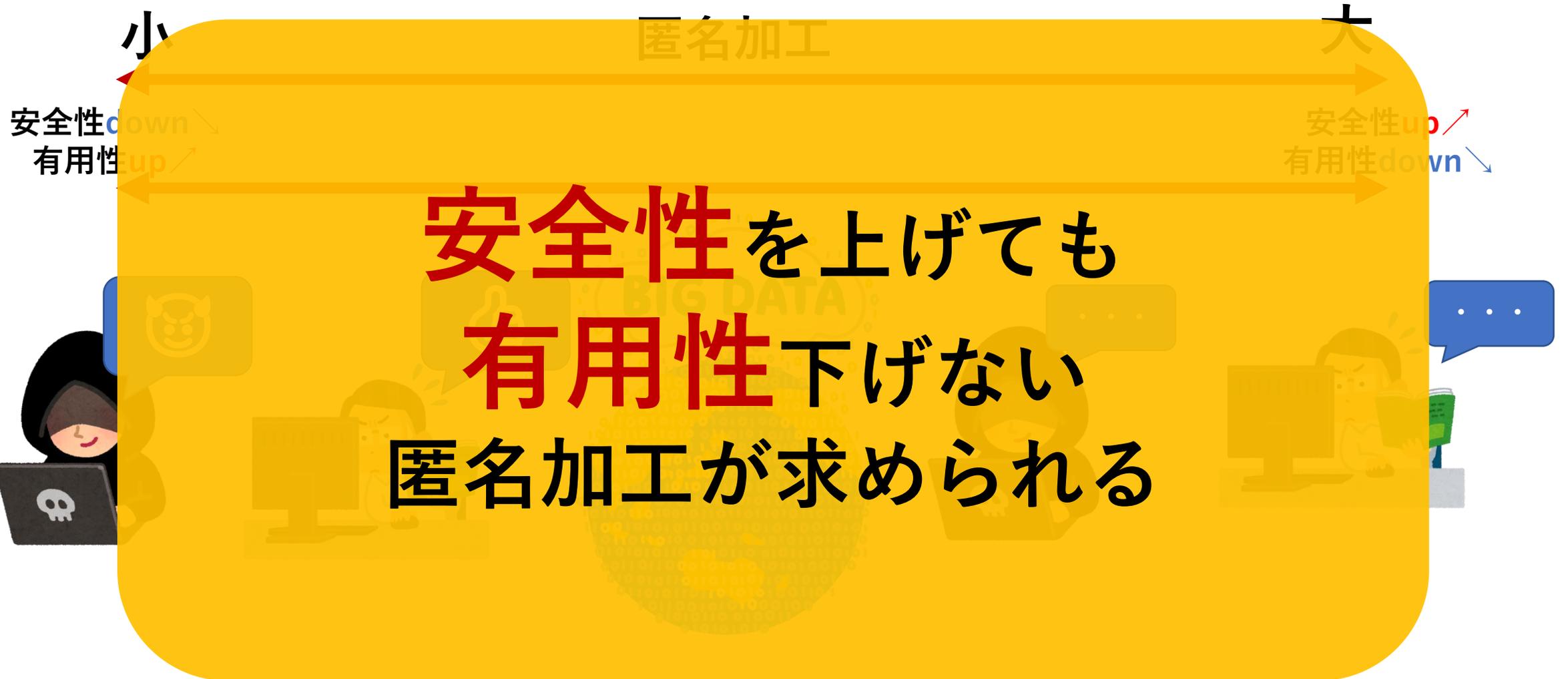
匿名加工後データ

性別、年齢、身長が同じ人が2人もいるからAくんがどれか分からないよ…



人を特定できる列が同じ値の人が常にk以上にするような加工を
k-匿名化と呼ぶ
(表のは2-匿名化)

どんな匿名化が優れているのか



先行研究の匿名化(レコード削除匿名*)

手法：識別列が同じ値がkより少ないレコードは全て削除することでk匿名を満たす

メリット：理解が容易

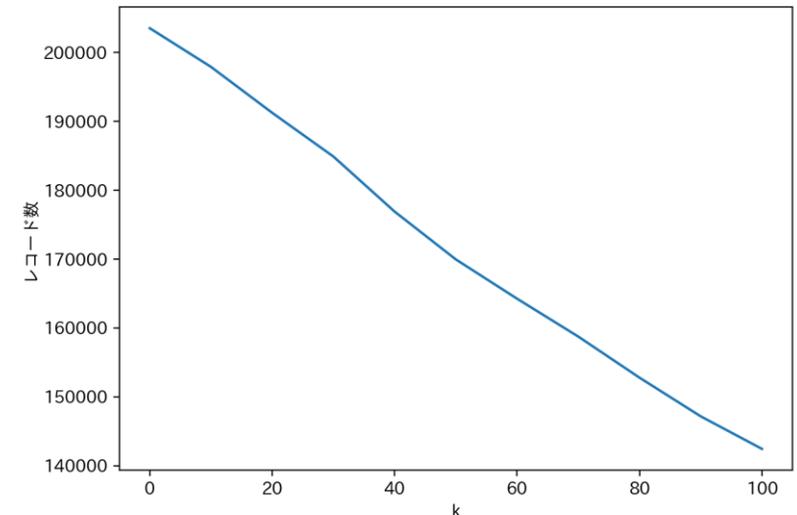
デメリット：**レコードが減少する(データとしての有用性が大きく落ちる)**

性別	年齢	身長	病気記録
男	20	169	○○○
男	20	170	○×△
男	20	171	○○△
男	20	171	○○○
男	21	170	○△×
男	21	170	×△×

→ **削除**

性別	年齢	身長	病気記録
男	20	171	○○△
男	20	171	○○○
男	21	170	○△×
男	21	170	×△×

先行研究の2-匿名化の例



Kの変化によるレコード数の変化

*PWS CUP2020 のサンプル加工で提供された匿名化手法

論文概要



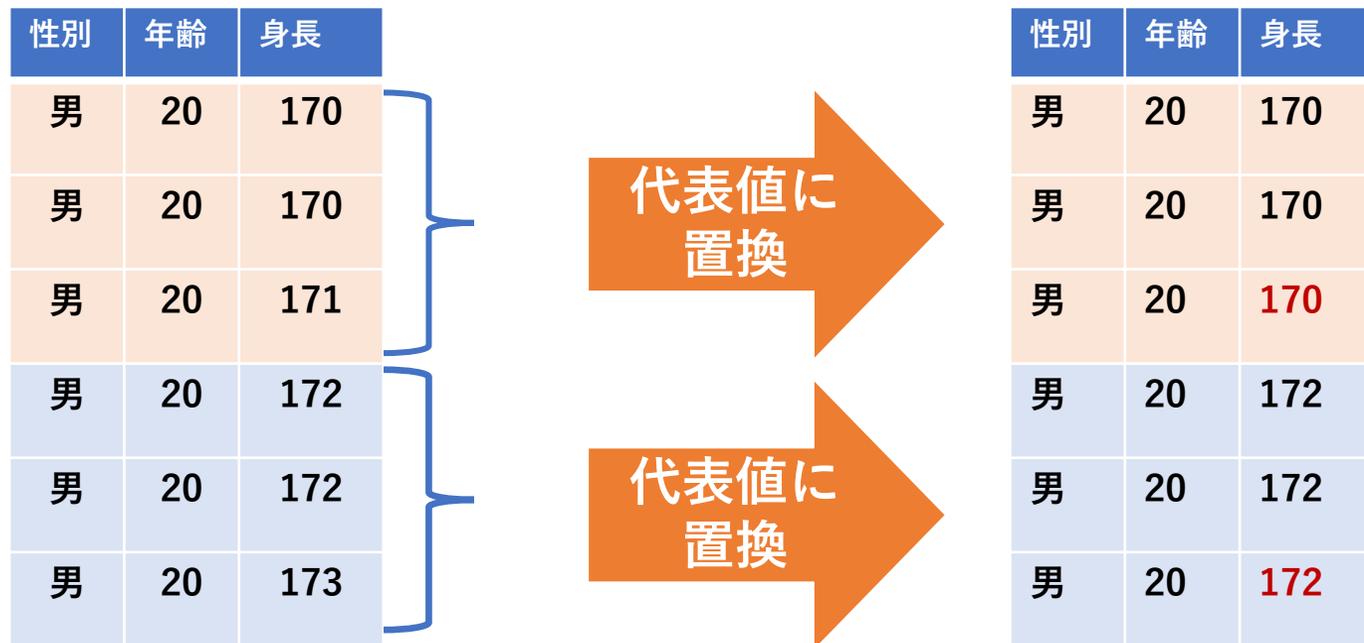
- レコードが削除されない**マイクロアグリゲーション***を2段階で用いた匿名加工手法の提案
- レコード削除匿名とともに評価し提案手法の有用性を評価

*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(2009)

マイクロアグリゲーション*とは

k個ごとにグルーピングして平均値等の代表値に置き換える→マイクロアグリゲーションによるk匿名

3-匿名化



*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(2009)

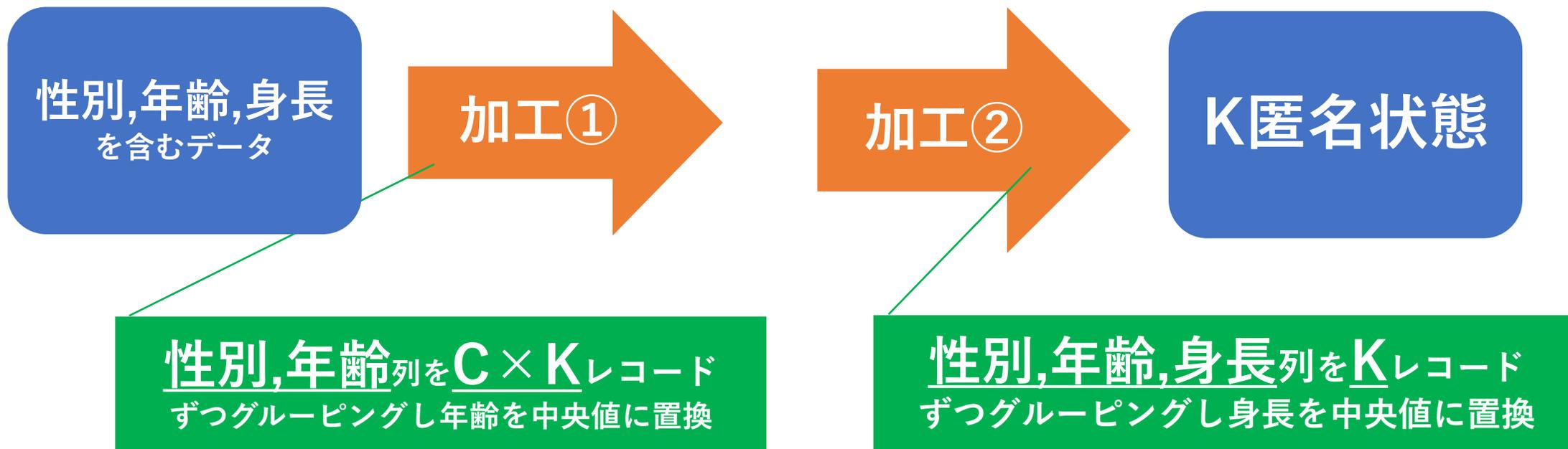
入沢流匿名化手法

加工①、加工②の2段階

パラメータ **C, K** を利用し、識別列を「性別、年齢、身長」とする。(Cは定数)

加工①…性別,年齢でマイクロアグリゲーション(年齢を加工)

加工②…性別,年齢,身長でマイクロアグリゲーション(身長を加工)



提案匿名化手法の特徴

- 加工後もレコード数を保持できる
- 使用データの中からQIを「性別、年齢、**身長**」に指定している
- Kに加えて**C(1~10の整数)**をパラメータに使用し2段階での加工をしている



Aくん(160cm)
Bくん(180cm)

成人の身長は
基本的に**不変**なため、
識別リスクが高い

C $\begin{matrix} + \\ \leftarrow \\ - \end{matrix}$

年齢加工**増** ↗
身長加工**減** ↘

年齢加工**減** ↘
身長加工**増** ↗

目的

Kごとに有用性の変化が
最も少ないCを使用することで
有用性が大きく下がることを減らす

有用性評価手法

評価手法

性別,年齢,身長の影響の傷病への変化から有用性を評価する

- 影響度は健康診断データから傷病列にロジスティック回帰した際のOR(オッズ比),p値を利用する
- 性別,年齢,身長のp値が0.05以下の傷病を利用する
- 加工前との比較を**RMSE**にて行う

性別,年齢,身長
含む

3年以内の罹患歴が
あれば1,なければ0

レコード数	健康診断データ列	傷病列	合計列数
203,521	53	1428	1481

使用データ

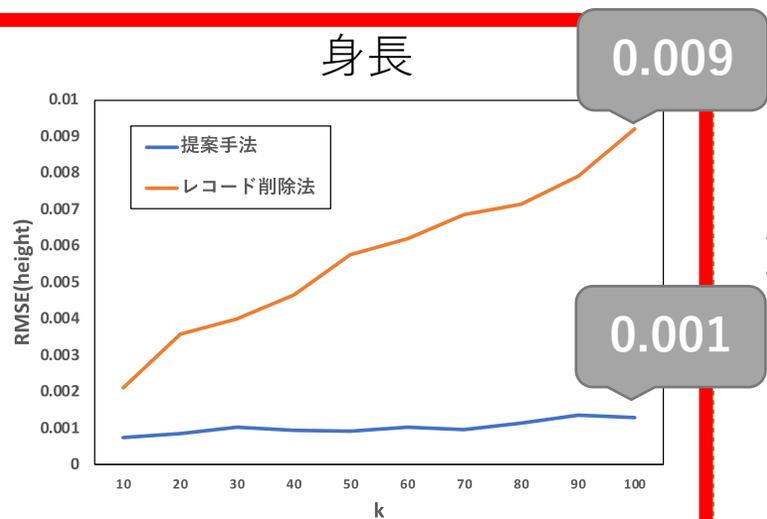
例:がんの罹患歴に対する年齢のOR

加工前	1.10
提案手法加工後	1.11
先行研究加工後	1.50

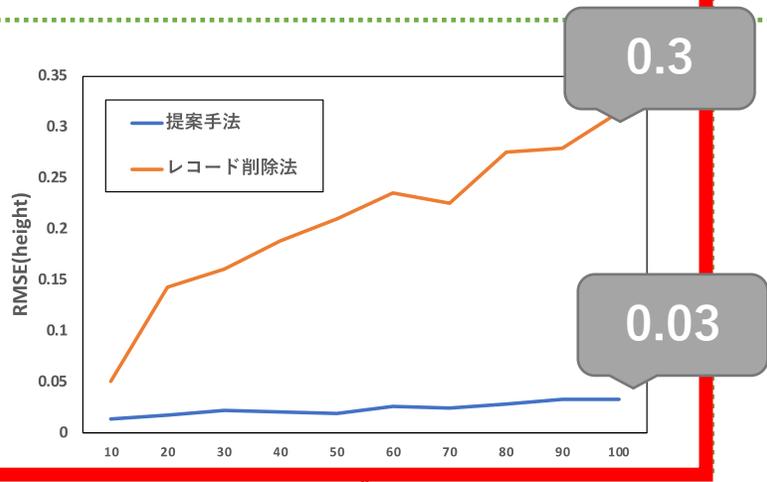
→提案手法の方が影響度の誤差が少ない

評価結果1

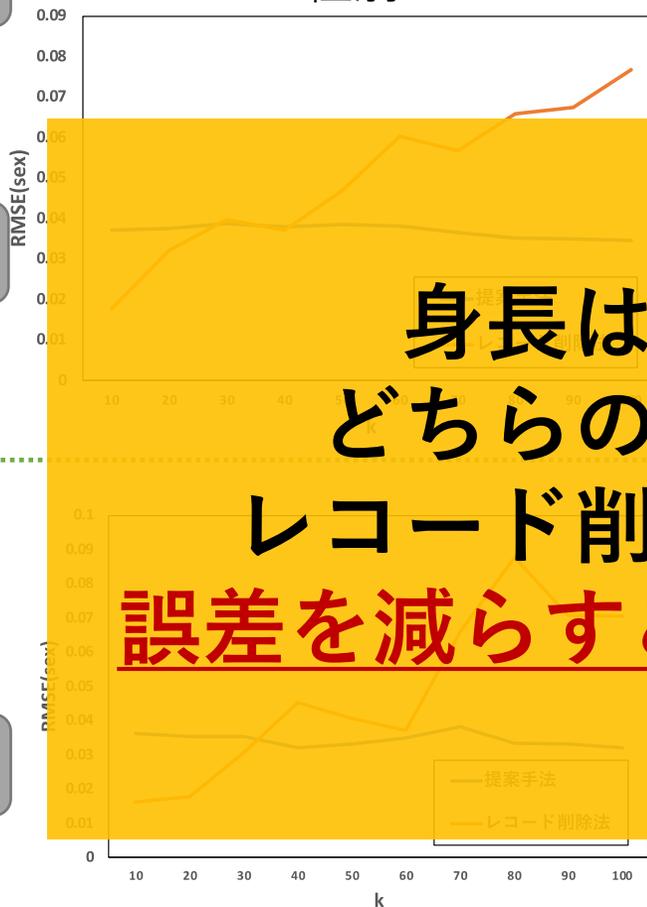
OR



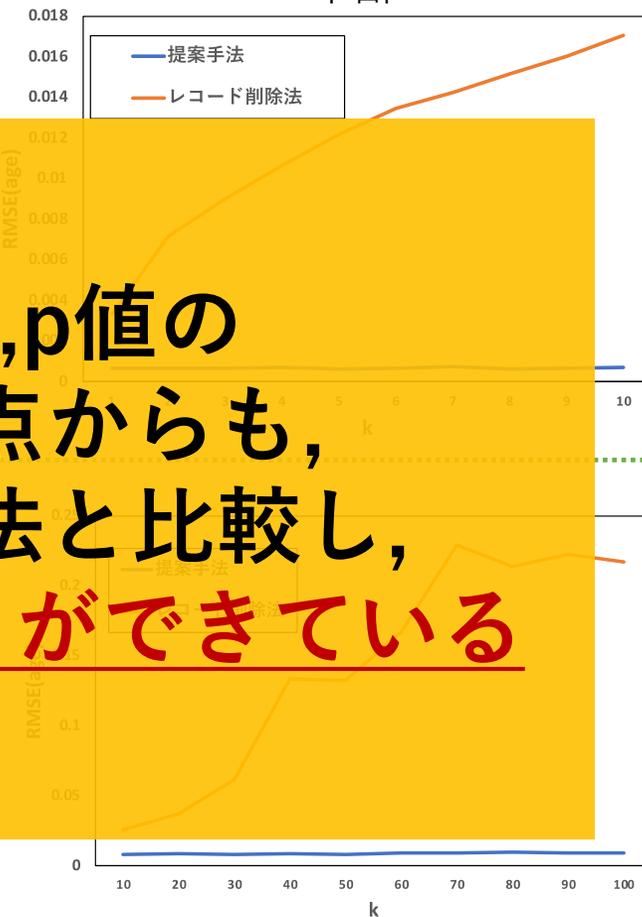
P値



性別



年齢



身長はOR,p値の
どちらの観点からも、
レコード削除法と比較し、
誤差を減らすことができています

評価結果2

OR



P値



OR,p値のどちらの観点からも、
レコード削除法と比較し、
誤差の少ない加工である

→有用性を保持できる加工である

考察と今後の課題



考察

- 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.

課題

- 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