# PWScup2024 team11

群馬大学情報学部 Gunmataro117(岡嶋佳歩,熊谷紫恩,江口誠,田中聖也,中曽根真衣,中島崚杜,矢島琴恵)

# 1.匿名化手法

#### 【方針】

クロス集計表で値が可能な限り保持されるように、ランダマイズを行う

#### 【加工方法】

- ① 10000行のレコードの中からランダムで1行選ぶ
- ② その行から更にランダムに1セル選ぶ
- ③ 下図の加エルールに従って値を変更する
- ④ ①~③を全てのレコードについて1つ以上のセルが変更されるまで繰り返

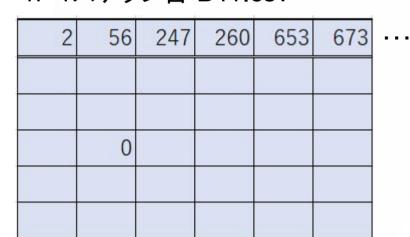
#### nステップ目 B11.csv



#### 【加エルール】

- ・nステップ目での集計表がオリジナルの集計表の値に近づく ようにratingの値を選択する
- "ZIP-code"など基本属性には重みづけをする

n+1ステップ目 B11.csv



"ZIP-code"×"56"の集計表の例







#### 【メリット】

- •集計表を使った分析に有用
- それぞれのセルの変更は独立してランダマイズされるため、レコード内のセル同 士の関連を断ち切れる

#### 【改善点】

- 重みづけの値を最適化する必要はある
  - 今回はセンシティブなデータ(ZIP-code -> Occupation -> Age)に それぞれ重みづけを行った.
  - → センシティブなデータがよりクロス集計表の値が保持される

#### ポイント

- ・クロス集計表が全部で1219個できる
  - → 全ての集計表を見ながら、
  - 46カラム×10000レコード=46万セルを変えていくので3~7日程度かかる チームメンバー全員各自のPCでプログラムを走らせた



# 2.攻擊手法

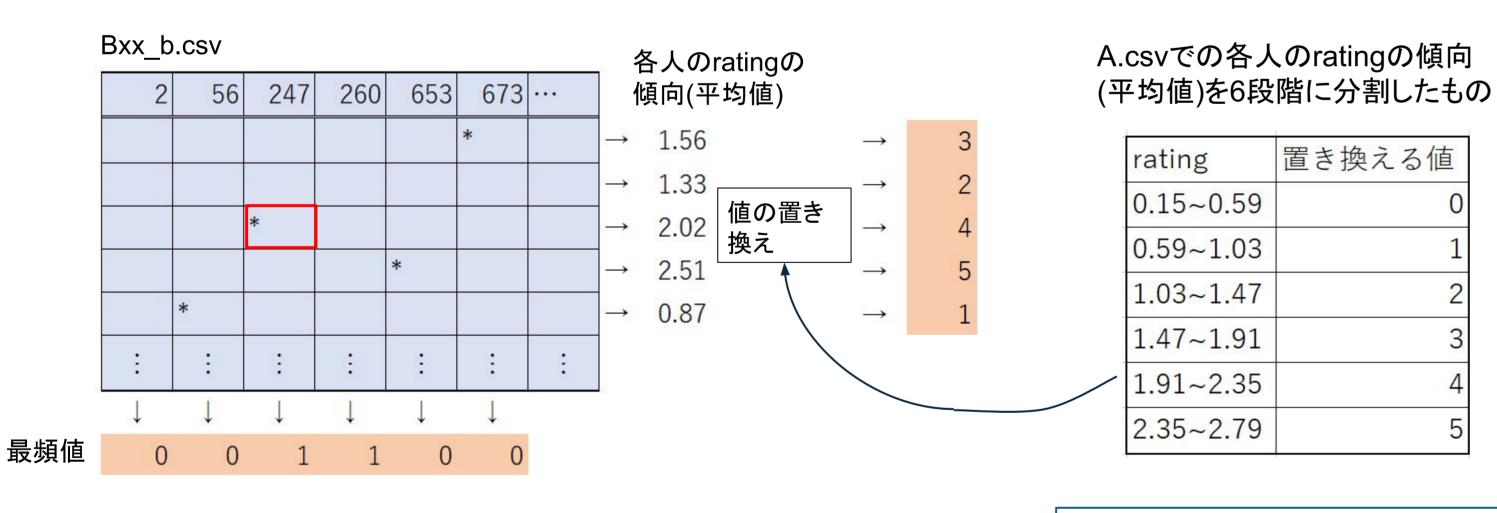
#### 【DB再構築攻撃 方針】

最頻値と、filmarksというレビューサイトのデータを用いた推定を組み合わせて行っ

※ チーム別に異なった手法を適用したが、今回は一例を紹介する

## 【方法】

- ① Bxx\_b、Cxx\_0~9について、各映画に対するratingの最頻値を出す
- ② 各映画について①の二つがどちらも0ならば黒塗りの箇所は0
- ③ それ以外の箇所は、下図のオレンジ色の値を元に推定する

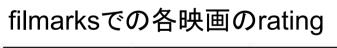


#### 【filmarksのデータ】

filmarks上での映画それぞれに対しての評価をそのまままとめた。 現実の攻撃ではこうしたレビューサイトの情報も活用するかと考えた。







映画ID	2	56	247	260	653	673	•••
評価	3.7	発見できず	3.6	3.9	3.4	3.4	•••

## 【計算例(赤枠)】

(1+4+3.6)/3=2.8666 → 3と推定

# 3.実際に分析に利用

## ↓「Age」と「ZIP-code」の組合せごとの、各映画を高評価(4と5)した人の度数

0.31 0.28

# 匿名加工前 Animation(B) 142 127

Animation(C) 241 161 230

匿名加工後

# 割合(医名加工後データ)

Animation	1	18	25	35	45	50	56					
0	0.29	0.3	0.29	0.29	0.27	0.26	0.24					
100	0.28	0.27	0.28	0.25	0.3	0.32	0.3					
200	0.25	0.31	0.28	0.3	0.28	0.29	0.25					
300	0.26	0.29	0.3	0.29	0.28	0.31	0.31					
400	0.28	0.26	0.29	0.29	0.26	0.27	0.28					
500	0.34	0.3	0.28	0.29	0.33	0.34	0.28					
600	0.32	0.23	0.3	0.27	0.29	0.31	0.28					
700	0.26	0.31	0.27	0.29	0.29	0.28	0.29					
800	0.27	0.29	0.29	0.28	0.27	0.3	0.28					

0.29 0.29 0.27

## ←Animation映画を高評価した人の割合

## 分析結果

・度数が多いわりに年齢区分が25と35の人にはアニメは人 気ではない

・ZIP-codeが「500」の人には比較的アニメが高評価の傾向 がある

# 匿名加エデータの有用性について

匿名加工後のクロス集計表【「Age」と「ZIP-code」の組合せと各映画(46作品)の高評価(4 と5)の維持が確認できた。

# 分析結果利用方法(ユースケース)

## •地域

- -地域活性化(聖地巡礼, サイン会のスケジュール)
- -映画館の再上映計画

## •年齡

- 地域の年齢分布からマーケティング
- (高齢化地域にはこんな映画を…など)
- -50代に人気 → 健康意識向上プログラムとのコラボ 50代に人気のキャラを使って... (例:運動、食生活)
- -映画の上映時間の検討
- 観客の属性によってストーリーが変わる映画
- (ロマンスが人気の属性にはロマンス展開多め等、他ジャンルとの組合せアニメ+ロマンス+コメ ディなど)

実は!アラフィフにアニメが人気!?