

2021.5.10

1.

当  $l=2$  时, 若攻击者已知 "Bob does not have heart disease", 则攻击者知道 Bob 可能患的病为频率为  $r_2, \dots, r_m$  对应的病.

当  $l>2$  时, 共有  $r_1, r_2, r_3, \dots, r_{l-1}$ , 和  $\sum_{i=l}^m r_i$  1 个频次较大的集合, 若攻击者已知  $l-2$  条形如 "Bob does not have heart disease" 的信息, 则攻击者仍不知道 Bob 患病的病属于哪个集合。但  $l-1$  条时, 则攻击者可排除 1 个集合, 则可以确定 Bob 所患的病属于哪个集合。

2. Race:  $R_1$  ZIP:  $Z_0$  Race:  $R_1$  ZIP:  $Z_1$  Race:  $R_1$  ZIP:  $Z_2$

person 94138 person 9413\* person 941\*\*

person 94138 person 9413\* person 941\*\*

person 94138 person 9413\* person 941\*\*

person 94138 person 9413\* person 941\*\*

person 94142 person 9414\* person 941\*\*

person 94142 person 9414\* person 941\*\*

person 94142 person 9414\* person 941\*\*

person 9414\* person 941\*\*

Race:  $R_0$  ZIP:  $Z_2$ 

asian 941\*\*

asian 941\*\*

asian 941\*\*

asian 941\*\*

black 941\*\*

black 941\*\*

black 941\*\*

3.

$$\forall (a) D[P, Q] = \frac{1}{8} \{$$

$$(a) \text{ For } r_1 = \frac{1}{3} - \frac{1}{9} = \frac{2}{9} \quad r_2 = r_4 = r_5 = r_6 = r_8 = r_9 = -\frac{1}{9} \quad r_3 = \frac{1}{3} - \frac{1}{9} = \frac{2}{9} \quad r_7 = \frac{1}{3} - \frac{1}{9} = \frac{2}{9}$$

$$D[P_1, Q] = \frac{1}{8} \left( \frac{2}{9} + \frac{1}{9} + \frac{2}{9} + \frac{2}{9} + \frac{1}{9} + 0 + \frac{2}{9} + \frac{1}{9} + 0 \right) = \frac{1}{6} = \frac{12}{72}$$

$$r_1 = r_2 = r_3 = r_5 = r_7 = r_8 = -\frac{1}{9} \quad r_4 = r_6 = r_9 = \frac{2}{9}$$

$$D[P_2, Q] = \frac{1}{8} \left( \frac{1}{9} + \frac{2}{9} + \frac{3}{9} + \frac{1}{9} + \frac{2}{9} + 0 + \frac{1}{9} + \frac{1}{9} + 0 \right) = \frac{11}{72}$$

$$r_2 = r_5 = r_8 = \frac{2}{9} \quad r_1 = r_3 = r_4 = r_6 = r_7 = r_9 = -\frac{1}{9}$$

$$D[P_3, Q] = \frac{1}{8} \left( \frac{1}{9} + \frac{1}{9} + 0 + \frac{1}{9} + \frac{1}{9} + 0 + \frac{1}{9} + \frac{1}{9} + 0 \right) = \frac{6}{72}$$

~~t = max~~

$$D[P, Q] = 0$$

$$t = \max_i [D[P_i, Q]] = \frac{12}{72} = \frac{1}{6}$$

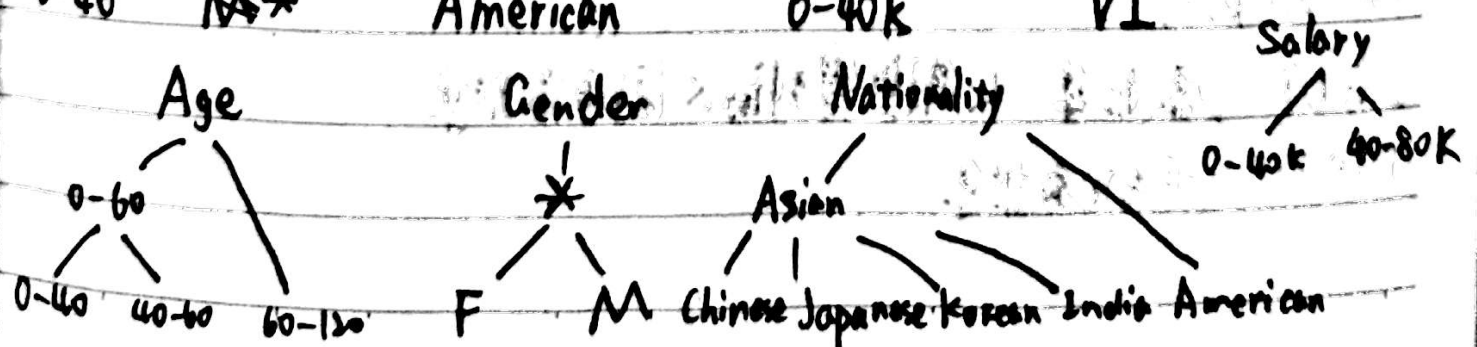
4.

(a)

Age, Gender, Nationality, Salary

(b)

Age	Gender	Nationality	Salary	Condition
0-40	F	Asian	0-40k	VI
0-40	<del>AA</del> *	American	0-40k	Flu
<del>0-60</del> 40-60	*	Asian	40-80k	HD
0-40	M	Asian	0-40k	Flu
<del>0-60</del> 40-60	*	Asian	40-80k	HD
0-40	<del>AA</del> *	American	0-40k	Flu
0-40	M	<del>American</del> Asian	0-40k	Flu
0-40	F	Asian	0-40k	Cancer
0-40	<del>AA</del> *	American	0-40k	VI
0-40	F	Asian	0-40k	Cancer
0-60	*	American	40-80k	VI
<del>0-60</del> 0-60-60	*	American	40-80k	HD
0-40	<del>AA</del> *	American	0-40k	VI



No.

Date.

the loss for attribute

Age  $\rightarrow \frac{1}{2} \rightarrow 4 \times \frac{1}{2} \times \frac{1}{13} = \frac{2}{13}$

Gender  $8 \times 1 \times \frac{1}{13} = \frac{8}{13}$

Nationality  $7 \times \frac{3}{4} \times \frac{1}{13} = \frac{21}{52}$

Salary 0

$LM = \frac{2}{13} + \frac{8}{13} + \frac{21}{52} = \frac{61}{52}$

(c)

K-Optimize( $k$ , head set  $H$ , tail set  $T$ , best cost  $c$ )

$T \leftarrow \text{PRUNE-USELESS-VALUES}(H, T)$

$c \leftarrow \text{MIN}(c, \text{COMPUTE-COST}(H))$

$T \leftarrow \text{PRUNE}(H, T, c)$

$T \leftarrow \text{REORDER-TAIL}(H, T)$

while  $T$  is non-empty do

$v \leftarrow$  the first value in the ordered set  $T$

$H_{\text{new}} \leftarrow H \cup \{v\}$

$T \leftarrow T - \{v\}$

参考 KACA 算法

KACA( $D, k$ )

INPUT:  $D$  数据集  $k$  期望得到的  $k$ -匿名数据集

OUTPUT:  $k$ -匿名的数据集

1. 由数据集  $D$  生成初始等价类, 等价类中各个元组的  $k$  维标识符的值相等
2. 若不存在元组个数小于  $k$  的等价类, 转 8
3. 否则转 3
3. 随机选择一个大小小于  $k$  的等价类  $C$
4. 计算  $C$  与其它所有等价类的距离
5. 选找到距  $C$  最近的等价类  $C'$
6. 将  $C$  和  $C'$  并为一类, 并泛化  $C$  和  $C'$
7. 转 2
8. 返回  $k$ -匿名后的数据集

注: 距离计算参考 EMD 和 Ordered Distance

泛化过程中选取  $LM$  小的泛化方式, 因数据量因  $k$  有限故复杂度不会太高.

5.

(a)

prior

$$20\% + 80\% \times \frac{1}{1001} \approx 0.2008$$

$$P(R_1(X)=0 | X=0) = 0.01 \times 20\% + (1-0.01) \times \frac{1}{1001} \times 80\%$$

$$P(R_1(X)=0 | X=i) = (1-0.01) \times 80\% \times \frac{1}{1001} \approx 0.0008$$

$$i \in \{200, \dots, 800\}$$

$$P(R_2(X)=0 | X=0) = \frac{1}{201} \approx 5 \times 10^{-3}$$

$$P(R_2(X)=0 | X=i) = 0$$

$$i \in \{200, \dots, 800\}$$

$$P(R_3(X)=0 | X=0) = 50\% \times \frac{1}{201} + 50\% \times \frac{1}{1001} \approx 3 \times 10^{-3}$$

$$P(R_3(X)=0 | X=i) = 50\% \times 0 + 50\% \times \frac{1}{1001} \approx 5 \times 10^{-4}$$

$$i \in \{200, \dots, 800\}$$

$$(b) P(R_1(X)=0) = 20\% \times \frac{1}{1001} + 80\% \times \frac{1}{1001} = 20\% \times 0.01 + 80\% \times \frac{1}{1001} \approx 2.80 \times 10^{-3}$$

$$P(R_1(X)=i) = 20\% \times \frac{1}{1001} + 80\% \times \frac{1}{1001} = 20\% \times 0.00099 + 80\% \times \frac{1}{1001} \approx 1 \times 10^{-4}$$

$$P(R_2(X)=0) = 0.01 \times \frac{1}{201} + 0.00099 \times 200 \times \frac{1}{201} \approx 1.03 \times 10^{-3}$$

$$P(R_2(X)=i) = 0.00099 \times \frac{1}{201} + 0.00099 \times 200 \times \frac{1}{201} \approx 9.9 \times 10^{-4}$$

$$P(R_3(X)=0) = 0.5 \times P(R_2(X)=0) + 0.5 \times \frac{1}{1001} \approx 1.01 \times 10^{-3}$$

$$P(R_3(X)=i) = 0.5 \times P(R_2(X)=i) + 0.5 \times \frac{1}{1001} \approx 9.9 \times 10^{-4}$$

(b)

方法(c)较好, 因为该方法的先验概率和后验概率差距较小, 能更好地保护隐私.



6.

假设  $\gamma \leq \frac{\beta}{\alpha} \cdot \frac{1-\alpha}{1-\beta}$  ~~且  $R$  满足~~ ~~且  $R$  是  $\gamma$ -amplifying BS~~

且  $\exists u \in D_U \exists v \in D_V$  s.t.  $P_f(\Psi(u_1)) \leq \alpha$  and  $P_f(\Psi(u_1) | R(u_1)=v) \geq \beta$

$$\Rightarrow P(R(u_1)=v) \leq \frac{\alpha}{\beta}$$

又有  $P_f(\Psi(u_2)) \geq 1-\alpha$  and  $P_f(\Psi(u_2) | R(u_2)=v) \leq 1-\beta$

$$\Rightarrow P(R(u_2)=v) \geq \frac{1-\alpha}{1-\beta}$$

$$\Rightarrow \gamma \geq \frac{P(R(u_2)=v)}{P(R(u_1)=v)} \geq \frac{\beta}{\alpha} \cdot \frac{1-\alpha}{1-\beta} \quad \text{矛盾 (因 } R \text{ 是 } \gamma\text{-amplifying BS)}$$

假设  $\gamma \leq \frac{\beta}{\alpha} \cdot \frac{1-\alpha}{1-\beta}$  ~~且  $R$  是  $\gamma$ -amplifying~~

且  $\exists u \in D_U \exists v \in D_V$  s.t.  $P_f(\Psi(u_1)) \geq \beta$  and  $P_f(\Psi(u_1) | R(u_1)=v) \leq \alpha$

$$\Rightarrow P(R(u_1)=v) \geq \frac{\beta}{\alpha}$$

又有  $P_f(\Psi(u_2)) \leq \frac{1-\alpha}{1-\beta}$  and  $P_f(\Psi(u_2) | R(u_2)=v) \geq 1-\alpha$

$$\Rightarrow P(R(u_2)=v) \leq \frac{1-\beta}{1-\alpha}$$

$$\Rightarrow \gamma \geq \frac{P(R(u_1)=v)}{P(R(u_2)=v)} \geq \frac{\beta}{\alpha} \cdot \frac{1-\alpha}{1-\beta} \quad \text{矛盾}$$

所以若  ~~$\gamma \leq \frac{\beta}{\alpha} \cdot \frac{1-\alpha}{1-\beta}$~~  当  $R$  是  $\gamma$ -amplifying BS

若  $\gamma \leq \frac{\beta}{\alpha} \cdot \frac{1-\alpha}{1-\beta}$ , 则  $R$  满足  $(\alpha, \beta)$ -privacy