

# Assignment-1

①

	Zip Code	Age	Nationality	Disease
1	476**	2*	*	Heart Disease
2	476**	2*	*	Viral Infection
3	476**	2*	*	Cancer
4	476**	2*	*	Cancer
5	4790*	≥40	*	Viral Infection
6	4790*	≥40	*	Heart Disease
7	4790*	≥40	*	Viral Infection
8	4790*	≥40	*	Cancer
9	476**	3*	*	Cancer
10	476**	3*	*	Cancer
11	476**	3*	*	Viral Infection
12	476**	3*	*	Heart Disease

a) Value of K:-

Group1: Zip code 476\*\* and Age 2\* - 4 records

Group2: Zip code 4790\* and Age ≥40 - 4 records

Group3: Zip code 476\*\* and Age 3\* - 4 records

Therefore  $K=4$

b) l1-diversity:-

Group1: { Heart Disease, Viral Infection, Cancer } - 3 Unique diseases

Group2: { Viral Infection, Heart Disease, Cancer } - 3 Unique diseases

Group3: { Cancer, Viral Infection, Heart Disease } - 3 Unique diseases

Here, there are 3 Unique diseases, hence the value of  $l1=3$

c) l2-diversity:-

	Heart Disease	Viral Infection	Cancer
Group1 -	1/4	1/4	2/4
Group2 -	1/4	2/4	1/4
Group3 -	1/4	1/4	2/4

$$\text{Entropy} = - \sum p_i \log p_i$$

$$\text{Group1} = - \left[ \frac{1}{4} \log\left(\frac{1}{4}\right) + \frac{1}{4} \log\left(\frac{1}{4}\right) + \frac{2}{4} \log\left(\frac{2}{4}\right) \right]$$

$$\log l = - (-1.0397) = 1.0397$$

$$l = \frac{1}{e} (1.0397) = 2.826$$

$$\text{Group 2: } -\left[\frac{1}{4}\log\left(\frac{1}{4}\right) + \frac{2}{4}\log\left(\frac{2}{4}\right) + \frac{1}{4}\log\left(\frac{1}{4}\right)\right]$$

$$\log l = -(-1.0397) = 1.0397$$

$$l = e^{1.0397} = 2.826$$

$$\text{Group 3: } -\left[\frac{1}{4}\log\left(\frac{1}{4}\right) + \frac{1}{4}\log\left(\frac{1}{4}\right) + \frac{2}{4}\log\left(\frac{2}{4}\right)\right]$$

$$\log l = -(-1.0397) = 1.0397$$

$$l = e^{1.0397} = 2.826$$

Therefore, all groups have same value  $l_2 = 2.826$

d) t<sub>1</sub>-closeness:

$$\begin{array}{l} \text{Global Probabilities of} \\ \text{HD} = 3/12 \\ \text{VI} = 4/12 \\ \text{Cancer} = 5/12 \end{array} \quad \left| \right.$$

	Heart Disease	Viral Infection	Cancer
Group 1	1/4	1/4	2/4
Group 2	1/4	2/4	1/4
Group 3	1/4	1/4	2/4

$$D_{KL}(P||Q) = \sum p(x) \log\left(\frac{p(x)}{q(x)}\right)$$

$$= \frac{1}{4} \log\left(\frac{1/4}{3/12}\right) + \frac{1}{4} \log\left(\frac{1/4}{4/12}\right) + \frac{2}{4} \log\left(\frac{2/4}{5/12}\right)$$

$$= 0 + \frac{1}{4} \log\left(\frac{3}{4}\right) + \frac{2}{4} \log\left(\frac{1}{2} * \frac{6}{5}\right)$$

$$= \frac{1}{4} \log\left(\frac{3}{4}\right) + \frac{1}{2} \log\left(\frac{6}{5}\right) = 0.01924$$

$$D_{KL}(P||Q) = \sum p(x) \log\left(\frac{p(x)}{q(x)}\right)$$

$$= \frac{1}{4} \log\left(\frac{1/4}{3/12}\right) + \frac{2}{4} \log\left(\frac{2/4}{4/12}\right) + \frac{1}{4} \log\left(\frac{1/4}{5/12}\right)$$

$$= \frac{1}{2} \log\left(\frac{3}{2}\right) + \frac{1}{4} \log\left(\frac{3}{5}\right) = 0.0750$$

$$D_{KL}(P||Q) = \sum p(x) \log \left( \frac{p(x)}{q(x)} \right)$$

$$= \frac{1}{4} \log \left( \frac{1/4}{3/12} \right) + \frac{1}{4} \log \left( \frac{1/4}{4/12} \right) + \frac{2}{4} \log \left( \frac{2/4}{5/12} \right)$$

$$= 0 + 0.25 * \log \left( \frac{3}{4} \right) + \frac{2}{4} \log \left( \frac{3}{5} \times \frac{2}{4} \right)$$

$$= 0.0192$$

$$KL\text{-Divergence} = \max(D_{KL}G_1, D_{KL}G_2, D_{KL}G_3)$$

$$= \underline{0.0750}$$

$$2) X \in \mathbb{R}^{100 \times 2}$$

$$\Rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix} \rightarrow 100 \text{ Rows}$$

$$q(x) = \frac{1}{100} x^T x$$

$$\|x_i\|_2 = 1$$

Laplacian :-

$$a) f(x) = \frac{1}{2b} \exp \left( -\frac{|x|}{b} \right); \quad b = \frac{\Delta}{\epsilon}$$

$$L_1\text{-Sensitivity} = \|f(x) - f(x')\|$$

$$q(x) = \frac{x^T x}{100} = \frac{1}{100} \sum_{i=1}^{100} x_i^T \cdot x_i$$

$$q(x') = \frac{x'^T \cdot x}{100} = \frac{1}{100} \sum_{i=1}^{100} x_i'^T x_i'$$

$$\Delta f (L_1\text{-Sensitivity}) = \|q(x) - q(x')\|$$

$$= \left\| \frac{1}{100} \sum_{i=1}^{100} x_i^T \cdot x_i - \frac{1}{100} \sum_{i=1}^{100} x_i'^T \cdot x_i' \right\|$$

$$= \frac{1}{100} (x_m^T \cdot x_m - x_m'^T \cdot x_m') \quad \left[ \because m \text{ is the value that differs from 2 datasets} \right]$$

eg:- Let  $x_m = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$   
 $x_m^T = \begin{bmatrix} 0 & 1 \end{bmatrix}$

$x_m' = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$   
 $x_m'^T = \begin{bmatrix} 1 & 0 \end{bmatrix}$

$$x_m^T \cdot x_m = \begin{bmatrix} 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$

$$x_m'^T \cdot x_m' = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\Delta f = \frac{1}{100} (x_m^T \cdot x_m - x_m'^T \cdot x_m') = \frac{1}{100} \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} = \frac{1}{100} \begin{pmatrix} -1 & 1 & 1 & 0 & 1 & 0 \\ + & 1 & 1 & 1 \end{pmatrix} = 2/100 //$$

$$\therefore l_1\text{-sensitivity} = 2/100$$

$$\text{pdf of calibrated distributed} = \frac{1}{2b} \exp\left(-\frac{|x|}{b}\right)$$

$$b = \frac{\Delta}{\epsilon} = \frac{2}{100(\epsilon)}$$

$$\Rightarrow \frac{1}{2\left(\frac{2}{100\epsilon}\right)} * \exp\left(-\frac{|x|}{\frac{2}{100\epsilon}}\right) = \frac{25}{\epsilon} * \exp\left(-\frac{|x|50}{\epsilon}\right)$$

here 1-differential privacy  $\epsilon=1$

$$\Rightarrow \frac{25}{1} * \exp(-50|x|)$$

$$\begin{array}{ll} \text{b) Let } x_1 = [1 \ 0] & x_1 = [0 \ 1] \\ x_2 = [0 \ 1] & x_2 = [0 \ 1] \\ \underbrace{\hspace{2cm}}_{\text{dataset } x} & \underbrace{\hspace{2cm}}_{\text{dataset } x'} \end{array}$$

$$q(x) = \frac{1}{100} X^T X$$

$$= \frac{1}{100} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= \frac{1}{100} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$q(x') = \frac{1}{100} X'^T X$$

$$= \frac{1}{100} \begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}$$

$$= \frac{1}{100} \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\Delta_{2f} = \|q(x) - q(x')\|_2 = \frac{1}{100} \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} = \frac{1}{100} \sqrt{1^2 + 1^2} = \frac{\sqrt{2}}{100}$$

$$\text{Gaussian pdf} = \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

$$\sigma = \frac{\frac{\epsilon \Delta_{2f}}{\epsilon}}{\epsilon} = \frac{\sqrt{2 \ln\left(\frac{1.25}{8}\right)} \cdot \frac{\sqrt{2}}{100}}{\epsilon} = \frac{\sqrt{4 \ln\left(\frac{1.25}{10^{-5}}\right)}}{(1) 100} = 0.2347$$

$$\text{pdf} = \frac{1}{(0.2347) \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x-\mu}{0.2347}\right)^2}$$

$$= \frac{1}{1.0427} * e^{-\frac{x^2}{0.1101}}$$

3) Given  $S(X, p) = - \left| \sum_{i \in [n]} \text{Sign}(x_i - p) \right|$

(3)

a) Rationale of defined score function:-

For each individual, the salary is compared to median values

- i) if  $x_i > p$  then  $\text{Sign}(x_i - p) = 1$
- ii) if  $x_i = p$  then  $\text{Sign}(x_i - p) = 0$
- iii) if  $x_i < p$  then  $\text{Sign}(x_i - p) = -1$

This can be explained as

- If  $p$  is too small (many salaries greater than  $p$ )  $\Rightarrow$  the  $\text{sign}(x_i - p)$  is positive
- If  $p$  is too large (many salaries less than  $p$ )  $\Rightarrow$  the  $\text{sign}(x_i - p)$  is negative
- If  $p$  is the median, the  $\text{sign}(x_i - p)$  is close to zero,

The theme is to minimize  $\text{sign}(x_i - p)$ , which corresponds to find a candidate  $p$  that balances the no. of salaries greater than  $p$  and no. of salaries less than  $p$ . This is exactly the property of Median.

Optimal Values

- When at Median:- If  $p$  is the median, then approximately half the salaries  $x_i > p$ , half salaries  $x_i < p$ . The  $\text{sum}(x_i - p)$  will be close to '0'.
- When not a Median:- If  $p$  is not a median, the  $\text{sum}(x_i - p)$  is not zero, then  $\text{sign}(x_i - p)$  will be either positive/negative, where it would be too large or too small.

The score function is minimized when  $p$  is median.

b) Sensitivity of Score functions:-

Sensitivity of score function  $\Delta S = \max_{x, x'} |S(x, p) - S(x', p)|$



Possible scenarios:-

- if  $x_i > p$  and  $x_i' < p$  then  $\text{Sign}(x_i - p) - \text{Sign}(x_i' - p) = 1 - (-1) = 2$
- if  $x_i < p$  and  $x_i' > p$  then  $\text{Sign}(x_i - p) - \text{Sign}(x_i' - p) = -1 - 1 = -2$
- if  $x_i < p$  and  $x_i' < p$  then  $\text{Sign}(x_i - p) - \text{Sign}(x_i' - p) = -1 - (-1) = 0$
- if  $x_i > p$  and  $x_i' > p$  then  $\text{Sign}(x_i - p) - \text{Sign}(x_i' - p) = 1 - 1 = 0$
- if  $x_i = p$  and  $x_i' = p$  then  $\text{Sign}(x_i - p) - \text{Sign}(x_i' - p) = 0 - 0 = 0$

Therefore, maximum value = 2

Therefore, Sensitivity value of Score function = 2

$$c) \Pr\left(\left|\text{index}(\tilde{p}, x^\uparrow) - \frac{n}{2}\right| \geq 4 \frac{\ln(m/\alpha)}{\epsilon}\right) \leq \alpha$$

The above inequality explains:-

$\tilde{p}$  - output of exponential Mechanism  
 $X$  - database  $X$

$\text{index}(\tilde{p}, x^\uparrow)$  - position of  $\tilde{p}$  in the sorted database  $x^\uparrow$

$n/2$  - index of true median  $x^\uparrow$

$m$  - range of possible salary values  
 $x_i \in [m] = \{1, 2, \dots, m\}$

$\epsilon$  - privacy budget.

$\alpha$  - failure probability

Consider a tail bound exponential mechanisms:-

$$\Pr(S(x, \tilde{p}) \leq S(x, p^*) - t) \leq \exp\left(-\frac{\epsilon t}{2\Delta S}\right)$$

Where  $S(x, p)$  is score function

$\tilde{p}$  is the output of mechanism

$p^*$  is true median.

$\Delta S = 2$ , Sensitivity of Score function

$t$  is the deviation in the score

$$\text{Let } t = 4 \frac{\ln(m/\alpha)}{\epsilon}$$

$$\Pr(S(x, \tilde{p}) \leq S(x, p^*) -$$

$$4 \frac{\ln(m/\alpha)}{\epsilon}) \leq \exp\left(-\frac{\epsilon \cdot 4 \frac{\ln(m/\alpha)}{\epsilon}}{2 \cdot 2}\right)$$

$$\Rightarrow \Pr(S(x, \tilde{p}) \leq S(x, p^*) -$$

$$4 \frac{\ln(m/\alpha)}{\epsilon}) \leq \exp\left(-\frac{4 \ln(m/\alpha)}{4}\right)$$

$$\Rightarrow \exp(-\ln(m/\alpha))$$

$$\Rightarrow \alpha/m$$

$\Rightarrow$  There are  $m$  possible values of  $p \Rightarrow p \in [m]$ . Applying the union bound over all  $m$  values the probability that any  $p$  deviates by more than  $t$  is:

$$\Pr\left(\left|\text{index}(\tilde{p}, x^\uparrow) - \frac{n}{2}\right| \geq t\right)$$

$$\leq m \cdot \frac{\alpha}{m} = \alpha$$

4) a) Given  $y_i = \begin{cases} x_i & \text{with probability } 1/2 + \nu \\ 1-x_i & \text{with probability } 1/2 - \nu \end{cases}$

$$\text{HIV} + \Rightarrow y_i = \frac{1}{2} + \nu (x_i)$$

$$\text{HIV} - \Rightarrow y_i = \frac{1}{2} - \nu (1-x_i)$$

As random response achieves differential privacy:-

$$\frac{\Pr[Y=y|x]}{\Pr[Y=y|x']} \leq e^\epsilon \quad \text{when HIV} + \Rightarrow x_i = 1 \\ \text{HIV} - \Rightarrow x_i = 0$$

2 Cases:- Case 1:-

$$\frac{\Pr[Y_i=1|x_i=1]}{\Pr[Y_i=1|x_i=0]} \\ = \frac{\frac{1}{2} + \nu}{\frac{1}{2} - \nu}$$

Case 2:-

$$\frac{\Pr[Y_i=0|x_i=1]}{\Pr[Y_i=0|x_i=0]} \\ = \frac{\frac{1}{2} - \nu}{\frac{1}{2} + \nu}$$

To determine  $\epsilon$ , we consider

$$\frac{\Pr[Y=y|x]}{\Pr[Y=y|x']} = \max\left(\frac{\frac{1}{2} + \nu}{\frac{1}{2} - \nu}, \frac{\frac{1}{2} - \nu}{\frac{1}{2} + \nu}\right) = e^\epsilon$$

$$\Rightarrow \epsilon = \max\left(\ln\left(\frac{\frac{1}{2} + \nu}{\frac{1}{2} - \nu}\right), \ln\left(\frac{\frac{1}{2} - \nu}{\frac{1}{2} + \nu}\right)\right)$$

$$\left[ \because \frac{\frac{1}{2} + \nu}{\frac{1}{2} - \nu} > 1 \text{ and } \frac{\frac{1}{2} - \nu}{\frac{1}{2} + \nu} < 1 \right]$$

$$\epsilon = \ln\left(\frac{1+2\nu}{1-2\nu}\right)$$

$\therefore$  The value of  $\epsilon$  depends on  $\nu$

As  $\nu$  increases, less privacy  $[\nu=0, \epsilon=0 \Rightarrow \text{perfect privacy}]$   
 $\nu$  decreases, more privacy  $(\nu=1/2, \epsilon=\infty \Rightarrow \text{No privacy})$

b) Given  $\bar{p} = \frac{1}{N} \sum_{i=1}^N y_i$

$$f(\tilde{p}) = \frac{1}{N} \sum_{i=1}^N f(Y_i)$$

$$\begin{aligned} \therefore \text{here } E[Y_i] &= \left(\frac{1}{2} + v\right) x_i + \left(\frac{1}{2} - v\right)(1 - x_i) \\ &= \cancel{\frac{1}{2} x_i} + v x_i + \frac{1}{2} - \cancel{\frac{1}{2} x_i} - v + v x_i \\ &= \frac{1}{2} + 2v x_i - v \\ &= \frac{1}{2} + v[2x_i - 1] = \frac{1}{2} + \frac{2v}{n} \sum_{i=1}^N x_i - 1 \\ &= \frac{1}{2} + (2p - 1)v \end{aligned}$$

$$E(\tilde{p}) = \frac{1}{2} + (2p - 1)v$$

$$\frac{E(\tilde{p}) - 1/2}{v} = (2p - 1) \Rightarrow \left[ \frac{E(\tilde{p}) - 1/2 + 1}{2} \right] = p$$

$$\Rightarrow \left[ \frac{E(\tilde{p}) - 1/2 + 1/2}{2v} \right] = p$$

$$\therefore \tilde{p}_0 = \left[ \frac{\tilde{p} - 1/2 + 1/2}{2v} \right]$$

$\tilde{p}_0$  - is the unbiased term

$$E(\tilde{p}_0) = E \left[ \frac{\tilde{p} - 1/2 + 1/2}{2v} \right]$$

$$= \frac{E[\tilde{p}] - 1/2 + 1/2}{2v} = \frac{1/2 + (2p - 1)v - 1/2 + 1/2}{2v}$$

$$= \frac{(2p - 1)v}{2v} + \frac{1}{2} = \frac{2p - 1 + 1}{2} = \frac{2p}{2} = p //$$

$\therefore \tilde{p}_0$  is an unbiased estimator of 'p'

$$c) \text{Var}[\tilde{p}_0] \leq \frac{1}{16v^2} N$$

$$\text{As in above solution } \tilde{p}_0 = \frac{\tilde{p} - 1/2 + 1/2}{2v}$$

$$\text{Var}(\tilde{p}_0) = \text{Var} \left[ \frac{\tilde{p} - 1/2 + 1/2}{2v} \right] = \text{Var} \left[ \frac{\tilde{p} - 1/2}{2v} \right]$$



We know that  $\text{Var}(ax) = a^2 \text{Var}(x)$

$$\therefore \text{Var}\left(\frac{\tilde{p} - 1/2}{2v}\right) = \frac{1}{4v^2} \text{Var}(\tilde{p} - 1/2) \\ = \frac{1}{4v^2} \text{Var}(\tilde{p})$$

We know  $\tilde{p} = \frac{1}{N} \sum_{i=1}^N y_i$

Variance of  $\tilde{p}$   $\text{Var}(\tilde{p}) = \frac{1}{N^2} \sum_{i=1}^N \text{Var}(y_i)$

$$E[y_i] = 1/2 + v(2x_i - 1)$$

$$\text{Var}(y_i) = E(y_i)(1 - E(y_i)) \\ = \left(\frac{1}{2} + v(2x_i - 1)\right) \left(\frac{1}{2} - v(2x_i - 1)\right)$$

$$\text{Var}(y_i) = \frac{1}{4} - v^2(2x_i - 1)^2 \\ = \frac{1}{4} - v^2 \quad (\because (2x_i - 1)^2 = 1)$$

$$\text{Var}(\tilde{p}) = \frac{1}{N} \left(\frac{1}{4} - v^2\right)$$

$$\text{Var}(\tilde{p}_0) = \frac{1}{4v^2} \text{Var}(\tilde{p}) \\ = \frac{1}{4v^2} \left(\frac{1}{N}\right) \left(\frac{1}{4} - v^2\right) \\ = \frac{1}{4v^2 N} \left(\frac{1}{4} - v^2\right) = \frac{1}{4v^2 N \cdot 4} = \frac{1}{16v^2 N}$$

$$\therefore \text{Var}[\tilde{p}_0] \leq \frac{1}{16v^2 N}$$

d)  $\tilde{p}$  is computed directly from the randomized response  $y_i$  which already satisfy DP.  $\tilde{p}_0$  is just a mathematical adjustment of  $\tilde{p}$ . It doesn't use any additional info/randomness. As DP is immune to post processing means, Once data  $y_i$  is released with a DP guarantee, any further computation (like calculating  $\tilde{p}_0$ ) doesn't weaken/strengthen the privacy guarantee. Hence  $\tilde{p}_0$  and  $\tilde{p}$  have exact same DP.

c) We have  $\tilde{P}_0 = \frac{\tilde{P} - 1/2}{2\sqrt{N}} + 1/2$

$$E[\tilde{P}_0] = P$$

$$\text{Var}[\tilde{P}_0] \leq 1/16 v^2 N$$

Apply Cheby Shev's Inequality:-

$$X = \tilde{P}_0; \ell = P; \sigma^2 = \text{Var}[\tilde{P}_0]$$

$$\Pr[|\tilde{P}_0 - P| \geq k \sigma] \leq 1/k^2$$

$$\left[ \because \sigma = \sqrt{\text{Var}[\tilde{P}_0]} \leq \frac{1}{4\sqrt{N}} \right]$$

$$\Pr\left[|\tilde{P}_0 - P| \geq k \frac{1}{4\sqrt{N}}\right] \leq 1/k^2$$

Let  $k = 1/s$

$$\Pr\left[|\tilde{P}_0 - P| \geq 1/s \left(\frac{1}{4\sqrt{N}}\right)\right] \leq s^2$$

$$|\tilde{P}_0 - P| \leq \frac{1}{4\sqrt{N}s} \Rightarrow |\tilde{P}_0 - P| = o\left(\frac{1}{\sqrt{N}}\right)$$

$$\boxed{\because |\tilde{P}_0 - P| = o\left(\frac{1}{\sqrt{N}}\right)}$$