# KFUPM
# College of Computer Science and Engineering
# Computer Engineering Department
# COE 449: Privacy Enhancing Technologies

Fall 2019 (191)

Assignment 2: Due date Thursday 26/10/

Student: **Faris Hijazi s201578750**

## Tasks

### Question1: Laplace Mechanism (20 pts)

Prove that a Laplace mechanism that adds random noise us- ing Laplace distribution with $\frac{GS_F}{\epsilon}$ is $\epsilon$-deferentially private.

# 449 HW2 DP

## Laplass

$$pdf = \frac{1}{2\lambda} e^{-\frac{|x-\mu|}{\lambda}}$$  varience: $2\lambda^2$

$$\lambda = \frac{GS_F}{\epsilon}$$

## Proving laplace is DP

$f$ is the query result

$$\frac{P[A(D_1) = 0]}{P[A(D_2) = 0]} \leq e^{\epsilon}$$

- Probability density $(0)$ & Prob density of added noise; $e^{-\frac{|F-0|}{\lambda}}$

- $\frac{P[F(D_1) + Lap(\lambda) = 0]}{P[F(D_2) + Lap(\lambda) = 0]} \leq \frac{e^{-\frac{|F(D_1)-0|}{\lambda}}}{e^{-\frac{|F(D_2)-0|}{\lambda}}} \leq e^{-\frac{|F(D_1)-F(D_2)|}{\lambda}} \xleftarrow{GS_F}$

$$\leq e^{-\frac{GSF}{\lambda}} \leq e^{\epsilon} \quad \left( \begin{array}{c} because \\ \lambda = GS_F/\epsilon \end{array} \right)$$

## Question2: Composition (30 pts)

- (a) Assume that you have $k$ randomized algorithms $A_1, \ldots, A_k$, each is $\epsilon$ differentially private. Prove that if you execute $A_1, \ldots, A_k$, sequentially then the composed mechanisms is $k\epsilon$-differentially private.
- (b) Assume that you have a database partitioned in $k$ disjoint subsets $D_i$. For each subset, we run a randomized mechanims $A_i$ that is $\epsilon$-differentially private. Prove that executing the $k$ mechanisms in parallel over $D$ is $\epsilon$-differentially private.

## Q2 Composition

a) sequential $A_1 \ldots A_k$ is $k\epsilon$ DP

We must assume that $A_i$ is
and input

suppose $\boxed{A(D) = Z_k}$

$z_i = A_i(D; z_1, z_2 \ldots z_{i-1})$

Modifying the $A_i$ function, $A_i$ must take an input of the previous $A_{i-1}$, for propogating the expected value $\underline{z_i}$

$$Pr\left[A(D_1) = Z_k\right] = Pr\overset{\text{chaining}}{\left[A_1(D_1) = Z_1\right]} Pr\left[A_2(D_1, \check{Z}_1) = Z_2\right] \ldots$$

$$\leq \exp(k\epsilon) \prod_{i=1}^{K} Pr\left[A_i(D_2, z_1, z_2 \ldots) = z_i\right]$$

$$= \exp(k\epsilon) \, Pr\left[A(D_2) = Z_k\right]$$

b) __Parallel__ prove $\epsilon$ - DP

$j$ partitions

suppose $A(D) = Z_k$

$$Pr\left[A(D) = Z_x\right] = \prod_{j=1}^{K} Pr\left[A_i(D_{ij}, z_1, z_2 \ldots z_{i-1}) = z_i\right]$$

$$\leq \exp(\epsilon) \, Pr\left[A_j(D_2; z_1 \ldots z_j) = z_j\right] \prod_{\substack{i=1 \\ i \neq j}}^{K} Pr\left[A_i(D, z \ldots = z_j\right]$$

$$= \exp(\epsilon) \, Pr\left[A(D_2) = z_x\right]$$

## Question3: Differential Privacy Implementation (50 pts)

The goal of this task is to understand the concept of differential privacy by implementing a differentially private histogram using Laplace mechanism. The objective is to generate and draw noisy histogram bins.

### Input

Your program takes the following as inputs

- Input dataset files
- Privacy budget $\epsilon$
- Number and sizes of bins

### Dataset description

The dataset used in this task is the IPUMS data extracted from the 2001 US Census. The dataset has 8 attributes as described in Table 1. The size of the dataset is 20,000 tuples (rows). All attributes include numerical values only. For example, Gender attributes can be either 1 or 2, which represents Male and Female, respectively. The Income attribute is the annual income in thousand USD, for example, an income of 20 means 20,000 (20K).

| . | . | . | . | . | . | . | . |
|---|---|---|---|---|---|---|---|
| Age | Gender | Marital | Race status | Birth place | Language | Occupation | Income (K) |

Table 1: Scheme of Census dataset

### Output

The output of your program is a differentially private his- togram for each of the following two dimensions in IPUMS dataset

- Age
- Salary

In [1]:

```python
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

N_BINS = 10 # number of bins
epsilon = 0.1 # privacy budget E
```

```
data = pd.read_csv('ipums.csv')
data.head()
```

Out[2]:

| | Age | Gender | Marital | Race status | Birth place | Language | Occupation | Income (K) |
|---|---|---|---|---|---|---|---|---|
| 0 | 33 | 1 | 6 | 2 | 1 | 1 | 10 | 144 |
| 1 | 40 | 2 | 4 | 1 | 1 | 1 | 6 | 830 |
| 2 | 21 | 2 | 6 | 1 | 1 | 1 | 3 | 992 |
| 3 | 39 | 1 | 4 | 1 | 1 | 1 | 6 | 673 |
| 4 | 55 | 2 | 4 | 1 | 1 | 1 | 10 | 470 |

In [3]:

```
age = data['Age']
salary = data['Income (K)']

age_hist, age_bins = np.histogram(age, bins=N_BINS)
salary_hist, salary_bins = np.histogram(salary, bins=N_BINS)
```

```python
def plot_hist_from_bins(hist, bins, **kwargs):
    fig, ax = plt.subplots()
    # hist, bins = np.histogram(x, **kwargs)
    width = 0.9 * (bins[1] - bins[0])
    center = (bins[:-1] + bins[1:]) / 2

    ax.bar(center, hist, align='center', width=width)
    return fig

print('Unperturbed histograms (RAW)')

f = plot_hist_from_bins(age_hist, age_bins)
plt.title('age (RAW)')

f = plot_hist_from_bins(salary_hist, salary_bins)
plt.title('salary (RAW)')
```
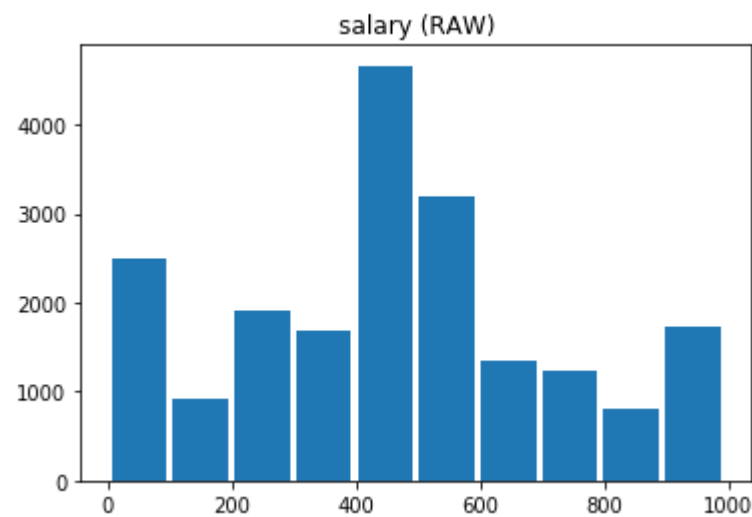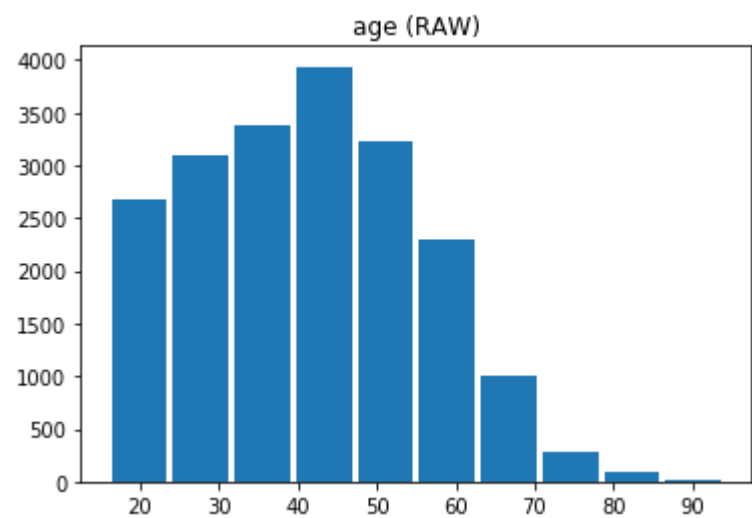
Unperturbed histograms (RAW)

Text(0.5, 1.0, 'salary (RAW)')

age (RAW)



salary (RAW)

```
# adding the Laplace noise to each bin
lambda_ = 1/epsilon
print('lambda = ', lambda_)

age_hist_laplace += np.random.laplace(scale=lambda_, size=age_hist.shape)
salary_hist_laplace += np.random.laplace(scale=lambda_, size=salary_hist.shape)

f = plot_hist_from_bins(age_hist_laplace, age_bins)
plt.title('age (Laplace)')

f = plot_hist_from_bins(salary_hist_laplace, salary_bins)
plt.title('salary (Laplace)')
```
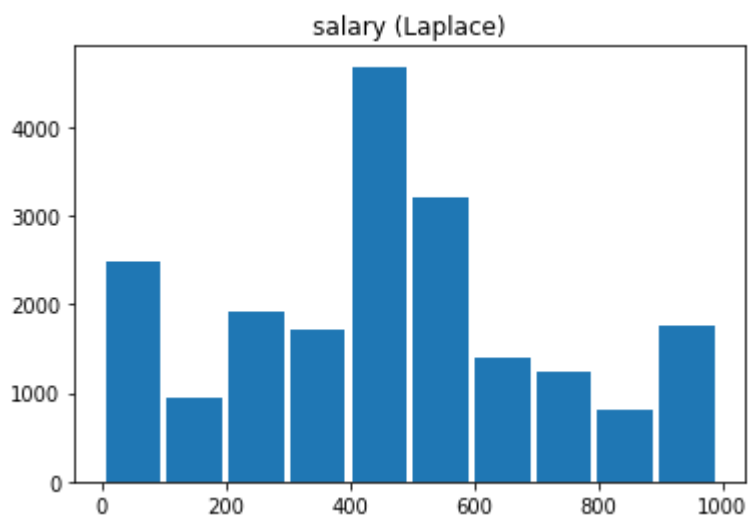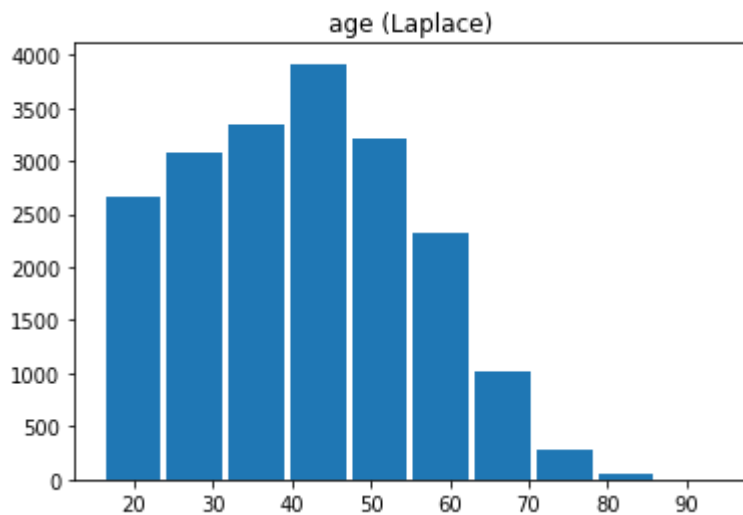
lambda =  10.0

Out[17]:

Text(0.5, 1.0, 'salary (Laplace)')

```python
# buckets/bins deltas (new vs old)
print('Age bin deltas:', age_hist_laplace - age_hist)
print()
print('Salary bin deltas:', salary_hist_laplace - salary_hist)
```

```
Age bin deltas: [ -5.19483814 -21.90559102 -19.27379443 -21.51569508 -11.0
6537617
   7.99750608  15.08609657  -2.61397685 -36.36093136 -13.91743756]

Salary bin deltas: [-0.11418419 31.42963447   5.40823898 17.81573584 21.231
6597   4.37342763
 38.67691332 13.10542777 13.36646206 36.2452841 ]
```

## Submission

The due date of this assignment is 11:59PM 26/10/2019. Please upload all files on the assignment page on BlackBoard. You need to submit the following:

1. A brief report including your answers to Questions 1 and 2. In addition, the report should include the obtained differentially private histograms from Question 3, and a comparison with the original hostograms (without privacy). Explain and discuss, if any, the algorithmic optimizations you have used in your implementation. Discuss the experiences and lessons you have learned from the implementation.
2. A zip file that contains the PDF brief report, the source files of your code (in any language), and a README explaining how to compile/run your program

## References

1. https://people.eecs.berkeley.edu/~stephentu/writeups/6885-lec20-b.pdf (https://people.eecs.berkeley.edu/~stephentu/writeups/6885-lec20-b.pdf)

In [ ]: