/burl@stx null def /BU.S /burl@stx null def def /BU.SS currentpoint /burl@lly exch def /burl@llx exch def burl@stx null ne burl@endx burl@llx ne BU.FL BU.S if if burl@stx null eq burl@llx dup /burl@stx exch def /burl@endx exch def burl@lly dup /burl@boty exch def /burl@topy exch def if burl@lly burl@boty gt /burl@boty burl@lly def if def /BU.SE currentpoint /burl@ury exch def dup /burl@urx exch def /burl@endx exch def burl@ury burl@topy lt /burl@topy burl@ury def if def /BU.E BU.FL def /BU.FL burl@stx null ne BU.DF if def /BU.DF BU.BB [ /H /I /Border [burl@border] /Color [burl@bordercolor] /Action « /Subtype /URI /URI BU.L » /Subtype /Link BU.B /ANN pdfmark /burl@stx null def def /BU.BB burl@stx HyperBorder sub /burl@stx exch def burl@endx HyperBorder add /burl@endx exch def burl@boty HyperBorder add /burl@boty exch def burl@topy HyperBorder sub /burl@topy exch def def /BU.B /Rect[burl@stx burl@boty burl@endx burl@topy] def /eop where begin /@ldeopburl /eop load def /eop SDict begin BU.FL end @ldeopburl def end /eop SDict begin BU.FL end def ifelse

# CIS 6930: Privacy & Machine Learning (Fall 2019)
# Homework 1 — Data Privacy

Name: Your Name Here

September 27, 2019

**This is an individual assignment!**

## Instructions

Please read the instructions and questions carefully. Write your answers directly in the space provided. Compile the tex document and hand in the resulting PDF.

For this homework, you will solve several data privacy problems. The fourth problem asks you to implement a differential privacy mechanism using Python. Use the code skeleton provided and submit the completed source file(s) alongside with the PDF.[1] *Note: bonus points you get on this homework \*do\* carry across assignments/homework.*

# Problem 1: Syntactic Metrics (20 pts)

Consider the data set depicted in **??**. Answer the following questions. (Justify your answers as appropriate.)

| Age | Zip Code | Sex | Credit Score | Yearly Income | Loan |
|-----|----------|-----|--------------|---------------|------|
| 30-39 | 32607 | M | 678 | 90k | Approved |
| 30-39 | 32607 | M | 799 | 120k | Approved |
| 40-49 | 32611 | F | 451 | 35k | Declined |
| 20-29 | 32607 | F | 783 | 30k | Approved |
| 20-29 | 32607 | F | 560 | 70k | Declined |
| 40-49 | 32611 | M | 725 | 22k | Declined |

Table 1: Anonymized Data Set.

1. (4 pts) What are the quasi-identifiers? What are the sensitive attributes?

   *Your answer here.*

2. (4 pts) What is the largest integer $k$ such that the data set satisfies $k$-anonymity? What is the largest integer $l$ such that the data set satisfies $l$-diversity?

   *Your answer here.*

3. (6 pts) Modify the data set using generalization and suppression to ensure that it satisfies 3-anonymity and 2-diversity. Here we are looking for a solution that minimally affects the utility of the data. Write the modified data set below.

   *Your answer here.*

4. (6 pts) Your student friend Alice (who is not in the anonymized data set) was recently declined for a loan despite her 30k yearly income. She thinks she may have been discriminated against.

   The bank who declined Alice's loan has published the following transparency report about their loan approval model.

   - If `yearly_income` $\geq$ 50k then: return `APPROVED`
   - If `yearly_income` $\geq$ 25k:
     - If `student`:
       * If `credit_score` $\geq$ 550 then: return `APPROVED`
       * Else: return `DECLINED`
     - Else (not `student`) if `credit_score` $\geq$ 500 then: return `APPROVED`
   - If `yearly_income` $\geq$ 20k and `credit_score` $\geq$ 650 then: return `APPROVED`
   - return `DECLINED`

   What can you infer about Alice assuming that the transparency report accurately reflects the loan approval model? What do you conclude about the possible tension between algorithmic transparency and privacy? (Explain your answer.)

   *Your answer here.*

# Problem 2: Randomized Response & Local Differential Privacy (25 pts)

Social science researchers at the University of Florida want to conduct a study to explore the prevalence of crime among students. Specifically they want to ask questions of the form: *have you ever committed crime X?* (Here X stands for a specific crime or crime category.)

Researchers are ethical so they want to carefully design the study to ensure that participants respond truthfully and that privacy is protected. They reached out to you, a data privacy expert, to evaluate their methods.

Consider a participant that is asked the question have you ever committed crime X? This question admits a yes or no answer. Before answering the participant is instructed to use the following algorithm to compute a "noisy" answer given their true answer and only report the noisy answer to the researchers.

---

NoisyAnswer(`true_answer`, $p \in (0, 1)$):

- If `true_answer` is `NO`, then:
    - With probability $p$ return `NO`
    - With probability $1 - p$ return `YES`
- Else (if `true_answer` is `YES`), then: return `YES`

---

Answer the following questions.

1. (10 pts) Suppose the researchers obtain noisy answers $z_1, z_2, \ldots, z_n$ from the $n$ study participants. You can assume that `YES` is encoded as 1 and `NO` is encoded as 0. Explain how the researchers can estimate the true proportion of `YES` from the noisy answers. Specifically, give formulae for (1) the expected number of `YES` answers and (2) the variance (or error) of the estimate.

    *Your answer here.*

2. (5 pts) Consider the following definition of (Local) Differential Privacy.

    **Definition 1.** *A randomized algorithm $\mathcal{F}$ which takes input in some set $X$ satisfies $\varepsilon$-differential privacy (for some $\varepsilon > 0$) if for any two input records $x \in X$, $x' \in X$ and any output $z \in \text{Range}(\mathcal{F})$:*

    $$\Pr\{z = \mathcal{F}(x)\} \leq e^{\varepsilon} \Pr\{z = \mathcal{F}(x')\} \ .$$

    Does the noisy answer algorithm satisfy **??**? Produce a proof or a counter-example. If it does, also give an expression for $\varepsilon$ in terms of $p$.

    *Your answer here.*

3. (5 pts) Now consider the following (more general) variant of the algorithm.

---

GeneralizedNoisyAnswer(`true_answer`, $p, p' \in (0, 1)$):

- If `true_answer` is `NO`, then:
    - With probability $p$ return `NO`
    - With probability $1 - p$ return `YES`
- Else (if `true_answer` is `YES`), then:
    - With probability $p'$ return `NO`
    - With probability $1 - p'$ return `YES`

---

Prove that this general variant satisfies $\varepsilon$-(local) differential privacy (**??**). Give an expression for $\varepsilon$ in terms of $p$ and $p'$.

*Your answer here.*

4. (5 pts) Suppose we can arbitrarily set $p$ and $p'$. Explain the trade-off between minimizing $\varepsilon$ and minimizing the error between the true answers and the one estimated from noisy answers.

   *Your answer here.*

# Problem 3: Privacy & Sampling (30 pts)

Consider the data set shown in **??** and the function $f(\mathbf{x}) = \text{mode}(\mathbf{x})$. Here the mode of a dataset is 0 if the proportion of individuals who are HIV negative (-) is higher than 0.5, 1 otherwise.

| | HIV | |
|---|---|---|
| | + | - |
| # of individuals | 7 | 23 |

Table 2: Data set.

We are interested in various ways of designing a differentially private mechanism to compute $f$ on an arbitrary data set of the same form as the one in **??**.

1. (2 pts) What is the *local* sensitivity of $f$ (with respect to the data set shown in **??**)?

   *Your answer here.*

2. (2 pts) What is the global sensitivity of $f$? (Justify your answer.)

   *Your answer here.*

3. (3 pts) Consider the mechanism defined by $\mathcal{F}(\mathbf{x}) = f(\mathbf{x})$. Does this mechanism satisfy $\varepsilon$-differential privacy? Why or why not?

   *Your answer here.*

Now consider your answers to the previous questions. We are interested in using Laplace noise to obtain a $\varepsilon$-differentially private mechanism for $f$.

4. (3 pts) Explain how you could add Laplace noise to obtain $\varepsilon$-differential privacy for $f$. Call the resulting mechanism $\mathcal{F}$.

   *Your answer here.*

5. (5 pts) Now consider the following post-processing step (after adding Laplace noise as you explained): return 0 if $\mathcal{F}(\mathbf{x}) < 0.5$ and 1 otherwise. If the data set $\mathbf{x}$ is such that $f(\mathbf{x}) = 1$, what is the probability that (after the post-processing step) the output is 0? What do you conclude about this mechanism?

   *Your answer here.*

6. (5 pts) Can you come up with a different $\varepsilon$-differential privacy mechanism for $f$ that adds Laplace noise but provides more accurate outputs?

   *Your answer here.*

7. (10 pts) Finally, consider the following mechanism.

   > **SampleAndComputeMode(data set $\mathbf{x}$, $p \in (0, 1)$):**
   > - Let $\mathcal{M}$ be a $\varepsilon$-differentially private mechanism to compute the mode of a data set.
   > - Let $\mathbf{s}$ be the data set obtained by independently selecting each record of $\mathbf{x}$ with probability $p$. (For each record, we flip a coin with probability of heads $p$, if heads then we add this record to $\mathbf{s}$, otherwise we do not.)
   > - return $\mathcal{M}(\mathbf{s})$.

Prove that SampleAndComputeMode() satisfies $\varepsilon'$-differential privacy and give an expression for $\varepsilon'$ in terms of $p$ and $\varepsilon$.

*Your answer here.*

# Problem 4: Implementing DP Mechanisms (25 pts)

For this problem you will implement several differential privacy mechanisms we talked about in class. Please use the comments in the Python files provided to guide you in the implementation.

For this question, we will use the dataset `data/ds.csv`. It contains pairs of age and yearly income for several individuals. For the purpose of calculating sensitivity, assume that the age range for any individual is [16, 100] and the yearly income range is [0, 1000000].

1. (5 pts) Fill in the implementation of `laplace_mech()`, `gaussian_mech()`. Also fill in the (global) sensitivity in the `mean_age_query()` function.

   You can test your implementation by running: 'python3 hw1.py problem4.1'.

   How close are the noisy answers to the true answer?

   *Your answer here.*

2. (5 pts) Complete the implementation of the `dp_accuracy_plot()` and run it for $\varepsilon = 0.1, 0.5, 1.0, 5.0$ on `mean_age_query()`. Paste the plots below.

   To run the code: 'python3 hw1.py problem4.2 <epsilon>'. By default, figures are saved in `./plots` and named based on the value of $\varepsilon$.

   What do you conclude?

   *Your answer here.*

3. (5 pts) Implement the function called `budget_plot()`. Use it to produce a plot of the budget of naive composition and advanced composition (refer to the course materials for details) when using `gaussian_mech()` to perform `mean_age_query()` $m > 1$ times. Plot the naive composition and advanced composition budgets (i.e., total $\varepsilon$) for varying $m$ from 1 to 100 keeping $\delta \leq 2^{-30}$. Paste the plot below. For what values of $n$ is naive composition better than advanced composition? (Justify your answer.)

   *Your answer here.*

4. (10 pts) Finally, suppose we want to compute the average ratio of yearly income and age in the dataset, i.e., how many extra dollars does one earn for an increase of one year of age (on average). Consider two ways of performing this query with differential privacy:

   (a) Compute the (global) sensitivity of this query (`income_per_age_query()`) and use the Laplace mechanism.

   (b) Use the Laplace mechanism to compute the mean yearly income. Use the Laplace mechanism to compute the mean age. Divide the two (noisy) results to obtain the ratio.

   Implement this functionality in `income_per_age_comp()`. Feel free to modify the signature of `income_per_age_comp()` and the corresponding code in `main()`. Set $\varepsilon = 1$. Paste the comparison plot below.

   What do you conclude?

   *Your answer here.*

# [Bonus] Problem 5: Privacy with Binomial Noise (20 pts)

Suppose we are interested in non-negative count functions $f$. For example $f$ is the number of records in the dataset which satisfy some property $P$. Consider the mechanism $\mathcal{F}(\mathbf{x}) = f(\mathbf{x}) + B$, where $B \sim \text{Binom}(n, p) - \mathbb{E}[\text{Binom}(n, p)]$. Here $n = |\mathbf{x}| > 0$ is the size of the dataset and $p \in (0, 1)$ is a parameter. In other words $\mathcal{F}$ adds noise from the binomial distribution but centered at 0.

1. (2 pts) How would you set the value of $p$? (Explain your answer.)

    *Your answer here.*

2. (3 pts) Under what condition(s) does $\mathcal{F}$ satisfy $\varepsilon$-differential privacy? (Justify your answer.)

    *Your answer here.*

3. (10 pts) Prove that $\mathcal{F}$ satisfies $(\varepsilon, \delta)$-differential privacy as long as $p \neq 0, 1$.

    *Your answer here.*

4. (5 pts) Characterize the trade-off between $\varepsilon$ and $\delta$.

    *Your answer here.*