

Assignment 6: Privacy
15-316 Software Foundations of Security and Privacy

Total Points: 30

1. **An attack (15 points).**

Suppose that you are tasked with releasing statistics about a dataset with the following schema.

Feature	<i>age</i>	<i>gender</i>	<i>marital_status</i>	<i>college_education</i>	<i>salary</i>
Encoding	$\text{int} \in [0, 100]$	$\text{int} \in \{0, 1\}$	$\text{int} \in \{0, 1\}$	$\text{int} \in \{0, 1\}$	$\text{int} \in [0, 10^6]$

Show that if you are allowed to query the dataset by counting the number of entries with a particular set of values, then it is possible to learn a person's salary. In particular, you have access to the following function, which you can query as many times as you like:

$$\text{count}(\textit{age}, \textit{gender}, \textit{education}, \textit{marital}, \textit{salary}) = |\{\# \text{ database rows matching given values}\}|$$

Your aim is to learn the salary of a particular individual for whom you know all attributes *except* salary, and you may assume knowledge of the rest of the dataset as described in lecture for the differential privacy threat model.

Your solution should describe (pseudocode is fine) a general procedure for learning the salary from the given information, regardless of the particular contents of the database.

2. **A fix (15 points).** Now your goal is to provide statistics about the average salary across gender and education level while satisfying ϵ -differential privacy.

- You have access to the database through variable X , which you should assume is an array containing N dictionaries that you can index by attribute name; i.e., $X[0][\text{"salary"}]$ returns the salary of the first row of the database.
- You may call a function `Laplace(b)`, which returns a single random sample from the zero-centered Laplace distribution with scale parameter b ; and a function `Uniform(a, b)` which returns a uniform random real number between a and b , inclusive. Note that your solution need not necessarily call both of these functions.
- You may assume that the breakdown of X by gender and education level is not private information.

Explain how to implement a 1-differentially private function `mean_by_gender_and_edu`, which returns a 4-tuple of floats containing the mean salary for each gender and education level in X . That is, this function privately computes the following statistics:

$$\text{mean_by_gender_and_edu} = (\text{mean}(\text{women}), \text{mean}(\text{men}), \text{mean}(\text{college}), \text{mean}(\text{nocollege}))$$

Be sure to state which composition principles your solution uses, if any. If it is easiest to present your solution as pseudocode then please do so, but you should explain how it works in words as well.