

Assignment 4
CS260B: Algorithmic Machine Learning, Spring 2021
Due: May 25, 10PM

Guidelines for submitting the solutions:

- The assignments need to be submitted on Gradescope. Make sure you follow all the instructions - they are simple enough that exceptions will not be accepted.
 - Start each problem or sub-problem on a separate page even if it means having a lot of white-space and write/type in large font.
 - The solutions need to be submitted by 10 PM on the due date. No late submissions will be accepted.
 - Please adhere to the code of conduct outlined on the class page.
1. Consider the following alternative to potentially achieve ϵ -differential privacy for releasing the mean of a database $X \in \{0, 1\}^n$. So we have n users each with one attribute and we want to release the count of users with this attribute so $f(X) = \sum_i X_i$. As studied in class, this query has $S_1(f) = 1$ so we can use Laplacian noise. But what if instead of Laplacian noise we consider the following strategy for achieving privacy: 1. Sample a uniform real-value in the interval $[-C/\epsilon, C/\epsilon]$ (for some constant C). 2. Release $f(X) + Z$.

Is the above mechanism ϵ -differentially private? (For some fixed constant C .) [4 points]

2. In class we stated without proof that the exponential mechanism gets a reasonable guarantee for preserving the utility of the released id. In particular, we stated in class without proof that (using notation from class), if we use exponential mechanism for a utility function with sensitivity Δu ,

$$\Pr[u(M_E(X, u, R)) \leq \text{OPT}_u(X) - \frac{\Delta u}{\epsilon} \cdot (\ln |R| + t)] \leq e^{-t}.$$

Prove the above statement. [4 points]

[Hint: Try to reason that every item with as small utility as above is quite unlikely, and then use that the total number of items is at most $|R|$.]

3. Suppose your database is the income numbers of n individuals. What scheme would you use to release the median income in the dataset to satisfy ϵ -differential privacy while achieving a good guarantee on error?

Concretely, suppose the each persons income is an integer in $\{0, 1, 2, \dots, N\}$. So the database $X \in [N]^n$ and the query we are trying to release privately is $Median(X)$. While this is not needed, if it helps for you, you can suppose that all incomes are distinct and that n is odd.

- What is the sensitivity of this query as a function of N ? [2 points]
- What would happen if you use Laplacian mechanism given this sensitivity bound? [2 points]
- Describe a *score function* or *utility function* for which the arg max would exactly be the median and one whose sensitivity does not depend on N, n . Note the remarkable gains you would get by now using exponential mechanism with this score function to release the median instead! [2 points]
- Describe a *score function* for which the arg max would exactly be the 90'th percentile income level among the people in the database and one whose sensitivity is still independent of N, n . [2 points]