

# CSC2412 F19: Assignment 1

Due: October 25, by 11:59pm

## Guidelines:

- Your assignment solution must be submitted as a *typed* PDF document. Scanned handwritten solutions, solutions in any other format, or unreadable solutions will **not** be accepted or marked. You are encouraged to learn the  $\text{\LaTeX}$  typesetting system and use it to type your solution.
- To submit this assignment, use the MarkUs system, at URL <https://markus.teach.cs.toronto.edu/csc2412-2019-09>
- This is an *individual assignment*. You may consult any of the reading material posted on the course website. However, your solutions should show your individual work.
- You may use any result discussed in class or covered in the assigned reading by just referring to it. You do not need to reproduce proofs that we have covered in the lectures.
- Unless stated otherwise, you should justify all your answers using rigorous arguments. Your solution will be marked based both on its completeness and correctness, and also on the clarity and precision of your explanation.

For all questions you can assume that the notion of neighbouring datasets used is the one based on replacement, and the size of the dataset is public. I.e. datasets  $X, X' \in \mathcal{X}^n$  are neighbouring if there exists an  $i \in \{1, \dots, n\}$  such that  $x'_j = x_j$  whenever  $j \neq i$ .

**Question 1.** (5 marks)

Suppose you are given a private dataset  $X \in \mathcal{X}^n$ , where  $\mathcal{X} = \{1, \dots, N\}$ . I.e. the dataset  $X$  consists of  $n$  integers  $x_1, \dots, x_n$  between 1 and  $N$ . Describe an  $\varepsilon$ -differentially private algorithm, based on the exponential mechanism, which outputs a number  $y \in \mathcal{X}$  such that if  $n \geq \frac{C_1}{\varepsilon} \ln(|N|/\beta) + C_2$ , then

$$\mathbb{P} \left( \min_{i=1}^n x_i \leq y \leq \max_{i=1}^n x_i \right) \geq 1 - \beta.$$

Above  $C_1$  and  $C_2$  are constants independent of  $n, N, \beta$ , and  $\varepsilon$ . Justify why your algorithm is  $\varepsilon$ -differentially private, and why it satisfies the property above. Specify the constants  $C_1, C_2$  in your answer.

**Question 2.** (10 marks)

For this question you can use the following identity: for a Laplace random variable  $w \sim \text{Lap}(b)$ , we have for any  $t \geq 0$

$$\mathbb{P}(|w| \geq t) = e^{-t/b}.$$

The goal in this question is to design an algorithm which estimates the mean of a dataset of numbers. The estimate should be accurate whenever the numbers are bounded, but the algorithm should be private even if the numbers are arbitrary.

**Part a.** (4 marks)

Suppose that the dataset  $X = (x_1, \dots, x_n)$  consists of integers, which can be positive or negative, and are not *a priori* bounded. Describe an  $\varepsilon$ -differentially private algorithm  $\mathcal{A}$  such that, if  $x_i \in [-B, +B]$  for every  $i$ , and  $n \geq \frac{2B \ln(1/\beta)}{\varepsilon \alpha}$ , then

$$\mathbb{P} \left( \left| \mathcal{A}(X) - \frac{1}{n} \sum_{i=1}^n x_i \right| \geq \alpha \right) \leq \beta. \quad (1)$$

You can assume that the parameter  $B$  is known to the algorithm. Note that while (1) needs to hold only if  $x_i \in [-B, +B]$  for all  $i$ , the algorithm  $\mathcal{A}$  needs to be  $\varepsilon$ -differentially private for **every** dataset, i.e. even if  $x_i \notin [-B, B]$  for some values of  $i$ . Justify your answer.

**Part b.** (6 marks)

Suppose that the dataset  $X = (x_1, \dots, x_n)$  consists of integers in  $[-N, +N]$ , where  $N$  is some large integer. Describe an  $\varepsilon$ -differentially private algorithm  $\mathcal{A}$  such that, if there exists some integer  $z \in [-N, +N]$  for which  $x_i \in [z - B, z + B]$  for every  $i$ , and

$$n \geq \max \left\{ \frac{C_1 B \ln(2/\beta)}{\varepsilon \alpha}, \frac{C_2 \ln(2|2N + 1|/\beta)}{\varepsilon} + C_3 \right\}$$

then

$$\mathbb{P} \left( \left| \mathcal{A}(X) - \frac{1}{n} \sum_{i=1}^n x_i \right| \geq \alpha \right) \leq \beta. \quad (2)$$

Above  $C_1, C_2$ , and  $C_3$  are constants independent of  $n, N, \beta, \varepsilon, B$ , and  $z$ . You can assume that the parameter  $B$  is known to the algorithm, but the parameter  $z$  is **not** known. Once again, the algorithm  $\mathcal{A}$  needs to be  $\varepsilon$ -differentially private for **every**  $X \in \{-N, \dots, +N\}^n$ . Justify your answer, and specify the constants  $C_1, C_2$ , and  $C_3$ .

HINT: You can use Question 1 to help you solve this subquestion.