

# HW 3: The Centralized Curator Model

CS 208 Applied Privacy for Data Science, Spring 2019

**Version 1.0: Due Tuesday, April 2, 11:59pm.**

**Instructions:** Submit a single PDF file containing your solutions, plots, and analyses. Make sure to thoroughly explain your process and results for each problem. Also include your documented code and a link to a public repository with your code (such as GitHub/GitLab). Make sure to list all collaborators and references.

1. **Tails, Trimming, and Winsorization:** In all of the parts below, the dataset is  $x \in \{0, 1, \dots, D\}^n$ . In all of the implementation parts, you should write code that takes as input  $D \in \mathbb{N}$ ,  $n \in \mathbb{N}$ ,  $x \in \{0, 1, \dots, D\}^n$ , and  $\varepsilon > 0$ .

- (a) Prove that the following algorithm for estimating a Trimmed mean is  $\varepsilon$ -DP and implement it in code:

$$M(x) = \frac{1}{.9n} \cdot \left( \sum_{P_{.05} \leq x_i \leq P_{.95}} x_i \right) + \text{Lap} \left( \frac{D}{\varepsilon n} \right),$$

where  $P_{.05}$  and  $P_{.95}$  are the 5th and 95th percentile of the dataset. That is, we are applying the Laplace mechanism after removing the bottom and top 5% of the dataset. (Hint: Think about Lipschitz constants.)

- (b) Prove that for large enough  $n$ , the analogous algorithm for the *Winsorized* mean is *not*  $\varepsilon$ -DP:

$$M(x) = \frac{1}{n} \cdot \sum_{i=1}^n [x_i]_{P_{.05}}^{P_{.95}} + \text{Lap} \left( \frac{D}{\varepsilon n} \right),$$

where  $[x]_a^b$  is defined as in Problem Set 2. In Winsorization, we clamp points rather than dropping them. (In class on 3/11, we incorrectly referred to dropping points as Winsorization.) Again, it may be useful to first think in terms of Lipschitz constants.

- (c) In class, we saw how to use the exponential mechanism to an estimate of the median,  $P_5$ . Describe and implement a version of the exponential mechanism that releases an estimate of the  $t$ th percentile  $P_t$  of a dataset  $x \in \{0, \dots, D\}^n$  any desired  $t \in [0, 100]$ .
- (d) Implement the following  $\varepsilon$ -DP algorithm for estimating a Trimmed mean of a dataset: use your algorithm from Part 1c to get  $\varepsilon/3$ -DP estimates  $\hat{P}_{.05}$  and  $\hat{P}_{.95}$  of the 5th and 95th percentiles, drop all datapoints that lie outside the range  $[\hat{P}_{.05}, \hat{P}_{.95}]$ , and then use the Laplace mechanism to compute an  $(\varepsilon/3)$ -DP mean of the trimmed data. That is, your code should compute and

$$M(x) = \frac{1}{.9n} \cdot \left( \sum_{i: \hat{P}_{.05} \leq x_i \leq \hat{P}_{.95}} x_i \right) + \text{Lap} \left( \frac{3(\hat{P}_{.95} - \hat{P}_{.05})}{\varepsilon n} \right).$$

- (e) Determine whether or not the following analogue for a Winsorized mean is  $\varepsilon$ -DP: use Part 1c to get  $\varepsilon/3$ -DP estimates  $\hat{P}_{.05}$  and  $\hat{P}_{.95}$  of the 5th and 95th percentiles, and output

$$M(x) = \frac{1}{n} \cdot \left( \sum_{i=1}^n [x_i]_{P_{.05}}^{P_{.95}} \right) + \text{Lap} \left( \frac{3(\hat{P}_{.95} - \hat{P}_{.05})}{\varepsilon n} \right).$$

You do not need to formally prove your answer, but you should at least provide an informal explanation.

- (f) The dataset `MaPUMS5full.csv` provides the 5% PUMS Census file for Massachusetts. For  $\varepsilon = 1$ , compare the RMSE for the algorithms from Parts 1a and 1d as well as the ordinary Laplace mechanism between the actual means of income in each PUMA in Massachusetts, and the DP released means. Also show box-and-whisker plots of the DP released means for each PUMA by these algorithms, noting the true means. You should probably order these by mean income, or perhaps skew of income, or anything you think reveals an interesting pattern. Use a range of  $D = 1,000,000$ .
2. **Composition:** Suppose you have a global privacy budget of  $\varepsilon = 1$  (and are willing to tolerate  $\delta = 10^{-9}$ ) and you want to release  $k$  count queries (i.e. sums of Boolean predicates<sup>1</sup>) using the Laplace mechanism with an individual privacy loss of  $\varepsilon_0$ . By basic composition, you can set  $\varepsilon_0 = \varepsilon/k$ . Using the advanced composition theorem (see statement in the slides or the Dwork–Roth text), you can set  $\varepsilon_0$  to be on the order of  $\varepsilon/\sqrt{k \ln(1/\delta)}$ . We will provide you with code from PSI for the “optimal” composition theorem for differential privacy that calculates the largest value of  $\varepsilon_0$  that ensures global  $(\varepsilon, \delta)$ -DP as a function of  $\varepsilon$ ,  $\delta$ , and  $k$ . For each of these choices, plot (on the same graph) the standard deviation of the Laplace noise added to each query as a function of  $k$ , and find the smallest values of  $k$  where the advanced and optimal composition theorems strictly improve upon the basic composition theorem.
3. **Synthetic Data:** Expanding the template from class, and using again `MaPUMS5Full.csv`, create a DP three-way histogram<sup>2</sup> release of income, education and age. You do not need to graph this histogram, just compute the release for each binned combination of the variables. From this, you should be able to generate synthetic data of these three variables. Run a linear regression as a post-process on your synthetic data, predicting income from education and your additional variable<sup>3</sup> using the equation:

$$\text{Income}_i = \beta_0 + \beta_1 \text{Education}_i + \beta_2 \text{Age}_i + \nu_i; \quad \nu_i \sim \mathcal{N}(0, \sigma^2) \quad (1)$$

Let  $\beta^* = \{\beta_0^*, \beta_1^*, \beta_2^*\}$  be the coefficients in the full sensitive data, while  $\hat{\beta}$  is the vector of coefficients in the some bootstrap of the data, and  $\tilde{\beta}$  the DP release we generate in that bootstrap sample. The mean-squared error of a DP release of  $\tilde{\beta}$  can be decomposed into the contributions of bias and variance as:

$$\text{MSE}(\tilde{\beta}) = \text{bias}(\tilde{\beta})^2 + \text{var}(\tilde{\beta}) = (\mathbb{E}[\beta^* - \tilde{\beta}])^2 + \mathbb{E}[(\tilde{\beta} - \tilde{\beta})^2] \quad (2)$$

<sup>1</sup>A Boolean predicate is a function that returns a 0 or a 1. An example of a count query might be the sum of bits for all college students.

<sup>2</sup>That is, a histogram representation counting the occurrences of having all possible combinations of the three binned variables.

<sup>3</sup>You will likely find that `log(income)` has a more linear relationship with your other two variables, so feel free to shift from `income` to `log(income)` if you prefer. However, you will need to decide how to treat zero values in income; one option is to clip the lower bound of income to some small positive value.

For this calculation, we are taking the (sensitive) regression coefficients  $\beta^*$  on the entire dataset as the true values of  $\beta$ . Show the contributions to MSE of the bias and variance of the DP-regression coefficients as a function of sample size for  $\epsilon = 0.5$ . Compare this to the MSE of  $\hat{\beta}$  which is simply due to finite sample size (and which is all variance and zero bias).<sup>4</sup>

4. **BONUS:** Using your developed understanding of differential privacy, and the described use case in the Gaboardi *et al.* PSI paper, reexamine the deployed instance of the PSI budgeting tool, available at <http://psiprivacy.org>. Provide any feedback that you think would make the interface easier for the intended non-expert “data owner” user to budget a DP-release, or would otherwise improve the system. (Note: Insightful, considered feedback will receive 1/2 point bonus, and feedback that strikes us as revelatory or particularly intriguing idea will receive 1 point bonus and a note of thanks in a future paper draft.)
5. **Final Project:** By **April 9**, submit a couple of pages giving a detailed description of what your final project will look like. You should be able to clearly state your research questions, briefly articulate how your project relates to what has been done in the past, describe the approach you are taking, give your timeline for completing various aspects of the project, and *discuss your fallback plan in case you don’t obtain the results that you’re hoping to obtain.*

---

<sup>4</sup>Assume we are using unclipped data to generate the sensitive values  $\beta^*$  and  $\hat{\beta}$ .