

CS480/680, Spring 2023

Assignment 4

Designer: Shufan Zhang; Instructor: Hongyang Zhang

Released: July 17; Due: July 31, noon

[IMPORTANT] Due to the due date changes, we adjust the marks for this assignment as follows:

- The marks for Question 1 are normalized to 100% for this assignment. If you answer Q1 correctly, you will get full marks for this assignment, which is 10 marks for the overall grades of this course.
- The marks for Question 2 change to bonus for all assignments in this course. If you get Q2 done correctly, you have 3 bonus marks for all the assignments, but the maximum marks you can get for the assignment are capped by 40 contributing to the overall grade of the course.

The new due date is **July 31, noon**.

1. Writing: Differentially Private Data Analytics [100 Marks]

Name	Age	Gender	Num. Matches
Henry	42	Male	8
Sarah	36	Female	16
Austin	22	Non-Binary	5
Adrian	44	Male	12
Natalie	30	Female	5
Chloe	23	Non-Binary	20
Tony	45	Male	13
Christine	28	Non-Binary	20
Olivia	39	Female	35

Table 1: Number of Matches Information

Tinker would like to use differential privacy to publish the data, since it is resilient to background knowledge. To do this, Tinker releases a differentially private histogram showing the number of users having a specific number of matches. To generate this histogram, Laplace noise is added to the true value. For instance, if 4 users in the dataset have 5 matches each, Laplace noise would be added to the total number of such users (4) to hide the true number of users with 5 matches.

The *histogram* representation of the dataset $x = (x_0, \dots, x_{n-1})$, where $n = 36$, is a 36-dimensional vector (ranging from a minimum of 0 matches to the maximum number of matches observed in the dataset, 35) where the j -th entry is the number of x 's rows whose number of matches is equal to j . For instance, according to the data in Table 1, $h_{20}(x) = 2$ since two users in the database (Chloe and Christine) have 20 Tinker matches, and thus:

$$h(x) := (h_0(x), \dots, h_{n-1}(x)), n = 36$$

Consider the following *noisy histogram algorithm* output:

$$\hat{h}(x) := (h_0(x) + L_0, \dots, h_{n-1}(x) + L_{n-1})$$

where every $L_j \sim \text{Laplace}(\lambda)$ is independent Laplace Noise.

Note: Wikipedia gives a good overview of differential privacy and differentially private mechanisms: https://en.wikipedia.org/wiki/Differential_privacy. You may also seek out additional resources to help answer this question.

Marking Scheme:

The student should specify if they follow bounded or unbounded DP. Full marks can be given for either definition if the answer is correct.

- (a) [15 Marks] What is the sensitivity of this query of releasing histograms?

Solution:

(Unbounded DP) Sensitivity is 1.

(Bounded DP) Sensitivity is 2.

- (b) [15 Marks] Tinker sets the parameters to $\epsilon = 0.01$, then what is λ in Laplace Mechanism?

Solution:

(Unbounded DP) $\lambda = 1/0.01 = 100$.

(Bounded DP) $\lambda = 2/0.01 = 200$.

- (c) [15 Marks] Please analyze the expected error of this mechanism. ($\mathcal{E} = \sum_{i=1}^d \mathbb{E}[(o_i - c_i)^2]$, where o_i is the i th entry of the noisy output, and c_i is the i th entry of the true answer.)

Solution:

$$\mathcal{E} = \sum_{i=1}^d \mathbb{E}[(o_i - c_i)^2] \quad (1)$$

$$= \sum_{i=1}^d \text{Var}(\text{Lap}(S(q)/\epsilon)) \quad (2)$$

$$= 2 * d * S(q)^2 / \epsilon^2 \quad (3)$$

(Unbounded DP) 720,000

(Bounded DP) 2,880,000

- (d) [25 Marks] Does this mechanism satisfy the definition of ϵ -differential privacy? Will the histogram output of this mechanism be useful? Justify.

Solution:

(a) The mechanism is ϵ -DP, because: 1) each data point can only go into one bin in the output histogram, 2) therefore for each counting query corresponding to a bin, the data is disjoint from others, 3) by Laplace mechanism, adding noise as per sensitivity of the histogram query will yield ϵ -DP.

(Mentioning “basic composition of 36 counting query” will lose marks.)

(b) (Useful) The sensitivity of this query is low. Therefore the noisy histogram mechanism is useful.

(Not useful) For this instance ($\epsilon = 0.01$), it is not useful because the database size is relatively small and it adds too much noise to the result that overwhelms the signals.

- (e) [30 Marks] Tinker would like to collect new user data with local differential privacy guarantee. Consider a domain $\Sigma = \{l_1, \dots, l_k\}$ of k locations, please design a randomized response R that takes in a true location $l \in \Sigma$ and randomly outputs a location $o \in \Sigma$. (Describe your algorithm and show that the algorithm achieves ϵ -local differential privacy.)

Solution:

The algorithm A is described as follows:

- The user reports the true location l with probability $p = \frac{e^\epsilon}{e^\epsilon + k - 1}$, where $k = |\Sigma|$.
- The user reports the other location in the location domain Σ with probability $q = \frac{1}{e^\epsilon + k - 1}$, where $k = |\Sigma|$.

Privacy proof.

The aforementioned algorithm satisfies ϵ -LDP because:

$$\frac{\Pr[A(l) = o]}{\Pr[A(l') = o]} = \frac{p}{q} = e^\epsilon$$

2. Coding: Private Data Synthesis [Bonus 3 Marks for All Assignments]

The US Census Bureau collects the geographic and demographic data of US residents. The data is anticipated to be used for performing several machine learning or analytic tasks to better understand the residents or incidents of residents, e.g., contact tracing, civic planning, natural disaster rapid response, etc. However, these types of data are considered highly sensitive that contain personally identifiable information (PII) of individuals. Due to privacy laws or regulations, some privacy enhancement techniques should be applied to guarantee individual privacy.

Now we consider using differential privacy (DP) as the means to protect individual privacy. One way to enforce differential privacy over the collected data is **private data synthesis**, meaning generating a synthetic dataset with DP guarantees. A **synthetic dataset** is a collection of artificially generated data that simulates real-world data. Instead of being collected from actual observations or measurements, synthetic data is created using various statistical and computational techniques to mimic the characteristics and patterns of the original data. Enforcing DP in data synthesis requires the data generation algorithm to be proved as differentially private.

Now you are given a dataset called *Adult* (which can be accessed via <https://archive.ics.uci.edu/dataset/2/adult>). This dataset has 48,842 rows and 14 attributes. These attributes (also called the *schema* of the dataset) include age, workclass, fnlwtg, education, education-num, marital-status, occupation, relationship, race, sex, capital-gain, capital-loss, hours-per-week, and native-country. Each row thereby denotes a person. A typical machine learning task on this dataset is to use these 14 attributes as features to predict whether a person makes over 50K a year.

Your tasks for this question are the following:

- [1 Marks] Design a DP synthetic dataset generation algorithm (pseudo-code with step-by-step brief description), and show why it is DP. You can use any generative models you like, including GAN, BayesianNet, and your self-created ones.
- [1 Marks] Generate two DP synthetic datasets (with $\epsilon = \{1, 5\}$ respectively) for the Adult dataset, which should contain the exact same number of rows and 15 columns (14 features + 1 prediction class). Submit the code and the generated datasets.
- [1 Marks] Choose any classification models (e.g., decision tree, etc.) to perform the learning task (on predicting if a person has income over 50K a year), and for 3 datasets (the ground truth and 2 synthetic datasets), train the model on the first 2/3 of the dataset and test it on the rest of the data. Report the accuracy or other reasonable utility metrics (e.g., ROC) for three testings. You can plot the results or simply report the numbers with a brief justification.