

HW 5: Approximate DP, Composition, and Contextual Integrity

CS 208 Applied Privacy for Data Science, Spring 2025

Version 1.0: Due Fri, Mar. 7, 11:59pm.

Instructions: Submit a PDF file that contains both your written responses as well as your code to the assignment on Gradescope. Read the section "Collaboration & AI Policy" in the syllabus for our guidelines regarding the use of LLMs and other AI assistance on the assignments.

1. **Approximate DP:** Below are several mechanisms that are not pure DP. For each, find the smallest value of δ such that they are (ϵ, δ) for a finite ϵ . In each case, also find a value of ϵ such that they are (ϵ, δ) -DP. (Your ϵ need not be the smallest possible value.) Both ϵ and δ may be functions of n . Briefly justify your answers, explaining how you obtained your results; formal proof is not necessary. In all cases, use d_{Ham} as your adjacency relation.

You may find the following characterizations of approximate DP useful. Consider a mechanism $M : \mathcal{X} \rightarrow \mathcal{Y}$ with adjacency relation \sim on \mathcal{X} :

- For $\epsilon \geq 0$, the smallest δ such that M is (ϵ, δ) -DP with respect to \sim is given by

$$\delta = \max_{x \sim x'} \sum_{y \in \mathcal{Y}} \max \{ \Pr[M(x) = y] - e^\epsilon \cdot \Pr[M(x') = y], 0 \}.$$

(In case of continuous mechanisms, the sum should be replaced with an integral and the pmfs replaced with pdfs.)

- For $\epsilon, \delta \geq 0$, M is (ϵ, δ) -DP if for all $x \sim x'$, with probability at least $1 - \delta$ over $y \leftarrow M(x)$, we have $\Pr[M(x) = y] \leq e^\epsilon \cdot \Pr[M(x') = y]$. (Note that this is only a sufficient condition for (ϵ, δ) -DP, not an exact characterization.) Reference for the proof of this sufficient condition: Lemma 1.4. <https://dpcourse.github.io/2021-spring/lecnotes-web/lec-09-gaussian.pdf>

Consider the following mechanisms:

- (a) The mechanism M that takes a dataset $x \in [0, 1]^n$ and returns an estimate of the mean $M(x) = \left(\frac{1}{n} \cdot \sum_{i=1}^n x_i \right) + [Z]_{-1}^1$, for $Z \sim \text{Lap}(2/n)$.
- (b) The mechanism M that takes a dataset $x \in [0, 1]^n$, computes $\bar{x} = \frac{1}{n} \cdot \sum_{i=1}^n x_i$, and outputs 1 with probability \bar{x} and 0 otherwise.
- (c) The mechanism M that takes a dataset $x \in \mathcal{S}^*$ where each record $x_i \in \mathcal{S}$ is an ascii string (e.g. the i 'th individual's favorite surf break) and does the following:
 - i. Calculate the $s \in \mathcal{S}$ that maximizes $n_s = |\{i : x_i = s\}|$ (i.e. the surf break that is the favorite of the most surfers), breaking ties arbitrarily.
 - ii. If $n_s + \text{Lap}(1/\epsilon) \geq 2n/3$ (the top break wins by a supermajority), output s .
 - iii. Otherwise, output "I don't reveal secret spots."

2. Composition:

- (a) 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 predicates¹) using the Laplace mechanism with an individual privacy loss of ε_0 . By basic composition, you can set $\varepsilon_0 = \varepsilon/k$. Using the advanced composition theorem, you can set $\varepsilon_0 = \varepsilon/\sqrt{2k \ln(1/\delta)}$. For the two choices (basic and advanced composition), plot (on the same graph) the standard deviation of the Laplace noise added to each query as a function of k .
- (b) As we saw in class, Wikimedia used a variant of differential privacy called zCDP (Zero-concentrated Differential Privacy) to release statistics on Wikipedia page views. zCDP is tailored to analyzing the Gaussian mechanism and its compositions. The formal definition of zCDP is not needed for this problem, but only that zCDP has a single privacy-loss parameter $\rho \geq 0$ and has the following properties:
- The Gaussian mechanism with noise of variance $\sigma^2 = (\Delta q)^2/2\rho$ is ρ -zCDP, where Δq is the global sensitivity of the query q .
 - Suppose \mathcal{M}_1 satisfies ρ_1 -zCDP and \mathcal{M}_2 satisfies ρ_2 -zCDP. Then their composition $(\mathcal{M}_1, \mathcal{M}_2)$ satisfies $(\rho_1 + \rho_2)$ -zCDP.
 - If a mechanism \mathcal{M} satisfies ρ -zCDP, then for every $\delta > 0$, it satisfies (ε, δ) -DP for $\varepsilon = \rho + \sqrt{4\rho \ln(1/\delta)}$.

We can calculate the smallest value of σ that ensures $(\varepsilon = 1, \delta = 10^{-9})$ -DP when using the above properties to analyze the Gaussian mechanism for answering k counting queries. To see the benefit one gets from using zCDP, plot (on the same graph) the standard deviation of the Gaussian noise added to each query as a function of k using the composition of zCDP against that of basic and advanced composition for approximate DP (from part (a)). From your plot, for what value of k does the Gaussian mechanism outperform advanced composition (from part (a))?

3. **Applying Contextual Integrity:** Imagine a fictional technology company called Coachable. Coachable designs wearable fitness trackers for athletes. Coachable trackers collect data points about users' blood flow and temperature in order to measure their resting heart rate, heart rate variability, and respiratory rate throughout the day and night. These measurements are used to calculate metrics on users' sleep quality (including duration in bed, duration asleep, number of disturbances, length of time spent in different sleep stages, etc.), their level of physical and mental stress, their recovery rate, readiness for activity, and their overall cardiovascular health. In addition, users can log the following information in the Coachable journal to learn how different factors affect their training and performance:

- Alcohol and marijuana consumption
- Supplement use and dosage
- Caffeine consumption
- Medications and sleep aids

¹A Boolean predicate is a function that returns a 0 or a 1. An example of a count query might be the number of Harvard college students that live in the Quad.

- Screen time and bedtime routines
- Air travel
- Stretching and other recovery modalities
- Nutrition and diet plans
- Menstruation and pregnancy
- Sexual activities

Coachable users receive detailed reports on how the behavior logged in their journal affects their athletic training, along with personalized training plans, lifestyle tips, and audio-guided workouts.

- Using the data they collect, suppose Coachable is planning to compute various summary statistics about behaviors, traits of users, and athletic performance. They plan to use these statistics to micro-target ads to its users. For example, they may use these insights to target ads for products including, but not limited to, diet plans, supplements, workout equipment, and recreational activities. They also plan to make these statistics available to researchers, advertisers, and sports recruiters who are interested in the relationship between behaviors, traits of users, and athletic performance.

If we take the Coachable app to be playing a similar role as a human coach, explain how these additional data practices disrupt the informational norm(s) that operate in typical athlete-coach relationships. Explicitly identify the parameters of contextual integrity in your analysis.

- Evaluate the disruptions you identified above. What are the context-specific values and goals of an athlete-coach relationship? How do these disruptions support or undermine these goals? Then, based on your evaluation, state whether you think Coachable should do anything differently with respect to these data practices.
- Now imagine that Coachable plans to compute the summary statistics under differential privacy. How would your analysis and recommendations in Parts 2 and 3 change, if at all? How would deployment decisions, like how ϵ is set, impact your response?

Collaborators

Please list all collaborators for this problem set. ChatGPT and other AI tools should be treated similarly to collaboration with your peers in the class. You may use these tools to help you understand the material and as part of your brainstorming process, but you should not be asking the tools to solve the homework problems for you. If you do use such tools, you must cite them and list the prompts you entered and responses obtained below.