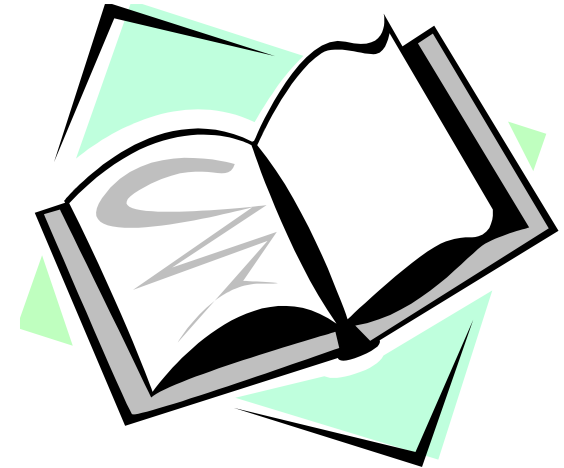


DATA SECURITY AND PRIVACY

Introduction to Differential Privacy

Optional Readings for This Lecture

- Differential Privacy: From Theory to Practice
 - Chapter 2: A Primer on Differential Privacy



Differential Privacy [Dwork et al. 2006]

■ **Definition:** A mechanism A satisfies ϵ -Differential Privacy if and only if

- for any **neighboring** datasets D and D'
- and any possible transcript $t \in \text{Range}(A)$,

$$\Pr[A(D)=t] \leq e^\epsilon \Pr[A(D')=t]$$

- For relational datasets, typically, datasets are said to be **neighboring** if they differ by a single record.

■ **Intuition:**

- Privacy is not violated if one's information is not included in the input dataset
- Output does not overly depend on any single record

Laplace Mechanism Calibrating noise to sensitivity

[DMNS'06]

Given a function $f: D \rightarrow \mathbb{R}^d$ over an arbitrary domain D , the *sensitivity* of f is

$$S(f) = \max_{A, B \text{ where } A \Delta B = 1} \|f(A) - f(B)\|_1$$

Examples:

1. Count: for $f(D) = |D|$, $S(f) = 1$.
2. Sum: for $f(D) = \sum d_i$, where $d_i \in [0, \Lambda]$, $S(f) = \Lambda$.

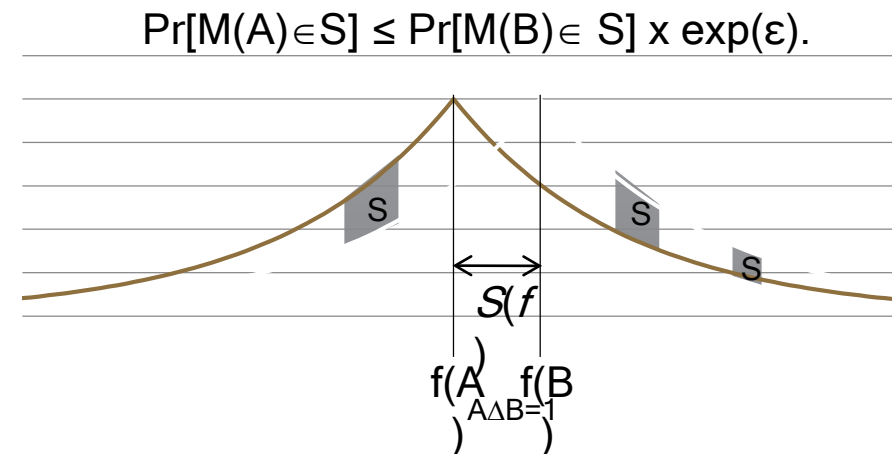
Given a function $f: D \rightarrow \mathbb{R}^d$ over an arbitrary domain D , the computation

$$M(X) = f(X) + (\text{Lap}(S(f)/\epsilon))^d$$

provides ϵ -differential privacy.

Examples:

1. NoisyCount(D) = $|D| + \text{Laplace}(1/\epsilon)$.
2. NoisySum(D) = $\sum d_i + \text{Laplace}(\Lambda/\epsilon)$.



Example of Laplace Mechanism

- Consider an example table of $N=23,450$ records with schema to the right?
- How many tuples are from IN?
 - True count: 546
- Answer while satisfying ϵ_1 -DP: $546 + \text{Lap}(\Delta/\epsilon_1)$
 - $\Delta = 1$
- How many people have score above 23?
- How many?

Name	Score	State
Alice	20	CA
Bob	23	CA
Carl	25	IN
David	18	NY
.....
Frank	20	TX
Jane	14	IN

Counting Queries

- In general, counting queries can be answered relatively accurately
 - Since one tuple affects the result by at most 1
 - A small amount of noise (following the Laplace distribution) can be added to achieve DP

Publishing a histogram

- Suppose we are interested only in the score distribution, then we want to publish the histogram to the right.
- Add $\text{Lap}(\Delta/\epsilon)$ to each of the cell
- What is the sensitivity Δ ?

Score=0	0
Score=17 1313
	2016
	3602
Score=20	1890
	1280
	612
Score=23	221
Score=24	56
Score=25	12

Difference Between Bounded and Unbounded DP

- In unbounded DP, D has one more record than D'
 - $\Delta(\text{histogram}) = 1$
- In bounded DP, D and D' have the same number of records, and only one of them differ
 - $\Delta(\text{histogram}) = 2$

Exponential Mechanism [MT'07]

Let $q: D' \times R \rightarrow \mathbb{R}$ be a query function that, given a database $d \in D'$, assigns a score to each outcome $r \in R$.

Then the exponential mechanism M , defined by

$$M(d, q) = \{\text{return } r \text{ with probability } \propto \exp(\epsilon q(d, r) / 2S(q))\},$$

maintains ϵ -differential privacy.

Reminder: $S(q) = \max_{A, B \text{ where } A \Delta B = 1} \|q(A) - q(B)\|_1$

Motivation: $\Pr(r) \propto \exp\left(\epsilon \frac{q(d, r)}{2S(q)}\right)$

Impact of changing a single record is within ± 1

Example - private vote what to order for lunch:

Option	Score (votes) Sensitivity=1	Sampling Probability		
		$\epsilon=0$	$\epsilon=0.1$	$\epsilon=1$
Pizza	27	0.25	0.4	0.88
Salad	23	0.25	0.33	0.12
Hamburger	9	0.25	0.16	10^{-4}
Pie	0	0.25	0.11	10^{-6}

Example of Exponential Mechanism

- What is the median score?
 - Define $q(D,x) = -|\# \text{ of students with score higher than } x - \# \text{ of students with score lower than } x|$
 - What is the sensitivity?
 - I.e., what is $\max(|q(D,x) - q(D',x)|)$?

Properties of DP

- Sequential Composability
 - If A_1 satisfies ε_1 -DP, and A_2 satisfies ε_2 -DP, then outputting both A_1 and A_2 satisfies $(\varepsilon_1 + \varepsilon_2)$ -DP
- Parallel Composability
 - If D is divided into two parts, applying A_1 and A_2 on the two parts satisfy $(\max(\varepsilon_1, \varepsilon_2))$ -DP
- Post-processing Invariance
 - If A_1 satisfies ε_1 -DP, then $A_2(A_1(\cdot))$ satisfies ε_1 -DP for any A_2

Privacy Budget

- When designing a multiple-step algorithm for ϵ -DP, one needs to divide ϵ into portions so that each step consumes some

Some queries are hard to answer

- Some queries are hard to answer
 - E.g., max, since it can be greatly affected by a single tuple

Four Settings of Satisfying DP

- **Local setting**
 - Do not trust server, perturb data before sending to server
- **Interactive setting**
 - Answer queries as they come, not knowing what the rest of the queries are
- **Single workload**
 - Learn a few parameters
- **Non-interactive publishing**
 - Able to answer a broad range of queries

Limitation of Interactive Setting

- Answering each query consumes some privacy budget
- After answering a pre-determined number of queries, one exhausts the privacy budget, and cannot answer any question anymore
- Problem especially intractable when dealing with multiple users of data

Privacy Preserving Data Publishing

- Design a mechanism A , such that given D , one publishes $T=A(D)$.
- Requirements
 - Privacy friendly
 - Preventing adversaries from learning (individual) information from $O=A(D)$ and A
 - Useful (fidelity-preserving)
 - Allow data users (researchers) to learn (aggregated) information from $O=A(D)$ and A