

Adam Smith and Jonathan Ullman

## 1 Defining “Privacy”

Having seen reconstruction attacks, we now want to get a handle on what it means that some set of statistics are actually ok to release—that they don’t expose individuals’ data (too much?) to attacks like the ones of the last two lectures. The question isn’t new. Researchers in statistics, computer science, and information theory have been tackling variations on it since the 1960’s [War65], and a many specific algorithms and techniques were developed that resist specific suites of attacks.

However, our goal will be to find a general criterion we can use to reason about many different kinds of released information, and about a broad class of attacks. In fact, what we really want is a clear sense in which we’ve prevented “all reasonable” attacks. *Differential privacy* provides one approach to this conundrum. Before we get to it, however, it is helpful to see an example of something that does *not* meet our desiderata.

*k*-Anonymity [Swe02] is one popular approach to reasoning about the privacy implications of publishing statistical tables. It applies only to specific kinds of information, called *generalized microdata*. This means a table of individual records, where each entry is either the original record’s entry (a specific person’s real age, for example) or a set of possible values for that entry (often an interval, like 30–34). Figure 1 gives an example of such a table with age and zip code data. The basic idea of *k*-anonymity is to divide attributes into ‘non-sensitive’ attributes—assumed to be available to an attacker—and ‘sensitive’ ones—assumed to be unknown—and to ensure that *every record matches at least  $k - 1$  others in the nonsensitive attributes*. That is, given an integer *k*, a table is *k*-anonymous if, when we delete the sensitive attributes and leave only the non-sensitive ones, each row appears at least *k* times.<sup>1</sup> The table of Figure 1 is 4-anonymous.

	Non-Sensitive			Sensitive
	Zip code	Age	Nationality	Condition
1	130**	<30	*	AIDS
2	130**	<30	*	Heart Disease
3	130**	<30	*	Viral Infection
4	130**	<30	*	Viral Infection
5	130**	≥40	*	Cancer
6	130**	≥40	*	Heart Disease
7	130**	≥40	*	Viral Infection
8	130**	≥40	*	Viral Infection
9	130**	3*	*	Cancer
10	130**	3*	*	Cancer
11	130**	3*	*	Cancer
12	130**	3*	*	Cancer

Figure 1: A 4-anonymous table.

The idea behind this criterion is that it makes linkage attacks (like the Netflix example from Lecture 1 [NS08]) harder to carry out—if an attacker has access to another table that contains some of the non-sensitive attributes, then each record in a *k*-anonymous table will match at least *k* of the records in

<sup>1</sup>Sweeney’s original notion [Swe02] was actually a bit more permissive: the condition did not have to hold simultaneously for all non-sensitive attributes, but only for those subsets of them that were likely to appear together in other tables. The simpler notion is good enough for our discussion.

the other table.

While  $k$ -anonymity is likely to make those specific attacks harder than they would be with raw data, a  $k$ -anonymous table can still leak lots of individual-level information. We can glean lots of information from the table in Figure 1: everyone in their 30's has cancer; our friend Alice, who's data we happen to know is in the table, cannot have visited the hospital because of a broken leg; etc. Of course, the example is simplistic (real hospital records don't look like the example in the table...) but it illustrates two important points: 1. Defending against one type of attack isn't sufficient, and 2. Criteria that limit the *form* of the output (in this case, the number of occurrences of each vector of non-sensitive attributes) do not necessarily constrain the *information* that is revealed.

**Composition**  $k$ -anonymity illustrates another important point, namely that when the same record is included in two (or more) data sets that are anonymized separately, the combination of the two might reveal far more than the two do individually [GKS08]. For example, consider the table of Figure 2. Suppose we know that Alice's record appears in both tables, and that she is 28 years old and lives in zip code 13012. Neither table on its own pins down her condition exactly (each one narrows it down to a few possibilities), but taken together they pin it down exactly.

This problem is known as *composition*—what happens when many different pieces of information are revealed about me? We return to this question in the next lecture.

	Non-Sensitive			Sensitive
	Zip code	Age	Nationality	Condition
1	130**	<35	*	AIDS
2	130**	<35	*	Tuberculosis
3	130**	<35	*	Flu
4	130**	<35	*	Tuberculosis
5	130**	<35	*	Cancer
6	130**	<35	*	Cancer
7	130**	≥35	*	Cancer
8	130**	≥35	*	Cancer
9	130**	≥35	*	Cancer
10	130**	≥35	*	Tuberculosis
11	130**	≥35	*	Viral Infection
12	130**	≥35	*	Viral Infection

Figure 2: A 6-anonymous table.

**Form versus process (or syntax versus semantics)** Perhaps the most important lesson we can draw from the examples above is that, to come up with a general approach to privacy of statistical data, it isn't enough to restrict the form of the outputs we generate.  $k$ -anonymity specifies a set of acceptable outputs, but doesn't substantially restrict *how* they are produced.

## 2 A First Example: Randomized Response

Let's recall the randomized response mechanism from Lecture 1. Suppose that our data set consists of a single bit  $x_i \in \{0, 1\}$  for each of  $n$  individuals. For each person, we'll generate a biased random bit  $Y_i$  as follows: roll a fair die. If the die comes up 1,2,3, or 4, we set  $Y_i = x_i$ . Otherwise, we set  $Y_i$  to be the opposite value to  $x_i$  (that is,  $Y_i = 1 - x_i$ ). The algorithm's output is the list of values  $(Y_1, \dots, Y_n)$ . Let  $RR_{\text{basic}}$  denote the resulting algorithm (spelled out in Algorithm 1).

What sort of privacy does  $RR_{\text{basic}}$  provide? There are many ways to answer the question, but one way is to think of a sort of *plausible deniability*. For any individual  $i$ , seeing a particular value of  $Y_i$  in the output doesn't give an outsider much information about whether  $x_i = 0$  or  $x_i = 1$ . For any particular

---

**Algorithm 1:** Basic Randomized Repsonse,  $RR_{\text{basic}}$ 

---

**Input:** Data set of  $n$  bits:  $\mathbf{x} = (x_1, \dots, x_n) \in \{0, 1\}^n$

**Output:** Bits  $Y_1, \dots, Y_n$

```
1 for  $i = 1$  to  $n$  do
2    $Y_i = \begin{cases} x_i & \text{w.p. } 2/3; \\ 1 - x_i & \text{w.p. } 1/3; \end{cases}$ 
3 return  $(Y_1, \dots, Y_n)$ 
```

---

output  $y_i$ , we have:

$$\frac{1}{2} \leq \frac{\mathbb{P}(Y_i = y_i \mid x_i = 1)}{\mathbb{P}(Y_i = y_i \mid x_i = 0)} \leq 2 \quad (1)$$

In other words, the outcome would have been roughly as likely if we had changed person  $i$ 's record from one to 0 or vice-versa (assuming everyone else's records were unchanged).

**Proposition 2.1.** *There is a procedure that, given the outputs  $Y_1, \dots, Y_n$  from randomized response on input  $x_1, \dots, x_n$ , returns an estimate  $A$  such that*

$$\sqrt{\mathbb{E}\left(\left(A - \sum_{i=1}^n x_i\right)^2\right)} = O(\sqrt{n}).$$

**Exercise 2.2.** Prove Proposition 2.1.

Now a factor of 2 maybe not be quite satisfactory. But we can get it to be arbitrarily close to 1 by changing the mechanism a bit. Suppose we want that odds ratio to be bounded by  $e^\varepsilon$  for small number  $\varepsilon > 0$ . (Recall that  $e^\varepsilon \approx 1 + \varepsilon$  when  $\varepsilon$  is close to 0.) Algorithm 2 gives a version which takes data in an arbitrary set  $X$ , along with a predicate  $\varphi$  that maps each record to a bit.

---

**Algorithm 2:** Randomized Response,  $RR_\varepsilon$ 

---

**Input:** Data set of  $n$  bits:  $\mathbf{x} = (x_1, \dots, x_n) \in X^n$ , predicate  $\varphi : X \rightarrow \{0, 1\}$ , and a parameter  $\varepsilon > 0$

**Output:** Bits  $Y_1, \dots, Y_n$

```
1 for  $i = 1$  to  $n$  do
2    $Y_i = \begin{cases} \varphi(x_i) & \text{w.p. } \frac{e^\varepsilon}{e^\varepsilon + 1}; \\ \varphi(1 - x_i) & \text{w.p. } \frac{1}{e^\varepsilon + 1}; \end{cases}$ 
3 return  $(Y_1, \dots, Y_n)$ 
```

---

Even though no particular  $\varphi(x_i)$  can be learned with confidence, when  $n$  is large we can use the  $Y_i$ 's to estimate the proportion of records that satisfy  $\varphi$ . Figure 5(a) shows the accuracy of randomized response for estimating the proportion of individuals in some population who satisfy some predicate, for different values of  $n$ . For  $\varepsilon = 0.5$ , the accuracy becomes reasonable for  $n$  in the low thousands.

**Proposition 2.3.** *There is a procedure that, given the outputs  $Y_1, \dots, Y_n$  from randomized response (Alg. 2) on input  $x_1, \dots, x_n$ , returns an estimate  $A$  such that*

$$\sqrt{\mathbb{E}\left(\left(A - \sum_{i=1}^n \varphi(x_i)\right)^2\right)} \leq \frac{e^{\varepsilon/2}}{e^\varepsilon - 1} \sqrt{n}.$$

For bounded  $\varepsilon$  (say, less than 1), this bound is  $\Theta\left(\frac{\sqrt{n}}{\varepsilon}\right)$ .

**Exercise 2.4.** Prove Proposition 2.3. (Hint: For which constants  $a, b$  do we have  $\mathbb{E}(aY_i - b) = x_i$ ?)

**Exercise 2.5.** Strengthen Proposition 2.3 as follows: show that there is a constant  $c > 0$  such that, for every  $t > 1$ , the probability that  $|A - \sum_{i=1}^n \varphi(x_i)| \geq t \frac{e^\varepsilon + 1}{e^\varepsilon - 1} \sqrt{n}$  is at most  $2 \exp(-ct^2)$ . (Hint: Write  $A$  as a sum of independent random variables and apply a Chernoff bound.) (The exact form of the function of  $\varepsilon$  isn't really important here. Anything expression that scales as  $\Theta\left(\frac{1}{\varepsilon}\right)$  will do.)

**Exercise 2.6.** It is typical to analyze a mechanism like randomized response in terms of how well it does at estimating the *average*  $\frac{1}{n} \sum_{i=1}^n \varphi(x_i)$ , instead of the sum (since if we add more data from the same population, the average should stay more or less the same). We can use  $\frac{A}{n}$  (where  $A$  is the estimate from Prop. 2.3) to estimate the average, and its standard deviation will  $\Theta\left(\frac{1}{\varepsilon\sqrt{n}}\right)$  (for  $\varepsilon \leq 1$ ).

Suppose we want the standard deviation of this estimate to be at most  $\alpha$ . For a privacy parameter  $\varepsilon$ , how large a dataset  $n(\alpha, \varepsilon)$  do we need? (Write an explicit function for  $n(\alpha, \varepsilon)$ , as well as a simple asymptotic expression for the setting where  $\varepsilon \leq 1$ ). If we halve  $\alpha$ , what will happen to  $n(\alpha, \varepsilon)$ ? What if we halve  $\varepsilon$  when  $\varepsilon$  is small?

### 3 Differential Privacy

Let  $\mathcal{X}$  be the set of all possible records for each individual. A dataset  $\mathbf{x}$  is thus a multiset<sup>2</sup> of values in  $\mathcal{X}$ . When the size  $n$  is fixed, we may think of it as a list  $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X}^n$ . (It is often convenient to view it as a *histogram*, that is, a function  $\mathcal{X} \rightarrow \mathbb{N}$  that counts the number of occurrences of each possible record in  $\mathcal{X}$ . We will return to this view later; for now, we'll stick with lists.)

**A Thought Experiment** The main idea of DP is to consider a thought experiment in which we compare how an algorithm behaves on a data set  $\mathbf{x}$  with the way it behaves on a hypothetical dataset  $\mathbf{x}'$  in which one person's record has been replaced with some other value.

We say *two data sets are neighbors if they differ in one individual's record*. A simple way to model this is to think of the size  $n$  of data sets as fixed, and to consider two data sets adjacent if one record has been replaced with a different value. For example, if they differ in index  $i$ , we would have:

$$\begin{aligned}\mathbf{x} &= (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) \\ \mathbf{x}' &= (x_1, x_2, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)\end{aligned}$$

Now consider a randomized algorithm  $A$ . For each possible input data set  $\mathbf{x}$ , its output is a random variable  $A(\mathbf{x})$ . We say an algorithm is differentially private if running the algorithm on two neighboring data sets yields roughly the same distribution on outcomes. Specifically, we'll ask that for every set  $E$  of possible outcomes—for example, those outputs from a healthcare study that lead to individual  $i$  being denied health insurance—the probability of an outcome in  $E$  should be the same under  $A(\mathbf{x})$  and  $A(\mathbf{x}')$ , up to a small multiplicative factor. In other words, the algorithm's outcomes should be about the same *whether or not individual  $i$ 's real data was used*.

**Definition 3.1** ( $\varepsilon$ -DP with fixed-size data sets). A randomized algorithm  $A : \mathcal{X}^n \rightarrow \mathcal{Y}$  taking inputs in  $\mathcal{X}^n$  is  $\varepsilon$ -differentially private for size  $n$  data sets if, for every pair of neighboring data sets  $\mathbf{x}, \mathbf{x}'$ , for all events<sup>3</sup>  $E \subseteq \mathcal{Y}$ :

$$\mathbb{P}(A(\mathbf{x}) \in E) \leq e^\varepsilon \cdot \mathbb{P}(A(\mathbf{x}') \in E). \quad (2)$$

<sup>2</sup>A multiset is a set where we keep track of how many times each element appears.

<sup>3</sup>In this course, it is generally fine to think of an event as any subset of the output set. In general, for uncountable output sets like  $\mathbb{R}$ , one restricts attention to a collection of "measurable" sets. Standard texts on probability discuss the issue in detail.

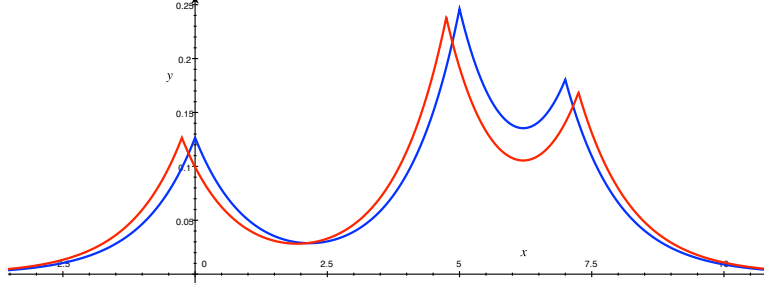


Figure 3: Two distributions  $P$  and  $Q$  that satisfy: for every event  $E$ ,  $P(E) \leq e^{1/4}Q(E)$  and  $Q(E) \leq e^{1/4}P(E)$ .

The definition of DP uses the parameter  $\epsilon$  to control how far apart the distributions of  $A(\mathbf{x})$  and  $A(\mathbf{x}')$  can be. For example, Figure 3 depicts two distributions that satisfy the criterion of Equation (2) with  $\epsilon = 1/4$ . As  $\epsilon$  gets smaller, the algorithm's output distributions can vary less. When  $\epsilon = 0$ , the algorithm leaks nothing at all—its output distribution must be the same for all inputs.

In earlier sections, we pretty much already proved that randomized response (Alg. 2) is differentially private, without using that terminology. Let's go through the argument again, filling in the missing pieces.

**Proposition 3.2.**  *$RR_\epsilon$  is  $\epsilon$ -differentially private.*

*Proof.* Fix two neighboring data sets  $\mathbf{x}$  and  $\mathbf{x}'$ , and let  $i$  be the position in which they differ (so that  $x_i \neq x'_i$  but  $x_j = x'_j$  for all  $j \neq i$ ). First, consider a particular outcome  $\mathbf{y} = (y_1, \dots, y_n)$ . Because we make selections independently for each  $i$ , we have

$$\mathbb{P}(RR_\epsilon(\mathbf{x}) = \mathbf{y}) = \mathbb{P}(Y_1 = y_1 \mid x_1) \cdot \mathbb{P}(Y_2 = y_2 \mid x_2) \cdots \mathbb{P}(Y_n = y_n \mid x_n) \quad (3)$$

When we compare this to the probability that  $RR_\epsilon(\mathbf{x}') = \mathbf{y}$ , only one of the terms in the product will change. We thus get that

$$\frac{\mathbb{P}(RR_\epsilon(\mathbf{x}') = \mathbf{y})}{\mathbb{P}(RR_\epsilon(\mathbf{x}) = \mathbf{y})} = \frac{\mathbb{P}(Y_i = y_i \mid x'_i)}{\mathbb{P}(Y_i = y_i \mid x_i)} \quad (4)$$

This ratio is at most  $\frac{e^\epsilon}{e^\epsilon + 1} \bigg/ \frac{1}{e^\epsilon + 1} = e^\epsilon$ .

Now let's take any subset  $E \subseteq \mathcal{Y} = \{0, 1\}^n$ . The probability that  $RR_\epsilon(\mathbf{x})$  lies in  $E$  is just the sum over  $\mathbf{y} \in E$  of the probability that  $RR_\epsilon(\mathbf{x}) = \mathbf{y}$ . We thus get

$$\mathbb{P}(RR_\epsilon(\mathbf{x}) \in E) = \sum_{\mathbf{y} \in E} \mathbb{P}(RR_\epsilon(\mathbf{x}) = \mathbf{y}) \stackrel{\text{Eq. (4)}}{\leq} \sum_{\mathbf{y} \in E} e^\epsilon \cdot \mathbb{P}(RR_\epsilon(\mathbf{x}') = \mathbf{y}) = e^\epsilon \cdot \mathbb{P}(RR_\epsilon(\mathbf{x}') \in E). \quad (5)$$

This completes the proof.  $\square$

The proof that randomized response is differentially private uses a useful trick that is true quite generally. For discrete output spaces, to prove that an algorithm is  $\epsilon$ -DP, it is enough to show that neighboring data sets lead to each individual output with about the same probability (that is, it suffices to show the condition of Equation (2) for singleton sets). If the condition holds for all singletons, then it holds for all events by summing the probabilities of all of the singletons. Similarly, if the distributions of  $A(\mathbf{x})$  and  $A(\mathbf{x}')$  both have probability densities (on  $\mathbb{R}$ , say), then it suffices to have  $f_{\mathbf{x}}(y) \leq e^\epsilon f_{\mathbf{x}'}(y)$  for all possible outputs  $y$  (where  $f_{\mathbf{x}}(y)$  and  $f_{\mathbf{x}'}(y)$  are the two probability densities).

**Exercise 3.3.** (i) Suppose  $\mathbf{x}$  and  $\mathbf{x}'$  are neighbors. Let  $A$  be a randomized algorithm that is  $\epsilon$ -differentially private for  $\epsilon = \ln(5/4) \approx 0.223$ . Suppose  $A$  outputs real numbers, and suppose  $\mathbb{P}(A(\mathbf{x}) \geq 14) = 0.2$ . What range of values for  $\mathbb{P}(A(\mathbf{x}') \geq 14)$  is possible? (ii) How would the answer change if you knew instead that  $\mathbb{P}(A(\mathbf{x}) \geq 14) = 0.5$ ? [Hint: Consider the constraints on  $\mathbb{P}(A(\mathbf{x}) < 14) = 1 - \mathbb{P}(A(\mathbf{x}) \geq 14)$ .]

## 4 A Second Example: The Laplace Mechanism

Another natural way to add randomness to a computation is to simply add noise to the output of some function  $f$  evaluated on the data. When does this satisfy differential privacy? How does the choice of the function  $f$  we evaluate affect the amount of noise we must add? One way to satisfy DP is to add noise that is scaled to the *global sensitivity* of  $f$ :

**Definition 4.1.** A query  $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$  has global sensitivity  $\Delta$  in the  $\ell_1$  norm if for all  $x_1, \dots, x_n \in \mathcal{X}$ , all indices  $i$ , and all  $x'_i \in \mathcal{X}$ :

$$\|f(x_1, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)\|_1 \leq \Delta$$

The notation  $\|\cdot\|_1$  refers to the  $\ell_1$  norm of a vector, which is sum of the absolute values of the vector's entries. For example,  $\|(1, 0, -3)\|_1 = 4$ , and  $\|(1, 1, \dots, 1)\|_1 = d$ . In 1 dimension, the  $\ell_1$  norm is

$d$  times

just the absolute value.

For example, a proportion  $f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$ , where  $\varphi : \mathcal{X} \rightarrow \{0, 1\}$ , has sensitivity  $\frac{1}{n}$ . The same is true if  $\varphi$  maps records to numbers in the interval  $[0, 1]$ .

To take another example, consider a *histogram*: given a data set  $\mathbf{x} \in \mathcal{X}^n$  and a partition of  $\mathcal{X}$  into  $d$  disjoint sets  $B_1, \dots, B_d$  (think of these as “bins” or “types” of items in  $\mathcal{X}$ ), we count how many records there are of each type.  $f(\mathbf{x}) = (n_1, n_2, \dots, n_d)$  where  $n_j = \#\{i : x_i \in B_j\}$ . So, for example, if we wanted to compute the number of residents of each of the 50 US states from a census of the US population, we would be asking a histogram query. The global sensitivity of a histogram query is at most 2, regardless of how many bins there are.

**Exercise 4.2.** Suppose we want to compute the average of a set of numbers known to lie in the interval  $[0, R]$  on a data set of size  $n$ . What is the sensitivity of the average?

If  $f$  has sensitivity  $\Delta$ , we can satisfy  $\epsilon$ -DP by adding noise from a *Laplace distribution* with scale  $\frac{\Delta}{\epsilon}$ , independently to each entry of the output. The Laplace distribution, also called the double exponential distribution, is sort of a pointy Gaussian (Fig. 4). The mean-0, scale-1 Laplace has density  $h(y) = \frac{1}{2}e^{-|y|}$  for  $y \in \mathbb{R}$ . This distribution has expected absolute value 1, and standard deviation  $\sqrt{2}$ . We can scale the distribution by a positive number  $\lambda > 0$ , to get the general form  $\text{Lap}(\lambda)$  with density

$$\text{Lap}(\lambda) : \quad \text{a distribution on } \mathbb{R} \text{ with p.d.f. } h_\lambda(y) = \frac{1}{2\lambda} \exp(-|y|/\lambda)$$

Figure 4 illustrates the probability density function for a few values of  $\lambda$ . If we translate the distribution to have mean  $\mu$ , then the density becomes  $\frac{1}{2\lambda} \exp(-|y - \mu|/\lambda)$ .

The resulting algorithm (Algorithm 3) is called the Laplace mechanism<sup>4</sup>, and is a basic building block for the design of many other differentially private algorithm. If used to estimate a proportion, it produces far more accurate estimates than randomized response (for the same privacy budget). See Figure 5 and Exercise 4.6.

<sup>4</sup>Here and elsewhere in the course, the term “mechanism”, inherited from literature on game theory, just means “algorithm”.

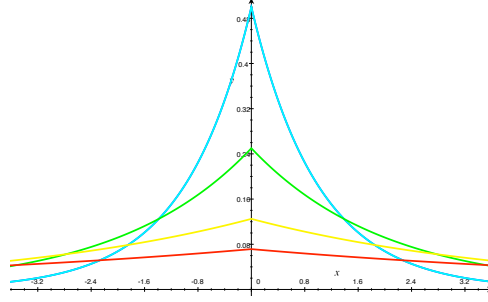


Figure 4: The probability density function of the Laplace distribution with  $\lambda \in \{1, 2, 4, 7\}$ .

---

**Algorithm 3:** Laplace mechanism  $A_{\text{Lap}}(\epsilon, \mathbf{x})$

---

**Input:** Data set  $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X}^n$  and parameter  $\epsilon > 0$ .

- 1 Receive a query  $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$  with global sensitivity  $\Delta$  from analyst;
  - 2 **return**  $f(\mathbf{x}) + (Z_1, \dots, Z_d)$  where  $Z_i \sim \text{Lap}(\frac{\Delta}{\epsilon})$  are i.i.d.
- 

**Theorem 4.3.** *The Laplace mechanism is  $\epsilon$ -differentially private.*

*Proof.* Fix two neighboring data sets  $\mathbf{x}$  and  $\mathbf{x}'$  in  $\mathcal{X}^n$ , and a query  $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$ . Let  $\Delta$  be the global sensitivity of  $f$ . Let  $\mu = f(\mathbf{x})$  and  $\mu' = f(\mathbf{x}')$ . We know that the  $\ell_1$  norm of  $\mu - \mu'$  is at most  $\Delta$ . Comparing the output distributions of  $A_{\text{Lap}}$  on  $\mathbf{x}$  and  $\mathbf{x}'$  thus means comparing two Laplace distributions that have been shifted relative to one another by  $\mu - \mu'$ .

Because we add noise independently to each entry of the output, the density of the output at vector  $\mathbf{y}$  on input  $\mathbf{x}$  can be written as a product:

$$h_{\mathbf{x}}(\mathbf{y}) = \frac{\epsilon}{2\Delta} e^{-\frac{\epsilon}{\Delta} |y_1 - \mu_1|} \times \frac{\epsilon}{2\Delta} e^{-\frac{\epsilon}{\Delta} |y_2 - \mu_2|} \times \dots \times \frac{\epsilon}{2\Delta} e^{-\frac{\epsilon}{\Delta} |y_d - \mu_d|}, \quad (6)$$

which can be simplified to  $h_{\mathbf{x}}(\mathbf{y}) = \left(\frac{\epsilon}{2\Delta}\right)^d e^{-\frac{\epsilon}{\Delta} \|\mathbf{y} - \mu\|_1}$ . If we look at the ratio of the densities for  $\mathbf{x}$  and  $\mathbf{x}'$  at the same output  $\mathbf{y}$ , we get

$$\frac{h_{\mathbf{x}'}(\mathbf{y})}{h_{\mathbf{x}}(\mathbf{y})} = e^{-\frac{\epsilon}{\Delta} (\|\mathbf{y} - \mu'\|_1 - \|\mathbf{y} - \mu\|_1)}. \quad (7)$$

By the triangle inequality, the difference  $\|\mathbf{y} - \mu\|_1 - \|\mathbf{y} - \mu'\|_1$  is at most  $\|\mu - \mu'\|_1$ , which is at most  $\Delta$ . We thus get:

$$\frac{h_{\mathbf{x}'}(\mathbf{y})}{h_{\mathbf{x}}(\mathbf{y})} \leq e^{\frac{\epsilon}{\Delta} \|\mu - \mu'\|_1} \leq e^{\frac{\epsilon}{\Delta} \cdot \Delta} = e^{\epsilon}. \quad (8)$$

And that is enough to conclude the mechanism is  $\epsilon$ -DP! (For any measurable set  $E$ , we have  $\Pr(A_{\text{Lap}}(\epsilon, \mathbf{x}) \in E) = \int_{\mathbf{y} \in E} h_{\mathbf{x}}(\mathbf{y})$ . So if the ratio of the densities is bounded everywhere by  $e^{\epsilon}$ , then so is the ratio of the probabilities of any given event  $E$ .)  $\square$

Thus, we can guarantee differential privacy by adding noise to the output of a function that scales with the function's sensitivity. Histograms, for example, fit the framework well: we can get away with adding error whose expected magnitude is  $2/\epsilon$  to each of the bin counts, regardless of the number of bins.

The following lemma provides useful bounds on the magnitude of the Laplace mechanism's error:

**Proposition 4.4.** *The 1-dimensional Laplace mechanism (with  $d = 1$ ) satisfies the following:*



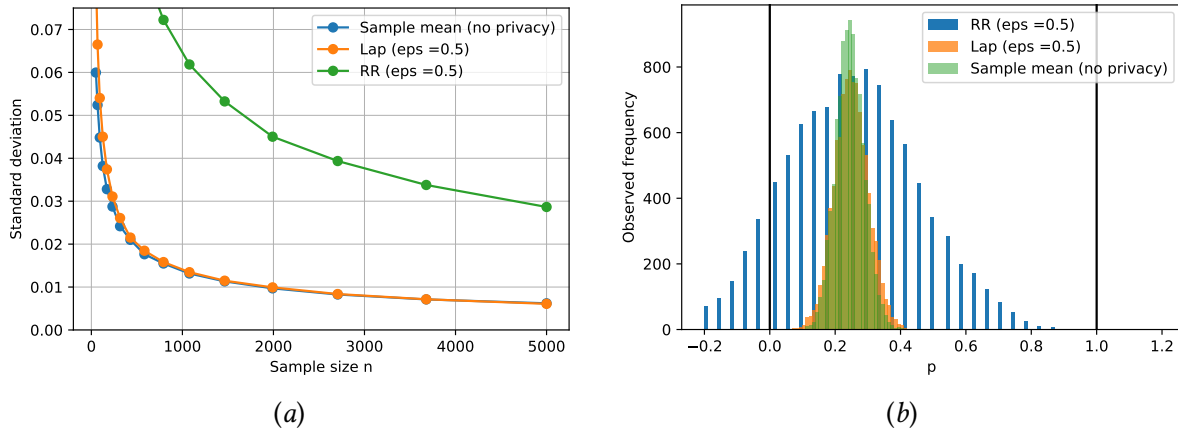


Figure 5: (a) The standard deviation of several different algorithms for estimating the bias of a coin from  $n$  independent flips of the coin. The non-private algorithm simply reports the sample mean (i.e., the proportion of 1's in the data). The other two curves are generated by applying randomized response and the Laplace mechanism with  $\epsilon = 0.5$ . The data are generated from a coin with bias 0.25. The standard deviations were computed from 10,000 independent runs of each algorithm, at sample sizes ranging from 50 to 5,000. (b) The empirical distribution from 10,000 executions of the three algorithms with a sample size of  $n = 100$ , using  $\epsilon = 0.5$ . Note that randomized response (using the simplest post-processing) generates estimates outside of  $[0, 1]$ .

1.  $\mathbb{E}(|A(\mathbf{x}) - f(\mathbf{x})|) = \Delta/\epsilon$
2.  $\sqrt{\mathbb{E}((A(\mathbf{x}) - f(\mathbf{x}))^2)} = \sqrt{2}\Delta/\epsilon$
3. For every  $t > 1$ :  $\mathbb{P}(|A(\mathbf{x}) - f(\mathbf{x})| > \frac{t\Delta}{\epsilon}) \leq \exp(-t)$ .

**Exercise 4.5.** Prove Proposition 4.4.

**Exercise 4.6.** Suppose we use the Laplace mechanism to estimate the fraction  $f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$  where  $\varphi : \mathcal{X} \rightarrow [0, 1]$  is a predicate. Give an expression for the data set size  $n(\alpha, \epsilon)$  at which the Laplace mechanism's root mean square error  $\sqrt{\mathbb{E}((A(\mathbf{x}) - f(\mathbf{x}))^2)}$  drops below  $\alpha$  (when run with privacy parameter  $\epsilon$ ). How does this compare to the analogous calculation for Randomized Response (Exercise 2.6)? What happens to  $n(\alpha, \epsilon)$  when  $\alpha$  is halved? When  $\epsilon$  is halved?

Notice that the accuracy of the Laplace mechanism is pretty bad when  $\epsilon$  is very small. Suppose we want to estimate a fraction (as in the Exercise 4.6). If we set  $\epsilon = 1/n$ , then the standard deviation of the Laplace mechanism is  $\sqrt{2} \cdot \frac{\Delta}{1/n} = \sqrt{2}$  (since  $\Delta = 1/n$  in this case). But the fraction can only take values between 0 and 1—the noise is thus of larger magnitude than the “signal” one is trying to release. This feature is inherent. As we will see next lecture, differentially private algorithms cannot yield any useful information when  $\epsilon < 1/n$ .

## Summary

### 4.1 Key Points

- To reason about information leakage and confidentiality we must look at the algorithms that process the data, not only the form of the output.



- Differential privacy is one way to quantify how much an algorithm leaks about individual inputs. It is parametrized by a positive real number  $\epsilon$ , which bounds the amount of leakage.
- Randomized response and the Laplace mechanism satisfy  $\epsilon$ -DP. For the same privacy budget, the Laplace mechanism adds far less error.

### Additional Reading and Watching

- The paper that defined differential privacy [DMNS06, DMNS16]
- A nontechnical introduction to DP [WAB<sup>+</sup>18]
- Dwork and Roth book, Chapter 2 [DR14]
- Videos from the MinutePhysics Youtube channel: “Protecting Privacy with MATH” and “When It’s OK to Violate Privacy”, 2019.
- Tutorial talks by
  - K. Ligett: [NeurIPS 2016 Tutorial](#)
  - A. Smith: [NASIT 2019 Tutorial](#)
- For further discussion of composition attacks, see [GKS08]

### References

- [DMNS06] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Conference on Theory of Cryptography*, TCC ’06, 2006.
- [DMNS16] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. *Journal of Privacy and Confidentiality*, 7(3), 2016.
- [DR14] Cynthia Dwork and Aaron Roth. *The Algorithmic Foundations of Differential Privacy*. NOW Publishers, 2014.
- [GKS08] Srivatsava Ranjit Ganta, Shiva Prasad Kasiviswanathan, and Adam Smith. Composition attacks and auxiliary information in data privacy. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’08, 2008. ACM.
- [NS08] Arvind Narayanan and Vitaly Shmatikov. Robust de-anonymization of large sparse datasets. In *IEEE Symposium on Security and Privacy*, 2008.
- [Swe02] Latanya Sweeney. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(5):557–570, 2002.
- [WAB<sup>+</sup>18] Alexandra Wood, Micah Altman, Aaron Bembenek, Mark Bun, Marco Gaboardi, James Honaker, Kobbi Nissim, David O’Brien, Thomas Steinke, and Salil Vadhan. Differential privacy: A primer for a non-technical audience. *Vanderbilt Journal of Entertainment & Technology Law*, 21(1):209, 2018.
- [War65] Stanley L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965.

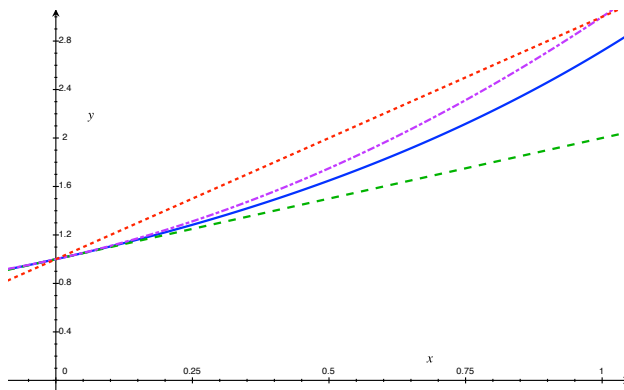


Figure 6: The function  $e^x$  (solid blue), seen here bounded below by  $1+x$  (green dashed) and bounded above on  $[0, 1]$  by  $1+x+x^2$  (purple dashed) and  $1+2x$  (red dashed).

## Appendix

### A The function $e^x$

We'll be working with the quantity  $e^x$  (often  $e^\varepsilon$  for DP algorithms) a lot. Here a few useful inequalities:

- For all  $x \in \mathbb{R}$ , we have  $e^x > 1+x$  (and thus  $e^{-x} \geq 1-x$ ).
- As  $x \rightarrow 0$  (either positive or negative), we have  $e^x = 1+x + \Theta(x^2)$ . As a consequence, we have:
  - $x \geq 1 - e^{-x} \geq x - O(x^2)$ ,
  - $\frac{1}{x} \geq \frac{1}{e^x - 1} \geq \frac{1}{x} - O(x^2)$ , and
  - $\frac{1}{x} \leq \frac{1}{1 - e^{-x}} \leq \frac{1}{x} + O(x^2)$ .

You can double check the direction of inequalities and a sense of specific constants by using a graphing app. For example, Figure 6 shows that  $1+x \leq e^x \leq 1+x+x^2 \leq 1+2x$  for  $x \in [0, 1]$ .