

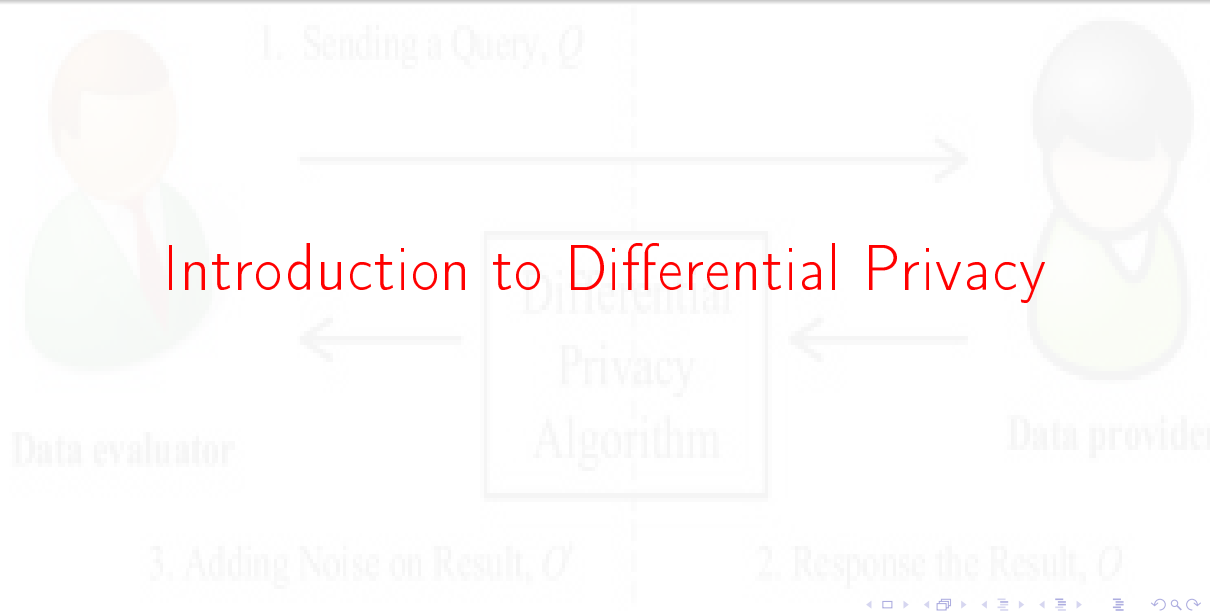
# A survey on Differential Privacy

Sarbajit Ghosh

Under Guidance of Dr. Srimanta Bhattacharya  
Indian Statistical Institute, Kolkata



# Introduction to Differential Privacy



# Security

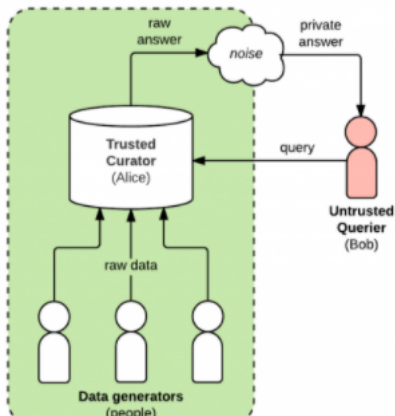
**What is differential privacy?** Differential privacy provides generic mechanism to anonymize non-private target functions viz. statistics, estimation procedure etc.



# Security

**What is differential privacy?** Differential privacy provides generic mechanism to anonymize non-private target functions viz. statistics, estimation procedure etc.

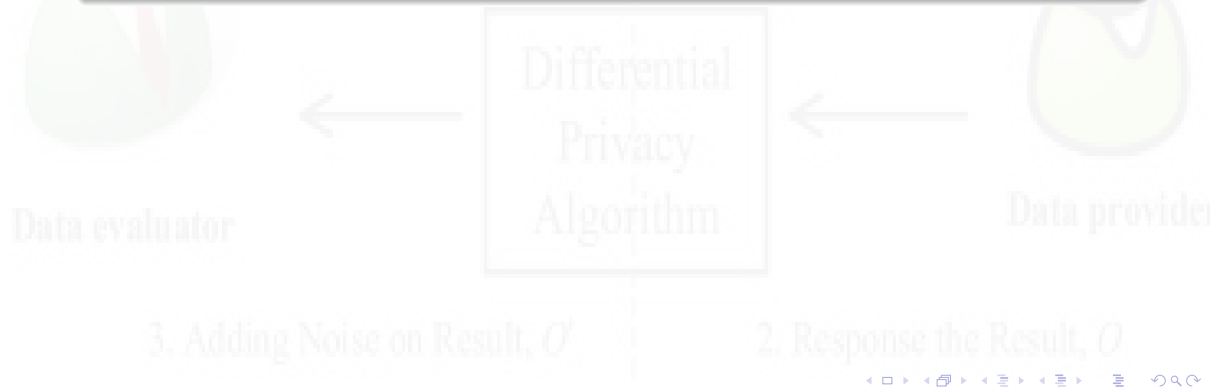
- **Security Set up:** Trusted curator model



# Definitions

## Definition

**Neighbouring Datasets:** Two datasets  $X$  and  $X'$  are said to be neighbouring datasets if they differ only in one row.



# Definitions

## Definition

**Neighbouring Datasets:** Two datasets  $X$  and  $X'$  are said to be neighbouring datasets if they differ only in one row.

## Definition

**$\epsilon$ -differential privacy:** Let  $X, X' \in \mathcal{X}^n$  be two neighbouring datasets also denoted by  $X \sim X'$  and  $M : \mathcal{X}^n \rightarrow \mathcal{Y}$  be an algorithm said to be  $\epsilon$ -differentially private, if for all neighbouring datasets  $X, X'$  and for all  $T \subseteq \mathcal{Y}$ ,

$$\Pr[M(X) \in T] \leq e^\epsilon \Pr[M(X') \in T]$$

Where  $\mathcal{X}^n$  and  $\mathcal{Y}$  is the set of all datasets and set of all queries respectively. This definition is due to Dwork, McSherry, Nissim and Smith in 2006.

## 1. Sending a Query, $Q$

### Some Points on the definition of differential privacy:

- $\epsilon$ -differential privacy is quantitative in nature.

Data evaluator



Data provider

## 3. Adding Noise on Result, $O'$

## 2. Response the Result, $O$

## 1. Sending a Query, $Q$

### Some Points on the definition of differential privacy:

- $\epsilon$ -differential privacy is quantitative in nature.
- Small value of  $\epsilon$  implies strong privacy.

Data evaluator



Data provider

## 3. Adding Noise on Result, $O'$

## 2. Response the Result, $O$



## Some Points on the definition of differential privacy:

- $\epsilon$ -differential privacy is quantitative in nature.
- Small value of  $\epsilon$  implies strong privacy.
- Grants privacy at an individual label. Thus by its nature it provides privacy for all  $X \sim X'$ . Hence privacy of outliers are protected.

## 1. Sending a Query, $Q$

### Some Points on the definition of differential privacy:

- $\epsilon$ -differential privacy is quantitative in nature.
- Small value of  $\epsilon$  implies strong privacy.
- Grants privacy at an individual label. Thus by its nature it provides privacy for all  $X \sim X'$ . Hence privacy of outliers are protected.
- For small value of  $\epsilon$ ,  $e^\epsilon \simeq (1 + \epsilon)$ . So for small value of  $\epsilon$  here we can visualise that both probabilities are very close on neighbouring databases.

Data evaluator

Algorithm

Data provider

## 3. Adding Noise on Result, $O'$

## 2. Response the Result, $O$

## Some Points on the definition of differential privacy:

- $\epsilon$ -differential privacy is quantitative in nature.
- Small value of  $\epsilon$  implies strong privacy.
- Grants privacy at an individual label. Thus by its nature it provides privacy for all  $X \sim X'$ . Hence privacy of outliers are protected.
- For small value of  $\epsilon$ ,  $e^\epsilon \simeq (1 + \epsilon)$ . So for small value of  $\epsilon$  here we can visualise that both probabilities are very close on neighbouring databases.
- The fact  $e^{\epsilon_1} \cdot e^{\epsilon_2} = e^{\epsilon_1 + \epsilon_2}$ , is helpful when we consider group privacy.

## Some Points on the definition of differential privacy:

- $\epsilon$ -differential privacy is quantitative in nature.
- Small value of  $\epsilon$  implies strong privacy.
- Guarantees privacy at an individual label. Thus by its nature it provides privacy for all  $X \sim X'$ . Hence privacy of outliers are protected.
- For small value of  $\epsilon$ ,  $e^\epsilon \simeq (1 + \epsilon)$ . So for small value of  $\epsilon$  here we can visualise that both probabilities are very close on neighbouring databases.
- The fact  $e^{\epsilon_1} \cdot e^{\epsilon_2} = e^{\epsilon_1 + \epsilon_2}$ , is helpful when we consider group privacy.
- The definition is symmetric, hence the role of the both datasets  $X$  and  $X'$  can be interchanged.

## Some Points on the definition of differential privacy:

- $\epsilon$ -differential privacy is quantitative in nature.
- Small value of  $\epsilon$  implies strong privacy.
- Grants privacy at an individual label. Thus by its nature it provides privacy for all  $X \sim X'$ . Hence privacy of outliers are protected.
- For small value of  $\epsilon$ ,  $e^\epsilon \simeq (1 + \epsilon)$ . So for small value of  $\epsilon$  here we can visualise that both probabilities are very close on neighbouring databases.
- The fact  $e^{\epsilon_1} \cdot e^{\epsilon_2} = e^{\epsilon_1 + \epsilon_2}$ , is helpful when we consider group privacy.
- The definition is symmetric, hence the role of the both datasets  $X$  and  $X'$  can be interchanged.

# Properties: Post-Processing

## Theorem

*If  $M : \mathcal{X}^n \rightarrow \mathcal{Y}$  be  $\epsilon$ -differentially private mechanism, and if  $G : \mathcal{Y} \rightarrow \mathcal{Z}$  be any randomized mapping then  $M \circ G$  is a  $\epsilon$ -differentially private mechanism.*



# Properties: Post-Processing

## Theorem

*If  $M : \mathcal{X}^n \rightarrow \mathcal{Y}$  be  $\epsilon$ -differentially private mechanism, and if  $G : \mathcal{Y} \rightarrow \mathcal{Z}$  be any randomized mapping then  $M \circ G$  is a  $\epsilon$ -differentially private mechanism.*

## Proof.

$G$  is a randomized function, we can consider it to be distributed uniformly over all deterministic function  $g$ . Now consider  $X, X'$  are two neighbouring databases and  $T \subseteq \mathcal{Z}$ ,

$$\begin{aligned} & \Pr[G(M(X)) \in T] \\ &= \mathbb{E}_{g \sim G}[\Pr[M(X) \in g^{-1}(T)]] \\ &\leq \mathbb{E}_{g \sim G}[e^\epsilon \Pr[M(X') \in g^{-1}(T)]] \\ &= e^\epsilon \Pr[G(M(X')) \in T] \end{aligned}$$

Hence the proof. □ 🔍 ↺

# Properties: Group Privacy I

## Theorem

If  $M : \mathcal{X}^n \rightarrow \mathcal{Y}$  be  $\epsilon$ -differentially private mechanism, and  $X, X'$  are two neighbouring databases, such that they differ in exactly  $k$ -positions. Then  $\forall T \subseteq \mathcal{Y}$ ,

$$\Pr[M(X) \in T] \leq e^{k\epsilon} \Pr[M(X') \in T]$$



# Properties: Group Privacy II

## Proof.

Since the data bases  $X, X'$  differ in  $k$  position so there should exist intermediate databases which differ in exactly one row. so we define  $X = X_0, X_1, \dots, X' = X_k$  a sequence of databases each has an extra row from the former database.

Then  $\forall T \subseteq \mathcal{Y}$  we have,

$$\begin{aligned} \Pr[M(X) \in T] &= \Pr[M(X_0) \in T] \\ &\leq e^\epsilon \Pr[M(X_1) \in T] \\ &\leq e^{2\epsilon} \Pr[M(X_2) \in T] \\ &\dots \\ &\leq e^{k\epsilon} \Pr[M(X_k) \in T] \\ &= e^{k\epsilon} \Pr[M(X') \in T] \end{aligned}$$

Hence the proof. □

# Properties: Group Privacy I

## Theorem

*Let us consider  $M = (M_1, M_2, \dots, M_k)$  is a sequence of  $\epsilon$ -differentially private mechanism, may have been chosen adaptively then,  $M$  is  $k\epsilon$ -differentially private.*

## Properties: Group Privacy II

### Proof.

Let us take two fixed databases  $X, X'$ , also a sequence of outputs  $y = (y_1, y_2, \dots, y_k)$ . We have the following,

$$\begin{aligned} & \frac{\Pr[M(X) = y]}{\Pr[M(X') = y]} \\ &= \prod_{i=1}^k \frac{\Pr[M_i(X) = y_i | (M_1(X), M_2(X), \dots, M_{i-1}(X)) = (y_1, y_2, \dots, y_{i-1})]}{\Pr[M_i(X') = y_i | (M_1(X'), M_2(X'), \dots, M_{i-1}(X')) = (y_1, y_2, \dots, y_{i-1})]} \\ & \leq \prod_{i=1}^k \exp(\epsilon) \\ & = \exp(k\epsilon) \end{aligned}$$

Hence the proof. □

# Relationship with Hypothesis testing I

Let output  $Y$  is generated from some differentially private mechanism  $M$  operating on some  $X \sim X'$ . An adversary wants to distinguish,

$$\begin{cases} H_0 : Y \text{ came from } X \\ H_1 : Y \text{ came from } X' \end{cases}$$

Now in statical set-up we may consider this as a hypothesis testing. Now intuitively differential privacy grantees that the adversary will not have any significant advantage than random guessing.

Now consider

$$\begin{cases} p := \Pr[\text{Adversary predicts } H_1 | H_0 \text{ true}] \\ q := \Pr[\text{Adversary predicts } H_0 | H_1 \text{ true}] \end{cases}$$

i.e. probability of false positive and false negative respectively.

## Relationship with Hypothesis testing II

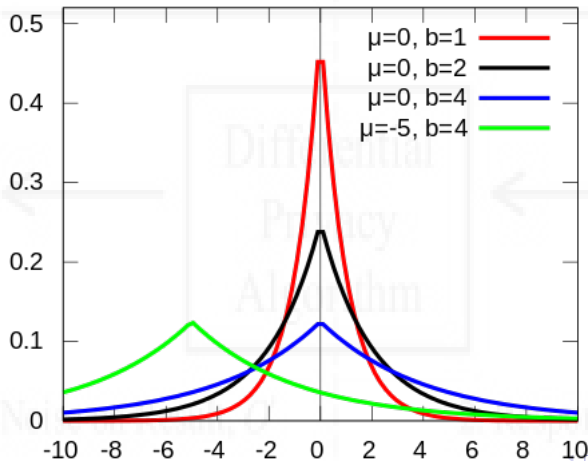
Now  $\epsilon$  differential privacy simultaneously impels that,

$$\begin{cases} p + e^\epsilon q \geq 1 \\ pe^\epsilon + q \geq 1 \end{cases}$$

The above equation is due to Wasserman and Zhou [Ref 3].

This equation intuitively tells that when  $\epsilon = 0$ , the advantage is also most negligible from blind guessing. But as the  $\epsilon$  increases the adversary has an advantage rather than blind guessing.

# Laplace Distribution



# Laplace distribution

## Definition

**Laplace Distribution:** The density of Laplace distribution with location and scale parameter  $\mu, b$  is given by

$$p(x) = \frac{1}{2b} \exp \left\{ -\frac{|x - \mu|}{b} \right\}$$

## Some properties of exponential distribution

- It is symmetric about its scale parameter, if scale  $\mu = 0$  it is symmetric about 0.
- It has variance  $2b^2$
- Tail probability of Laplace distribution is proportional to  $\exp \{-k|x|\}$ . So its probability is concentrated towards centre and it has decaying tail.

1. Sending a Query,  $Q$

# Laplace Mechanism

## Intuition behind the Laplace distribution

Assume for example

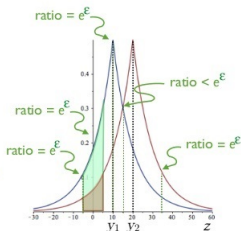
- $\Delta_f = |f(x_1) - f(x_2)| = 10$
- $y_1 = f(x_1) = 10, y_2 = f(x_2) = 20$

Then:

- $dP_{y_1}(z) = \frac{\epsilon}{2 \cdot 10} e^{\frac{|z-10|}{10} \epsilon}$
- $dP_{y_2}(z) = \frac{\epsilon}{2 \cdot 10} e^{\frac{|z-20|}{10} \epsilon}$

The ratio between these distribution is

- $= e^\epsilon$  outside the interval  $[y_1, y_2]$
- $\leq e^\epsilon$  inside the interval  $[y_1, y_2]$



Data provider

3. Adding Noise on Result,  $O'$

2. Response the Result,  $O$



# Laplace Mechanism

## Definition

$l_1$  **sensitivity**: Consider  $f : \mathcal{X}^n \rightarrow \mathbb{R}^k$ . The  $l_1$  sensitivity of  $f$  is

$$\Delta_f = \max_{X, X'} \|f(X) - f(X')\|,$$

where  $X$  and  $X'$  are neighbouring data bases and  $\|\cdot\|$  is  $l_1$  norm.

# Laplace Mechanism

## Definition

$l_1$  **sensitivity**: Consider  $f : \mathcal{X}^n \rightarrow \mathbb{R}^k$ . The  $l_1$  sensitivity of  $f$  is

$$\Delta_f = \max_{X, X'} \|f(X) - f(X')\|,$$

where  $X$  and  $X'$  are neighbouring data bases and  $\|\cdot\|$  is  $l_1$  norm.

**Why sensitivity?**

# Laplace Mechanism

## Definition

**$l_1$  sensitivity:** Consider  $f : \mathcal{X}^n \rightarrow \mathbb{R}^k$ . The  $l_1$  sensitivity of  $f$  is

$$\Delta_f = \max_{X, X'} \|f(X) - f(X')\|,$$

where  $X$  and  $X'$  are neighbouring data bases and  $\|\cdot\|$  is  $l_1$  norm.

**Why sensitivity?** Main object of differential privacy is to preserve individuals privacy. So sensitivity of a function is the most natural choice to analyse.

# Laplace Mechanism

## Definition

**$l_1$  sensitivity:** Consider  $f : \mathcal{X}^n \rightarrow \mathbb{R}^k$ . The  $l_1$  sensitivity of  $f$  is

$$\Delta_f = \max_{X, X'} \|f(X) - f(X')\|,$$

where  $X$  and  $X'$  are neighbouring data bases and  $\|\cdot\|$  is  $l_1$  norm.

**Why sensitivity?** Main object of differential privacy is to preserve individuals privacy. So sensitivity of a function is the most natural choice to analyse.

## Definition

**Laplace Mechanism:** Consider  $f : \mathcal{X}^n \rightarrow \mathbb{R}^k$ . Then Laplace mechanism is defined as

$$M(X) = f(X) + (Y_1, Y_2, \dots, Y_k),$$

where the  $Y_i$  are independent  $\text{Laplace}(\frac{\Delta}{\epsilon})$  random variables.

# Laplace Mechanism is differentially private I

## Theorem

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

## Laplace Mechanism is differentially private II

Proof.

Consider two neighbouring databases  $X$  and  $Y$ , also assume that  $p_X$  and  $p_Y$  be two pdfs corresponding to  $M(X)$  and  $M(Y)$  over the support  $\mathbb{R}^k$ .

$$\begin{aligned}\frac{p_X(z)}{p_Y(z)} &= \frac{\prod_{i=1}^k \exp(-\frac{\epsilon|f(X)_i - z_i|}{\Delta})}{\prod_{i=1}^k \exp(-\frac{\epsilon|f(Y)_i - z_i|}{\Delta})} \\&= \prod_{i=1}^k \exp(-\frac{\epsilon(|f(X)_i - z_i| - |f(Y)_i - z_i|)}{\Delta}) \leq \prod_{i=1}^k \exp(-\frac{\epsilon(|f(X)_i - f(Y)_i|)}{\Delta}) \\&= \exp(-\frac{\sum_{i=1}^k \epsilon(|f(X)_i - f(Y)_i|)}{\Delta}) = \exp(-\frac{\epsilon\|f(X) - f(Y)\|_1}{\Delta}) \leq \exp(-\epsilon)\end{aligned}$$

The third lines follows from triangle inequality, and the last line follows using the definition of  $l_1$  sensitivity of a function. □

# An example I

Effect of  $\epsilon$ -differential privacy on mean function.

Consider the mean function,

$$f = \sum_{i=1}^k X_i; X_i \in \{0, 1\}$$

Here we are considering  $X_i$ 's as indicator variable of a person's habit. Clearly the sensitivity of function  $f$  will be  $\Delta = \frac{1}{n}$  because presence and absence of a person will contribute 1 in the numerator.

Now if we want estimate  $p$  the ratio of persons having that habit in the database the by Laplace mechanism we have

$$\hat{p} = f(X) + Y$$

Where  $Y$  is a is  $\text{Laplace}(\frac{\Delta}{\epsilon})$  random variable and  $X$  is the  $k$ -tuple vector. Now  $(\frac{\Delta}{\epsilon}) = (\frac{1}{n\epsilon})$ .

## An example II

Now, we observe that  $\mathbb{E}(Y) = 0$ , since location parameter is zero. Also  $p = f(X)$  the original ratio from the database. Therefore we have,

$$\mathbb{E}(\hat{p}) = p$$

### Definition

**Tchebychev's inequality** If  $X$  is a RV with finite expectation  $\mu$  and variance  $\sigma^2 \geq 0$  then for all real number  $k > 0$ ,

$$Pr[|X - \mu| \leq k\sigma] \geq 1 - \frac{1}{k^2}$$

Let us calculate variance of  $\hat{p}$ ,



## An example III

$$\begin{aligned}\text{Var}[\hat{p}] &= \text{Var}[Y + f(X)] \\ &= \text{Var}[Y] \\ &= \frac{1}{\epsilon^2 n^2}\end{aligned}$$

Now if we apply Tchebychev's inequality with the random variable as  $\hat{p}$  we get,

$$\Pr[|\hat{p} - p| \leq \frac{k}{n\epsilon}] \geq 1 - \frac{1}{k^2}; \forall k > 0$$

That is we have,

$$|\hat{p} - p| \leq \mathcal{O}\left(\frac{1}{n\epsilon}\right)$$

with reasonable probability, and also it differentially private since Laplace mechanism is differentially private.

# Counting query & Histogram query I

**Single counting query:** In a database there is a binary column  $P$ , corresponds to a property entries of database[each row]. We want to know  $f = \sum_i X_i$ ,  $X_i \in \{0, 1\}$  if  $i^{th}$  row has property  $P$ . So here  $f(X) + \text{Laplace}(\frac{1}{\epsilon})$  is the privatization statistic, With error  $\mathcal{O}(\frac{1}{\epsilon})$ .

**Multiple counting query:** Consider  $k$  counting query,  $f = (f_1, f_2, \dots, f_k)$ , which are fixed before[non-adaptive].  $f(X) + Y$  will be privatization statistic. Where  $Y = (Y_1, Y_2, \dots, Y_k)$  are independent Laplace random variables.

**Which scale parameter should we use?** Each query  $f_i$  has sensitivity 1[as counting query], and the underlying database is same so changing an individual person will change the result of many counting queries.

## Counting query & Histogram query II

Lets take an example consider two persons one has no properties another has all the properties, changing them will cause  $l_1$  sensitivity to differ by  $k$ . So in this case the sensitivity is bounded by

$$\Delta_f = \sum_i |f_i(X) - f_i(X')| \leq k,$$

where  $X, X'$  are neighbouring databases. Thus we will use  $\Delta = k$  and we will add noise  $Y_i \sim \text{Laplace}(k/\epsilon)$  into each co-ordinates.

**Histogram query** is example of a structured query. Where each universe  $\mathcal{X}^n$  is partitioned into bins. So a person, at a time can belong to a single bin. Due to the above fact addition of one individual will change, the count of histogram query at most 1, this type of queries can be answered by adding independent draws from  $\text{Laplace}(1/\epsilon)$  to the original count of each cell.

1. Sending a Query,  $Q$

# Exponential Mechanism

Differential  
Privacy  
Algorithm

Data provider

2. Response the Result,  $O$

3. Adding Noise on Result,  $O'$

# Motivation beyond Exponential Mechanism

The exponential mechanism was designed for cases where we wish the best response, but adding some noise to the best value can dramatically change its value. Consider example of an auction where auction house has adequate supply of a product, and we have 3 bidders they want to bid on the product first two bidders are willing to pay \$1 and third bidder willing to pay \$3.01 for the product. Now consider the following situations,

- Price at the auction house is \$1 in this case revenue will be \$3.
- Price at the auction house is \$1.01 in this case revenue will drop to \$1.01.
- Price at the auction house is \$3.01 in this case revenue will be \$3.01.
- Price at the auction house is \$3.02 in this case revenue will be \$0.
- We can see that dramatic change in revenue in 2 and 4th point. While the change of price is very few. Hence here sensitivity is more.
- So here instead considering prices as numerical value we will consider prices as objects. So in this set up price \$1 and \$3.01 as "High quality object" and 1.01 and 3.02 as "Low quality object"

## Definition: Exponential Mechanism

Consider, a dataset  $X \in \mathcal{X}^n$ , a set of objects  $\mathcal{H}$  & a score function  $s : \mathcal{X}^n \times \mathcal{H} \rightarrow \mathbb{R}$

### Definition

**Sensitivity of the score function** for two neighbouring dataset  $X, X'$  is defined as,

$$\Delta_s = \max_{X, X'} \|s(X, h) - s(X', h)\|,$$

### Definition

The exponential mechanism  $M_E(X, h, s)$  selects and outputs some object  $h \in \mathcal{H}$ , with probability proportional to  $\exp(\frac{\epsilon s(X, h)}{2\Delta})$ .

**Note:** We assume that set of objects and the score function are public and we don't bother about their security or privacy. The only is the dataset  $X$  which is kept secret here.

Thus sensitivity of the score function is defined on the dataset only.

# Privacy proof I

## Theorem

*The exponential mechanism  $M_E$  is  $\epsilon$ -differentially private mechanism.*

## Privacy proof II

Proof.

Let us fix two neighbouring datasets  $X, X'$  and some  $h \in \mathcal{H}$ , Then we have

$$\begin{aligned} \frac{Pr[M_E(X) = h]}{Pr[M_E(X') = h]} &= \frac{\frac{\exp(\frac{\epsilon s(X, h)}{2\Delta})}{\sum_{h'} \exp(\frac{\epsilon s(X, h')}{2\Delta})}}{\frac{\exp(\frac{\epsilon s(X', h)}{2\Delta})}{\sum_{h'} \exp(\frac{\epsilon s(X', h')}{2\Delta})}} = \exp\left(\frac{\epsilon(s(X, h) - s(X', h))}{2\Delta}\right) \frac{\sum_{h'} \exp(\frac{\epsilon s(X', h')}{2\Delta})}{\sum_{h'} \exp(\frac{\epsilon s(X, h')}{2\Delta})} \\ &\leq \exp(\epsilon/2) \exp(\epsilon/2) \frac{\sum_{h'} \exp(\frac{\epsilon s(X, h')}{2\Delta})}{\sum_{h'} \exp(\frac{\epsilon s(X, h')}{2\Delta})} = \exp(\epsilon) \end{aligned}$$

The above inequality holds based on the definition of  $\Delta$ . Second inequality follows from the fact that,

$$\exp\left(\frac{\epsilon s(X', h')}{2\Delta}\right) \leq \exp(\epsilon/2) \exp\left(\frac{\epsilon s(X, h')}{2\Delta}\right)$$

Hence the proof. □



# Connection between Exponential and Laplace Mechanism

Laplace mechanism can be thought as a particular case of exponential mechanism. Let's take we are interested in computing the sensitivity of  $\Delta$  statistic  $f : \mathcal{X}^n \rightarrow \mathbb{R}$  on a dataset  $X$  by Laplace mechanism.

Now let us take set of objects  $\mathcal{H}$  is the real line  $\mathbb{R}$ , and the score function  $s(X, h) = -|f(X) - h|$ .

This yields the probability of a point  $h \in \mathbb{R}$  being output with probability proportional to  $\exp(-\frac{\epsilon|f(X)-h|}{2\Delta})$ . Which is exactly the density of Laplace distribution upto a factor 2, Which can be removed with more care.

# Approximate Differential Privacy I

Here we will discuss relaxation in the definition  $\epsilon$ -differential privacy. This relaxation was first proposed by Dwork, Kenthapadi, McSherry, Mironov, and Naor. The main idea for weakening privacy notion is that to achieve it less amount of noise.

## Definition

**Approximate Differential Privacy** A mechanism  $M : \mathcal{X}^n \rightarrow \mathcal{Y}$  is  $(\epsilon, \delta)$  differential private, if for all neighbouring datasets  $X, X' \in \mathcal{X}^n$  and for all  $T \subseteq \mathcal{Y}$ ,

$$\Pr[M(X) \in T] \leq e^\epsilon \Pr[M(X') \in T] + \delta$$

To interpret new definition we will require another notion, called **privacy loss random variable**.

# Approximate Differential Privacy II

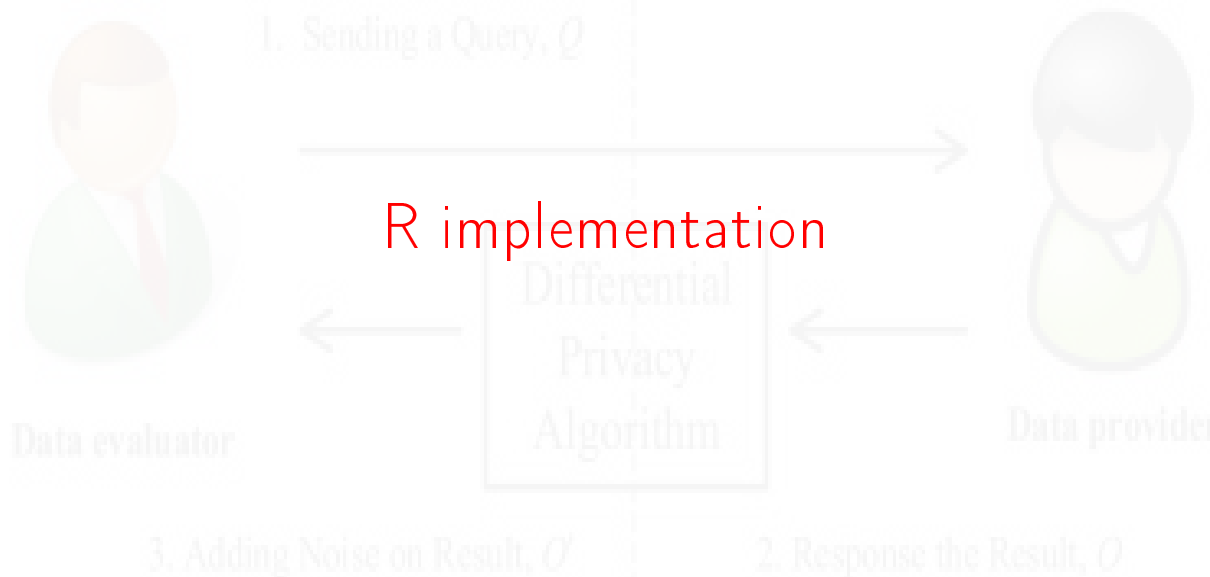
## Definition

Let  $Y$  and  $Z$  be two random variables, then a privacy loss random variable  $\mathcal{L}_{Y||Z}$  is distributed by drawing  $t \sim Y$ , and outputting  $\ln \frac{\Pr[Y=t]}{\Pr[Z=t]}$ .

**Remark 1:** In the above definition supports of  $Y$  and  $Z$  should be equal, otherwise the privacy loss random variable is undefined. In fact this holds for continuous random variables also.

**Remark 2:** We have defined privacy loss in terms of random variables but we will apply it as for  $Y, Z$  equal to  $M(X)$  and  $M(X')$  where  $X \sim X'$  are two neighbouring datasets. Intuition behind notion of privacy loss random variable is that how much it is likely to be input database is  $X$  compared to  $X'$  based on the observation of the realization of  $M(X)$ .

# R implementation



# R implementation of Laplace Mechanism

```
1 # Here we will demonstrate Laplace privatization of
2 # the sample mean on bounded data [0,1]
3 library(diffpriv)
4 func <- function(X) mean(X) ## target function
5 n <- 150 ## dataset size
6 mechanism <- DPMechLaplace(target = func, sensitivity = 1/n, dims = 1)
7 Database <- runif(n, min = 0, max = 1) ## the sensitive database in
   [0,1]^n
8 pparams <- DPParamsEps(epsilon = 1) ## desired privacy budget
9 r <- releaseResponse(mechanism, privacyParams = pparams, X = Database)
10 cat("Private response r$response:", r$response,
11     "\nNon-private response f(D): ", func(Database))
```

Output we get

Private response r\$response: 0.4688349

Non-private response f(D): 0.4733297

# R implementation of Exponential Mechanism I

```
1 # we demonstrate the exponential mechanism here
2
3 library(randomNames) ## a package that generates representative random
  names
4 oracle <- function(n) randomNames(n)
5 Database <- c("Object A", "Object B", "Object C", "Object D",
6             "Object E", "Object F", "Object G",
7             "Object H", "Object I", "Object J")
8 n <- length(Database)
9 func <- function(X) { function(r) sum(r == unlist(base::strsplit(X, ""))
10   ) }
11 rSet <- as.list(letters) ## the response set, letters a—z, must be a
  list
12 mechanism <- DPMechExponential(target = func, responseSet = rSet)
13 #
14 mechanism <- sensitivitySampler(mechanism, oracle = oracle, n = n,
  gamma = 0.1)
```

## R implementation of Exponential Mechanism II

```
14 pparams <- DPParamsEps(epsilon = 1)
15 r <- releaseResponse(mechanism, privacyParams = pparams, X = Database)
16 cat("Private response r$response: ", r$response,
17     "\nNon-private f(D) maximizer: ", letters[which.max(sapply(rSet,
    func(Database)))]])
```

Output we get

```
Private response r$response:  j
Non-private f(D) maximizer:  b>
```

# References

- ① Calibrating Noise to Sensitivity in Private Data Analysis. Cynthia Dwork, Frank McSherry, Kobbi Nissim, Adam Smith.
- ② The Composition Theorem for Differential Privacy, Peter Kairouz, Sewoong Oh, Pramod Viswanath.
- ③ A statistical framework for differential privacy. Larry Wasserman, Shuheng Zhou.
- ④ Revealing Information while Preserving Privacy. Irit Dinur, Kobbi Nissim.
- ⑤ The Algorithmic Foundations of Differential Privacy. Aaron Roth and Cynthia Dwork.
- ⑥ Our Data, Ourselves: Privacy Via Distributed Noise Generation. Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, Moni Naor.