

Overview of Differential Privacy part 2

Jingchen (Monika) Hu

Vassar College

Statistical Data Privacy

Outline

- 1 Introduction
- 2 The Laplace Mechanism
- 3 Properties of differential privacy
- 4 Summary and References

Outline

- 1 Introduction
- 2 The Laplace Mechanism
- 3 Properties of differential privacy
- 4 Summary and References

Recap from overview part 1

- Key idea: add **noise** to the **output** of **statistics** made to **databases**
- Added noise is random; depends on a predetermined **privacy budget** and the type of statistics
- Two important implications:
 - 1 the added noise is positively related to the sensitivity
 - 2 the added noise is negatively related to the privacy budget

Recap from overview part 1

- Key idea: add **noise** to the **output** of **statistics** made to **databases**
- Added noise is random; depends on a predetermined **privacy budget** and the type of statistics
- Two important implications:
 - 1 the added noise is positively related to the sensitivity
 - 2 the added noise is negatively related to the privacy budget
- How to add noise then?

Plan

- The Laplace Mechanism
 - ▶ a mechanism that satisfies ϵ -differential privacy
 - ▶ adds noise from a Laplace distribution to the statistic output
 - ▶ the parameters of the corresponding Laplace distribution depend on the **sensitivity** (Δf) and the **privacy budget** (ϵ)
- Properties of differential privacy
 - ▶ composition theorem
 - ▶ sequential composition
 - ▶ parallel composition
 - ▶ post-processing

Outline

- 1 Introduction
- 2 The Laplace Mechanism**
- 3 Properties of differential privacy
- 4 Summary and References

The Laplace distribution

- A random variable has a Laplace(μ, s) distribution if its probability density function is

$$f(x \mid \mu, s) = \frac{1}{2s} \exp\left(-\frac{|x - \mu|}{s}\right) \quad (1)$$

$$= \frac{1}{2s} \begin{cases} \exp\left(-\frac{\mu-x}{s}\right) & \text{if } x < \mu; \\ \exp\left(-\frac{x-\mu}{s}\right) & \text{if } x \geq \mu, \end{cases} \quad (2)$$

- ▶ μ is a location parameter
- ▶ $s > 0$ is a scale parameter
- ▶ when $\mu = 0, b = 1$, the positive half-line is an exponential distribution scaled by $1/2$

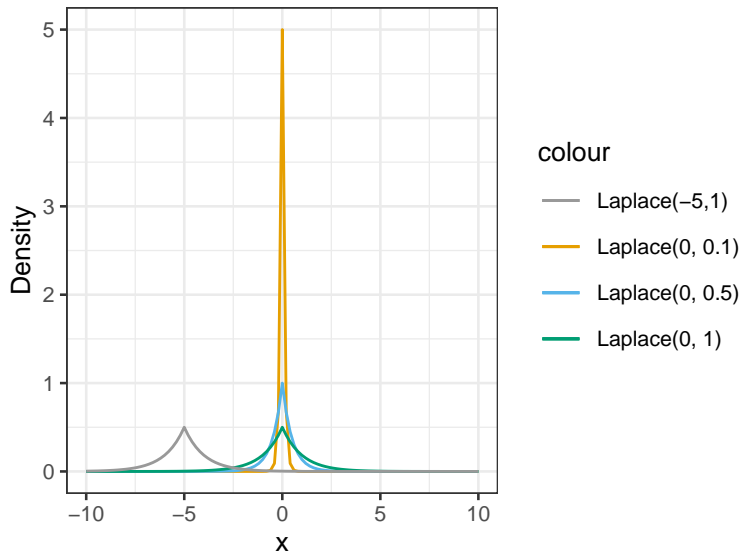
The Laplace distribution

- Like the normal distribution, the Laplace distribution is symmetric
- It is centered at its location parameter μ
- The scale parameter s controls its spread: larger s indicates bigger spread

The Laplace distribution

```
require(rmutil)
cbPalette <- c("#999999", "#E69F00", "#56B4E9", "#009E73",
               "#CC79A7", "#D55E00", "#F0E442", "#0072B2")
ggplot(data.frame(x = c(-10, 10)), aes(x)) +
  stat_function(fun = dlaplace, args = list(m = 0, s = 0.1),
               aes(color = "Laplace(0, 0.1)")) +
  stat_function(fun = dlaplace, args = list(m = 0, s = 0.5),
               aes(color = "Laplace(0, 0.5)")) +
  stat_function(fun = dlaplace, args = list(m = 0, s = 1),
               aes(color = "Laplace(0, 1)")) +
  stat_function(fun = dlaplace, args = list(m = -5, s = 1),
               aes(color = "Laplace(-5,1)")) +
  scale_colour_manual(values = cbPalette) + ylab("Density") +
  theme_bw(base_size = 10, base_family = "")
```

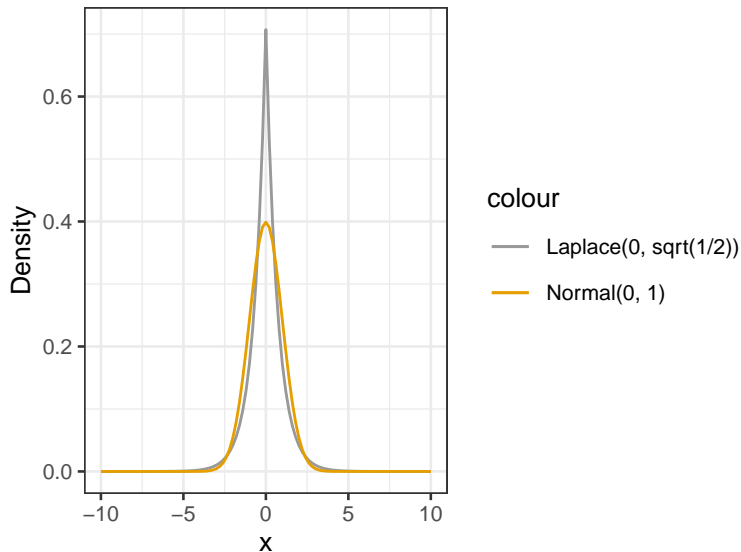
The Laplace distribution



The Laplace distribution vs the normal distribution

```
require(rmutil)
ggplot(data.frame(x = c(-10, 10)), aes(x)) +
  stat_function(fun = dlaplace, args = list(m = 0, s = sqrt(1/2)),
    aes(color = "Laplace(0, sqrt(1/2))")) +
  stat_function(fun = dnorm, args = list(mean = 0, sd = 1),
    aes(color = "Normal(0, 1)")) +
  scale_colour_manual(values = cbPalette) + ylab("Density") +
  theme_bw(base_size = 10, base_family = "")
```

The Laplace distribution vs the normal distribution



Laplace noise for privacy protection

- The Laplace Mechanism adds noise to the output with ϵ -differential privacy guarantee
- The noise is drawn from a Laplace distribution
- Two important implications:
 - 1 the added noise is positively related to the sensitivity (Δf)
 - 2 the added noise is negatively related to the privacy budget (ϵ)
- Also know that:
 - 1 the sensitivity (Δf) is dependent on the database and the statistic
 - 2 the privacy budget (ϵ) is independent of the database and the statistic

Laplace noise for privacy protection

- For given sensitivity Δf and privacy budget ϵ , the added noise to the output of a statistic sent to database \mathbf{x} , X^* is drawn from a Laplace distribution with mean 0, and scale $\frac{\Delta f}{\epsilon}$:

$$X^* \sim \text{Laplace} \left(0, \frac{\Delta f}{\epsilon} \right) \quad (3)$$

Laplace noise for privacy protection

- For given sensitivity Δf and privacy budget ϵ , the added noise to the output of a statistic sent to database \mathbf{x} , X^* is drawn from a Laplace distribution with mean 0, and scale $\frac{\Delta f}{\epsilon}$:

$$X^* \sim \text{Laplace}\left(0, \frac{\Delta f}{\epsilon}\right) \quad (3)$$

- The scale of a Laplace distribution controls its spread, and larger scale value indicates bigger spread
- If the added noise needs to be larger, we should draw it from a Laplace distribution with larger scale value, and vice versa

The scale $\frac{\Delta f}{\epsilon}$

$$X^* \sim \text{Laplace}\left(0, \frac{\Delta f}{\epsilon}\right)$$

- The scale $\frac{\Delta f}{\epsilon}$ is the ratio of the ℓ_1 -sensitivity and the privacy budget
- Connection to the two implications?
 - 1 the added noise is positively related to the sensitivity (Δf)
 - 2 the added noise negatively related to the privacy budget (ϵ)

The Laplace Mechanism

- Formally, given any function $f : \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathbb{R}^k$, the Laplace mechanism is defined as

$$\mathcal{M}_L(\mathbf{x}, f(\cdot), \epsilon) = f(\mathbf{x}) + (X_1^*, \dots, X_k^*), \quad (4)$$

where X_i^* are i.i.d. random variables drawn from $\text{Laplace}\left(0, \frac{\Delta f}{\epsilon}\right)$

Examples: count statistic and average statistic

$$\mathcal{M}_L(\mathbf{x}, f(\cdot), \epsilon) = f(\mathbf{x}) + (X_1^*, \dots, X_k^*),$$

- statistic f : How many rural CUs are there in this sample?
 - ▶ $\Delta f = 1$
 - ▶ $k = 1$ (i.e. statistic f is 1–dimension)
 - ▶ **Discussion question**: What is the Laplace distribution the noise should be drawn from?

Examples: count statistic and average statistic

$$\mathcal{M}_L(\mathbf{x}, f(\cdot), \epsilon) = f(\mathbf{x}) + (X_1^*, \dots, X_k^*),$$

- statistic f : How many rural CUs are there in this sample?
 - ▶ $\Delta f = 1$
 - ▶ $k = 1$ (i.e. statistic f is 1–dimension)
 - ▶ **Discussion question**: What is the Laplace distribution the noise should be drawn from?

- Another statistic f : What is the average income of this sample?
 - ▶ $\Delta f = \frac{b-a}{n}$
 - ▶ $k = 1$ (i.e. statistic f is 1–dimension)
 - ▶ **Discussion question**: What is the Laplace distribution the noise should be drawn from?

The Laplace Mechanism preserves ϵ -differential privacy

Proof: Let $\mathbf{x}, \mathbf{y} \in \mathbb{N}^{|\mathcal{X}|}$ and $\delta(\mathbf{x}, \mathbf{y}) = 1$, and let $f(\cdot)$ be some function $f: \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathbb{R}^k$. Let $p_{\mathbf{x}}$ and $p_{\mathbf{y}}$ denote the probability density functions of $\mathcal{M}_L(\mathbf{x}, f(\cdot), \epsilon)$ and $\mathcal{M}_L(\mathbf{y}, f(\cdot), \epsilon)$. We compare the two at some arbitrary output point $z \in \mathbb{R}^k$:

$$\begin{aligned}
 \frac{p_{\mathbf{x}}(z)}{p_{\mathbf{y}}(z)} &= \prod_{i=1}^k \left(\frac{\exp\left(-\frac{|f(\mathbf{x})_i - z_i|}{\Delta f / \epsilon}\right)}{\exp\left(-\frac{|f(\mathbf{y})_i - z_i|}{\Delta f / \epsilon}\right)} \right) \\
 &= \prod_{i=1}^k \exp\left(\epsilon \frac{|f(\mathbf{y})_i - z_i| - |f(\mathbf{x})_i - z_i|}{\Delta f}\right) \\
 &\leq \prod_{i=1}^k \exp\left(\epsilon \frac{|(f(\mathbf{x})_i - f(\mathbf{y})_i)|}{\Delta f}\right) \\
 &= \exp\left(\epsilon \frac{\|f(\mathbf{x}) - f(\mathbf{y})\|_1}{\Delta f}\right) \leq \exp(\epsilon)
 \end{aligned}$$

CE example of a count statistic: step 1

- Calculate the true count of rural CUs

```
require(readr)
CEdata <- read_csv("CEdata.csv")
CEdata$s[CEdata$UrbanRural == 2] <- 1
```

```
## Warning: Unknown or uninitialised column: `s`.
```

```
CEdata$s[CEdata$UrbanRural == 1] <- 0
```

```
n_rural <- CEdata %>%
  summarize_at(vars(s), sum) %>%
  pull()
n_rural
```

```
## [1] 337
```

CE example of a count statistic: step 2

- Add Laplace noise to the true count

```
Delta_f_count <- 1
```

```
require(rmutil)
set.seed(123)
epsilon1 <- 0.1
rlaplace(1, n_rural, Delta_f_count/epsilon1) %>%
  round()
```

```
## [1] 331
```

CE example of a count statistic: step 2

```
set.seed(123)
epsilon2 <- 1
rlaplace(1, n_rural, Delta_f_count/epsilon2) %>%
  round()
```

```
## [1] 336
```

- With the true count of 337 rural CUs, we can see that smaller privacy budget adds more noise, $337 - 331 = 6$ (when $\epsilon = 0.1$) versus $337 - 336 = 1$ (when $\epsilon = 1$)
- These outcomes are in line with our previously discussed implications, that when fixing the sensitivity value, the added noise is negatively related to the privacy budget

Outline

- 1 Introduction
- 2 The Laplace Mechanism
- 3 Properties of differential privacy**
- 4 Summary and References

Composition theorem

- Idea: taking together the independent use of an ϵ_1 -differentially private algorithm and an ϵ_2 -differentially private algorithm, results in $(\epsilon_1 + \epsilon_2)$ -differential privacy

Composition theorem

- Idea: taking together the independent use of an ϵ_1 -differentially private algorithm and an ϵ_2 -differentially private algorithm, results in $(\epsilon_1 + \epsilon_2)$ -differential privacy
- Formally, let $\mathcal{M}_1 : \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathcal{R}_1$ be an ϵ_1 -differentially private algorithm, and let $\mathcal{M}_2 : \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathcal{R}_2$ be an ϵ_2 -differentially private algorithm. Then their combination, defined to be $\mathcal{M}_{1,2} : \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathcal{R}_1 \times \mathcal{R}_2$ by the mapping: $\mathcal{M}_{1,2}(\mathbf{x}) = (\mathcal{M}_1(\mathbf{x}), \mathcal{M}_2(\mathbf{x}))$ is $(\epsilon_1 + \epsilon_2)$ -differentially private
- Proof omitted; check out Dwork and Roth (2014)

Composition theorem

- Idea: taking together the independent use of an ϵ_1 -differentially private algorithm and an ϵ_2 -differentially private algorithm, results in $(\epsilon_1 + \epsilon_2)$ -differential privacy
- Formally, let $\mathcal{M}_1 : \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathcal{R}_1$ be an ϵ_1 -differentially private algorithm, and let $\mathcal{M}_2 : \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathcal{R}_2$ be an ϵ_2 -differentially private algorithm. Then their combination, defined to be $\mathcal{M}_{1,2} : \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathcal{R}_1 \times \mathcal{R}_2$ by the mapping: $\mathcal{M}_{1,2}(\mathbf{x}) = (\mathcal{M}_1(\mathbf{x}), \mathcal{M}_2(\mathbf{x}))$ is $(\epsilon_1 + \epsilon_2)$ -differentially private
- Proof omitted; check out Dwork and Roth (2014)
- A generalization: let $\mathcal{M}_i : \mathbb{N}^{|\mathcal{X}|} \rightarrow \mathcal{R}_i$ be an ϵ_i -differentially private algorithm for $i \in [k]$. Then if $\mathcal{M}_{[k]} : \mathbb{N}^{|\mathcal{X}|} \rightarrow \prod_{i=1}^k \mathcal{R}_i$ is defined to be $\mathcal{M}_{[k]}(\mathbf{x}) = (\mathcal{M}_1(\mathbf{x}), \dots, \mathcal{M}_k(\mathbf{x}))$, then $\mathcal{M}_{[k]}$ is $(\sum_{i=1}^k \epsilon_i)$ -differentially private

Sequential composition

- If m statistics are sent to the same dataset, the privacy budget needs to be divided by m :

$$\epsilon_{new} = \frac{\epsilon}{m}. \quad (5)$$

Sequential composition

- If m statistics are sent to the same dataset, the privacy budget needs to be divided by m :

$$\epsilon_{new} = \frac{\epsilon}{m}. \quad (5)$$

- Example: adding Laplace noise to four statistics
 - ➊ What is the average income of this sample?
 - ➋ What is the average expenditures of this sample?
 - ➌ What is the variance of income in this sample?
 - ➍ What is the variance of expenditures in this sample?

```
epsilon <- 0.1
m <- 4
epsilon_new <- epsilon/m
rlaplace(1, income_average, Delta_f_average_income/epsilon_new)
rlaplace(1, expenditures_average, Delta_f_average_exp/epsilon_new)
rlaplace(1, income_variance, Delta_f_variance_income/epsilon_new)
rlaplace(1, expenditures_variance, Delta_f_variance_exp/epsilon_new)
```

Sequential composition

```
epsilon <- 0.1  
m <- 4  
epsilon_new <- epsilon/m
```

```
rlaplace(1, income_average, Delta_f_average_income/epsilon)
```

results in smaller noise (larger ϵ , same Δf)

```
rlaplace(1, income_average, Delta_f_average_income/epsilon_new)
```

results in larger noise (smaller ϵ , same Δf)

Parallel composition

- If m statistics are sent to the same database but on non-overlapping subsets of the dataset, the privacy budget does not need to be divided by m

Parallel composition

- If m statistics are sent to the same database but **on non-overlapping subsets** of the dataset, the privacy budget does **not** need to be divided by m
- Examples: adding Laplace noise to two statistics
 - 1 What is the average income of rural CUs?
 - 2 What is the average income of urban CUs?

$$\epsilon_{new} = \epsilon$$

Post-processing

- Dealing with contingency tables (e.g. counts of observations in each category of a categorical variable)
- The post-processing property indicates that for a contingency table with c cells:

$$\epsilon_{new} = \frac{\epsilon}{c - 1}. \quad (6)$$

- This is because knowing the noisy counts of $c - 1$ cells determines the count of the c -th cell

Post-processing cont'd

- Examples: adding Laplace noise to six statistics (i.e. $c = 6$)
 - ① How many CUs with reference person's race category as White are there in this sample?
 - ② How many CUs with reference person's race category as Black are there in this sample?
 - ③ How many CUs with reference person's race category as Native American are there in this sample?
 - ④ How many CUs with reference person's race category as Asian are there in this sample?
 - ⑤ How many CUs with reference person's race category as Pacific Islander are there in this sample?
 - ⑥ How many CUs with reference person's race category as Multi-race are there in this sample?

```
epsilon <- 0.1
c <- 6
epsilon_new <- epsilon / (c-1)
Race_6 <- n - Race_1 - Race_2 - Race_3 - Race_4 - Race_5
```

Outline

- 1 Introduction
- 2 The Laplace Mechanism
- 3 Properties of differential privacy
- 4 Summary and References**

Summary

- The Laplace Mechanism
- Properties of differential privacy

Summary

- The Laplace Mechanism
- Properties of differential privacy
- Next week:
 - ▶ Project presentations, 10-12min each team
 - ▶ Use presentation slides
 - ▶ Include a couple of slides on differential privacy methods for your data

References I

Dwork, C., and A. Roth. 2014. The Algorithmic Foundations of Differential Privacy. Foundations; Trends in Theoretical Computer Science.