

# CSE660

# Differential Privacy

November 1, 2017

**Marco Gaboardi**

Room: 338-B

[gaboardi@buffalo.edu](mailto:gaboardi@buffalo.edu)

<http://www.buffalo.edu/~gaboardi>

# Outline of the class

## ***Week 1***

Introduction, motivation and privacy limitations. Definition of Differential Privacy and the curator model.

## ***Week 2***

Basic mechanisms: Randomized Response, Laplace Mechanism,

## ***Week 3***

Basic properties following from the definition, Exponential Mechanism and comparison with the other basic mechanisms.

## ***Week 4***

The Report Noisy max algorithm.

## ***Week 5***

The Sparse Vector technique. Releasing Many Counting Queries with Correlated Noise. The smallDB algorithm.

## ***Week 6***

The MWEM algorithm.

# Outline of the class

## ***Week 7***

Revisiting MWEM, The DualQuery algorithm.

## ***Week 8***

Advanced Composition and variations on differential privacy: Renyi DP, zero-concentrated DP.

## ***Week 9***

Studying the experimental accuracy.  
The local model for differential privacy.

## ***Week 10***

More algorithms for the local model.

## ***Week 11***

PAC learning and private PAC learning

## ***Week 12***

Differentially Private Hypothesis Testing

## ***Week 13***

Differential Privacy and Generalization in Adaptive Data Analysis

## ***Week 14***

Project presentations

# Differential privacy

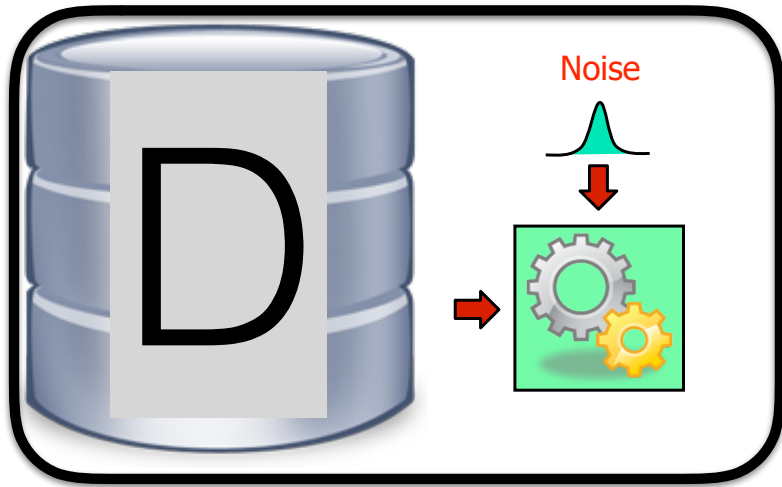
## Definition

Given  $\epsilon, \delta \geq 0$ , a probabilistic query  $Q: X^n \rightarrow R$  is  $(\epsilon, \delta)$ -differentially private iff

for all adjacent database  $b_1, b_2$  and for every  $S \subseteq R$ :

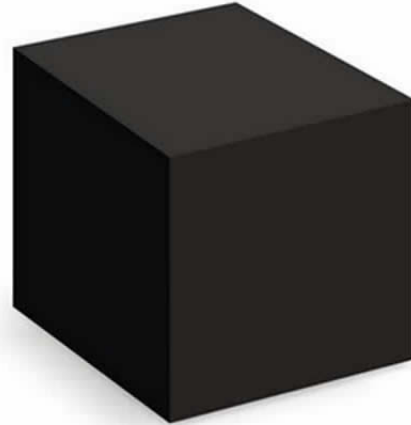
$$\Pr[Q(b_1) \in S] \leq \exp(\epsilon) \Pr[Q(b_2) \in S] + \delta$$

# Composition



The overall process is  $(\epsilon_1 + \epsilon_2 + \dots + \epsilon_n)$ -DP

# Composition



We always need to think before applying composition to whether we have other options!

# Composition

**Theorem 1.18** (Standard composition for  $\epsilon$ -differential privacy). Let  $\mathcal{M}_i : \mathcal{X}^n \rightarrow R_i$  be  $\epsilon_i$ -differentially private algorithms (for  $1 \leq i \leq k$ ). Then, their composition defined to be  $\mathcal{M}(D) = (\mathcal{M}_1(D), \mathcal{M}_2(D), \dots, \mathcal{M}_k(D))$  is  $\sum_{i=1}^k \epsilon_i$ -differentially private.

# Privacy Loss

In general we can think about the following quantity as the **privacy loss** incurred by observing  $r$  as output of  $\mathcal{M}$  on the databases  $D$  and  $D'$ .

$$\mathcal{L}_{\mathcal{M}}^{D \rightarrow D'}(r) = \ln \left( \frac{\Pr[\mathcal{M}(D) = r]}{\Pr[\mathcal{M}(D') = r]} \right) = -\mathcal{L}_{\mathcal{M}}^{D' \rightarrow D}(r)$$

The  $(\epsilon, 0)$ -differential privacy requirement corresponds to requiring that for every  $r$  and every adjacent  $D, D'$  we have:

$$\left| \mathcal{L}_{\mathcal{M}}^{D \rightarrow D'}(r) \right| \leq \epsilon$$



# $(\epsilon, \delta)$ -Differential Privacy

This corresponds to a privacy loss of the form:

$$\mathcal{L}_{\mathcal{M}}^{D \rightarrow D'}(r) = \ln \left( \frac{\Pr[\mathcal{M}(D) = r | E]}{\Pr[\mathcal{M}(D') = r | E']} \right)$$

The  $(\epsilon, \delta)$ -differential privacy requirement corresponds to requiring that for every  $r$  and every adjacent  $D, D'$  we have:

$$\Pr \left[ \left| \mathcal{L}_{\mathcal{M}}^{D \rightarrow D'}(r) \right| \leq \epsilon \right] \geq 1 - \delta$$

# Composition for $(\epsilon, \delta)$ -DP<sup>10</sup>

**Theorem 1.22** (Standard composition for  $(\epsilon, \delta)$ -differential privacy).  
Let  $\mathcal{M}_i : \mathcal{X}^n \rightarrow R_i$  be  $(\epsilon_i, \delta_i)$ -differentially private algorithms (for  $1 \leq i \leq k$ ). Then, their composition defined to be  $\mathcal{M}(D) = (\mathcal{M}_1(D), \mathcal{M}_2(D), \dots, \mathcal{M}_k(D))$  is  $(\sum_{i=1}^k \epsilon_i, \sum_{i=1}^k \delta_i)$ -differentially private.

# Advanced Composition

11

**Question:** how much perturbation do we have if we want to answer  $n$  queries under  $(\epsilon, \delta)$ -DP?

Using advanced composition we have as a max error

$$O\left(\frac{1}{\epsilon_{\text{global}} \sqrt{n}}\right)$$

If we don't renormalize this is of the order of

$$O\left(\frac{\sqrt{n}}{\epsilon_{\text{global}}}\right)$$

comparable to the sample error.

[DworkRothblumVadhan 10, SteinkeUllman 16]

# Advanced Composition

**Theorem 1.23** (Advanced composition). Let  $\mathcal{M}_i : \mathcal{X}^n \rightarrow R_i$  be  $(\epsilon, \delta)$ -differentially private algorithms (for  $1 \leq i \leq k$  and  $k < 1/\epsilon$ ). Then, their composition defined to be  $\mathcal{M}(D) = (\mathcal{M}_1(D), \mathcal{M}_2(D), \dots, \mathcal{M}_k(D))$  is  $(O(\sqrt{2k \ln(1/\delta')})\epsilon, k\delta + \delta')$ -differentially private for every  $\delta' > 0$ .

**Intuition:** some of the outputs have positive privacy loss (i.e. give evidence for dataset  $D$ ) and some have negative privacy loss (i.e. give evidence for dataset  $D'$ ). The cancellations gives a smaller overall privacy loss.

## Strategy:

- 1-considering the expected value of the privacy loss,
- 2-bound the privacy loss of all the variables together
- 3-compute the probability

# Answering multiple queries<sup>13</sup>

We have seen several methods to answer a single query:

- Randomized Response
- Laplace Mechanism
- Exponential Mechanism

And methods to answer multiple queries with small error:

- Standard composition - we can answer  $\sqrt{n}$  queries.
- Advanced composition - we can answer  $n$  queries.

# Gaussian Mechanism

We have seen several methods to answer a single query:

- Randomized Response
- Laplace Mechanism
- Exponential Mechanism

These mechanisms are  $(\epsilon, \delta)$ -differentially private for every  $\epsilon > 0$  and  $\delta \geq 0$ .

We will add today another one that is  $(\epsilon, \delta)$ -differentially private for every  $\epsilon > 0$  and  $\delta > 0$ .

# Another measure of global sensitivity

**Definition 1.13** (Global sensitivity in  $\ell_2$ ). The *global sensitivity in  $\ell_2$*  of a function  $q : \mathcal{X}^n \rightarrow \mathbb{R}$  is:

$$\Delta_2 q = \max \left\{ \sqrt{(q(D) - q(D'))^2} \mid D \sim_1 D' \in \mathcal{X}^n \right\}$$

# Gaussian Mechanism

---

**Algorithm 14** Pseudo-code for the Gaussian Mechanism

---

```
1: function GAUSSMECH( $D, q, \epsilon$ )  
2:    $Y \stackrel{\$}{\leftarrow} \text{Gauss}(0, \frac{2 \ln(\frac{1.25}{\delta})(\Delta_2 q)^2}{\epsilon^2})$   
3:   return  $q(D) + Y$   
4: end function
```

---

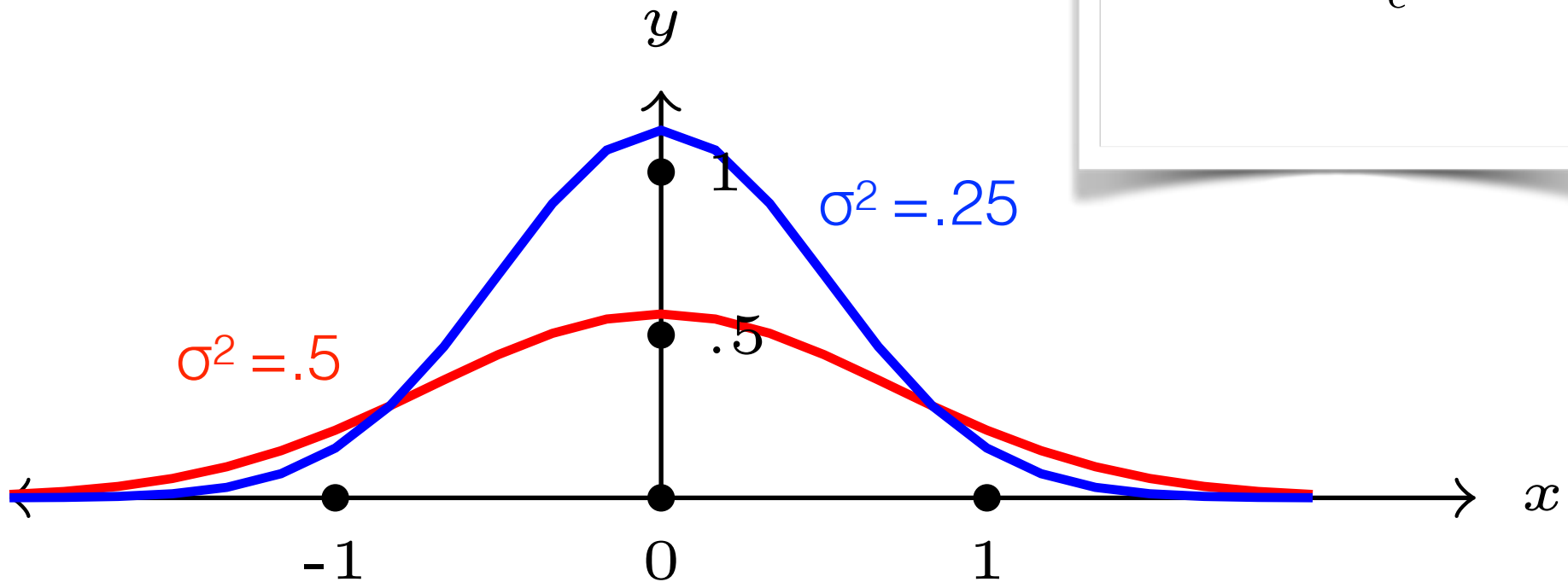


# Gaussian Distribution

$$\text{Gauss}(\mu, \sigma^2)(X) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}((X - \mu)^2)\right)$$

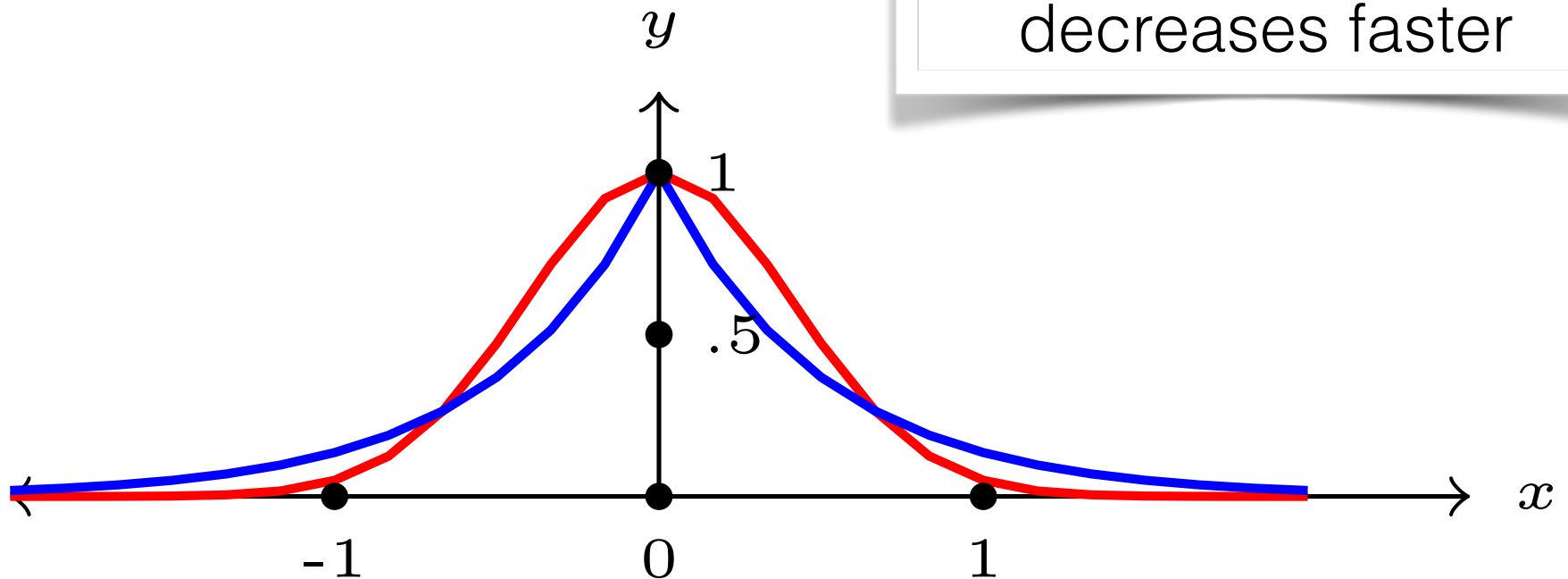
$\sigma^2$  regulates the  
the curve,  
in our case

$$\sigma^2 = \frac{2 \ln\left(\frac{1.25}{\delta}\right) (\Delta_2 q)^2}{\epsilon^2}$$



# Gaussian vs Laplace Distribution

Laplace has a better variance but the tail of the normal distribution decreases faster



# Gaussian Mechanism

## Theorem (Privacy of the Gaussian Mechanism)

The Gaussian mechanism is  $(\epsilon, \delta)$ -differentially private.

**Proof:** Intuitively

We need  $\delta$  to account for bigger differences in the tail

