# Leigett-2013-Tutorial

## Tutorial on Differential Privacy

<u>Talk webpage</u>

Presenter: <u>Katrina Ligett, Hebrew University</u>

## 1. Motivation about Differential Privacy

▼ Question: How to preserve the privacy of data from individuals?

Idea 1: make the answers of the queries randomized output (solves the *differential attacks*)

▼ Question: How to determine the noise level we would like to add to the answers of the queries?

Idea 2: if whether my data appears in the database or not will not affect the answer of the query, then the adversary cannot infer any information from me by looking at the answer of the query.

# 2. Formal Definition

*Database $D$*: the place where all data from individuals are stored, each row stands for data from one individual.

*Mechanism $M$*: the computation process that maps the database to the output space, e.g. the median-computing mechanism would map the database to the real numbers.
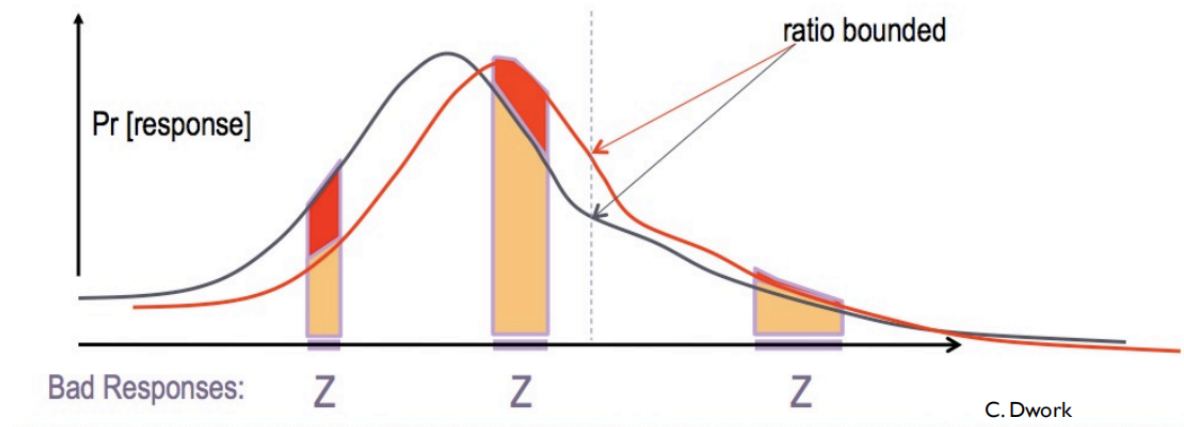
We would characterize the *differential privacy* as a statistical property of the <u>*mechanism $M$*</u>:

- unaffected by the auxiliary information

- independent of the adversary's computation power

## 2.1. $\varepsilon$-differential privacy

> [DinurNissim03, DworkNissimMcSherrySmith06] For any two neighboring database $D$ and $D'$, and any possible output $O \in \mathrm{Range}(M)$, we have
> $$\Pr\left[M\left(D\right) \in O\right] \leq e^\varepsilon \Pr\left[M\left(D'\right) \in O\right]$$

The differential privacy-preserving mechanism is robust to the neighborhoods of the database. Mathematically, we want the the distribution of answer $M(D)$ and $M(D')$ to be close to each other. This can be done by requiring the KL-divergence between $M(D)$ and $M(D')$ small.

$$\mathrm{KL}\left(M(D)\|M(D')\right) = \mathbb{E}\left[\ln\frac{\Pr[M(D) \in O]}{\Pr[M(D') \in O]}\right] \leq \varepsilon$$

But the KL-divergence only add restrictions on the expectation. Hence we replace the expectation with the total variation

$$\max_{O \in \mathrm{Range}(M)}\left[\ln\frac{\Pr[M(D) \in O]}{\Pr[M(D') \in O]}\right] \leq \varepsilon$$
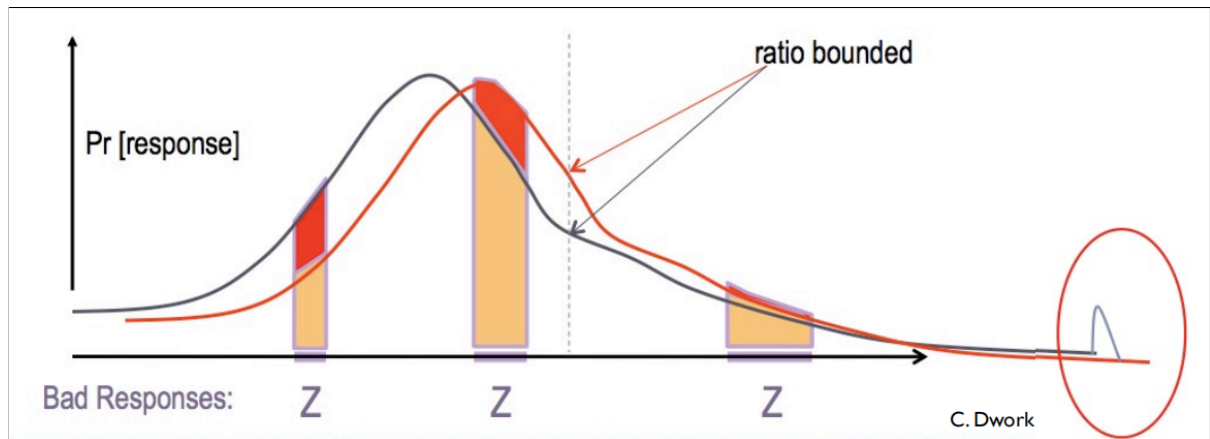
which yields the definition above.

## 2.2. $(\varepsilon, \delta)$-differential privacy

> 📖 [$(\varepsilon, \delta)$-differential privacy] For any two neighboring database $D$ and $D'$, and any possible output $O \in \mathrm{Range}(M)$, we have
> $$\Pr\left[M(D) \in O\right] \leq e^{\varepsilon}\Pr\left[M(D') \in O\right] + \delta$$

When the mechanism preserves the differential privacy, we can promise to the individuals that: if you leave the database, no outcome will change probability very much.

# 3. Composition

📖 an $\varepsilon_1$-DP mechanism, followed by an $\varepsilon_2$-DP mechanism, is $\varepsilon_1 + \varepsilon_2$-DP

📖 an $(\varepsilon_1, \delta_1)$-DP mechanism, followed by an $(\varepsilon_2, \delta_2)$-DP mechanism, is $(\varepsilon_1 + \varepsilon_2, \delta_1 + \delta_2)$-DP

📖 $k$ runs of $(\varepsilon, \delta)$-DP mechanism gives $(\varepsilon', k\delta + \delta')$-DP where $\varepsilon' = (2k \log (1/\delta'))^{1/2} \varepsilon + k\varepsilon(e^\varepsilon + 1)$

- vector-valued query of dimension $d$

📖 can apply composition and add noise $\mathrm{Laplace}(d\Delta f/\varepsilon)$ in each component of output, where $\Delta f$ is the sensitivity of each component.
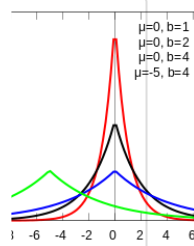
# 4. Noise Addition Mechanism

These mechanism just add random noise to the true answer of the query to achieve DP.

## 4.1. Laplace mechanism

Laplace distribution - Wikipedia

In probability theory and statistics, the Laplace distribution is a continuous probability distribution

W  https://en.wikipedia.org/wiki/Laplace_distribution

When your query is to compute some number $f(D)$, then adding a scaled symmetric Laplace noise is enough.

📖 [DMNS06]: on query $f$, can add scaled symmetric noise $\mathrm{Laplace}(b)$ with $b = \Delta f / \varepsilon$ where $\Delta f = \max_{||D_1 - D_2||_0 \leq 1} ||f(D_1) - f(D_2)||$ , to achieve $\varepsilon$-differential privacy.

## 4.2. Gaussian mechanism

📖 [DKMMN06]: Gaussian noise gives $(\varepsilon, \delta)$-DP with $\sigma \geq (2\log(2/\delta))^{1/2} / \varepsilon (\max L_2 \text{ distance})$

## 4.3. How much queries can we handle with noise addition mechanism?



a brief history of synthetic data (theory)

BLR08: ε-DP, error $\log^{1/3} |Q| \, n^{2/3}$

DNRRV09: (ε, δ)-DP, error $|Q|^{o(1)} n^{1/2}$

DRV10: (ε, δ)-DP, error polylog $|Q| \, n^{1/2}$

HR10: (ε, δ)-DP, error $\log|Q| \, n^{1/2}$

HLM12: simple & matches best bounds

# 4. More Sophisticated Mechanism

### 4.1. Exponential mechanism [MT07]

> 1. specify a scoring function $u(D, O)$
> 2. select the output $O \in \mathrm{Range}(M)$ with probability $\sim$
> $\exp\{\varepsilon u(D, C)/(2\Delta u)\}$

### 4.2. Interactive mechanism

- so far, have discussed creating synthetic data where must know query set in advance

- tools exist to answer similar number of queries on the fly (correlating randomness across queries)

# 5. Robustness

### 5.1. Bootstrapping

TBC

### 5.2. Propose-test-release

TBC