# An introduction to differential privacy

Édouard Pauwels and Sébastien Gerchinovitz

Draft of May 4, 2016

**Abstract**

Tentative notes for the presentation. This note is based on the survey of Dwork [2], the monograph of Dwork and Aaron [3] and the series of talks given by Rothblum at CIRM in January 2016.

## Contents

# 1 Introduction

## 1.1 Content of this note

Differential privacy was introduced in 2006 by Dwork *et al.* [1]. This is a data/output perturbation approach which relies on randomization. The general intuition behind this mechanism is that information released should be distributed similarly whether a specific

individual takes part in a statistical study or not. We will introduce this notion in light of the elements of context given below. We will only consider privacy and will mention utility only briefly.

## 1.2   Context

**What is privacy (vie privée, intimité)?**   For us today: "control over others' use and access to information about me". A major factor of increasing privacy concerns is the progress of technology: photography, press, sound and video recording, computer networks, portable connected devices, biotechnologies ... Not semantic security: access to data does not allow to learn something that cannot be learnt without access (Dalenius, smocking causes cancer).

**De-anonymization, linkage attack, reconstruction attack:**

- In the years 2000, Sweeney, showed that ZIP, birth date and sex were enough to uniquely identify 87% of the american population. She purshased the voter registration list for Cambridge Massachusetts for 20$, linked it to a patient database (free for academics, charged for industry) which includes records such as ethnicity, visit date, diagnosis, procedure, medication and total charge. She could identify the medical record of William Weld, the governor of Massachusetts at that time[10]. This is an example of linkage attack.

- In 2006 Narayanan and Shmatikov showed that only a few movie ratings were enough to uniquely identify anonymized users from Netflix prize dataset by linking with IMDB publicly available database [9]. As a consequence, they gained information about movies that IMDB users saw and did not want to tell the world about. This was a subject of lawsuits and one important factor in the cancelation of a second challenge. This is another example of linkage attack.

- In 2008 Homer *et al.* showed that it is possible to identify the contribution of an individual to a mixture of different persons' DNA from aggregate counts related to specific genome regions (SNPs). This is possible even when the individual's contribution is as small as 0.1% [5]. Genome wide association studies are based on aggregate counts of those SNPs, for groups of case and control patients. US National Institutes of Health used to fund these types of studies on the condition that aggregate counts are released publicly. This policy was modified in 2008 in response to the study. This is an example of reconstruction attack from aggregate statistics.

- In 2011 Zang and Bolot showed that it is possible to identify individuals from anonymized raw location data such as those collected by smartphones. In 2013, de Montjoye *et al.* showed that human mobility data have identifiability features of the same flavor as human fingerprints [11].

- In 2013 Gymrek *et al.* showed that it is possible to infer surname from a few motifs of the Y chromosome by linking with recreational genetic genealogy databases. Triangulating with metadata such as age and state allows to identify individuals [4].

**Consequences:**

- Anonymization is not sufficient. "Just" releasing statistics is providing to much information.
- It is very hard to decide wether some information is sensitive or not for a particular individual.
  - It depends on what an adversary knows.
  - Refusing to answer a given query is already providing a lot of information.
  - It is provably computationally hard [6].

**What is specific to differential privacy?**

- Rigorous privacy guarantees in a model of information communication.
- Simple composition / multiple queries rules.
- Robust to what an adversary knows or doesn't know (linkage attack).
- Maintain some utility, it should still be possible to infer global trends.

# 2 Differential privacy

## 2.1 Communication model

A trusted curator gathers sensitive information from a large number of individuals in a database $D$. He is the only one in the world who has access to $D$. His goal is to release statistics about $D$ to a data analyst while preserving privacy. The curator receives a finite number of statistical queries to which he communicates an answer by executing a randomized algorithm. The analyst is not trusted, a way to model this is to consider that, once released, the answers are in the public domain and anybody has access to it. An important aspect in this model is that after this process, the only publicly available information about $D$ is limited to what was released by curator. This quote from [1] illustrates this fact,

*since the data will never be used again, the curator can destroy the data (and himself) once the statistics have been published.*

## 2.2 Notations

We consider a universe $\mathcal{U}$ of database rows[1]. A database is an element of $\mathcal{U}^n$ for some $n \in \mathbb{N}$. The curator will answer a given query by applying a randomized algorithm, $A$, which argument is any database and which range is a probability space $\mathcal{Y}$ (*e.g.* for counting querries, $\mathcal{Y} = \{1, \ldots, n\}$). That is for any database $D$, $A(D)$ is a random variable over $\mathcal{Y}$. To fix the ideas, one can think of a deterministic algorithm with two input arguments, one of which is a random variable. Two databases $D_1, D_2 \in \mathcal{U}^n$ are called adjacent if they differ in a single row[2].

---

[1]This is the usual terminology, used without further specification

[2]There are variations on this.

## 2.3 Defining differential privacy

Differential privacy is a property of the distribution of the answers of the curator to statistical queries. Since the curator answers statistical queries by calling randomized algorithms (seen as a function), it is therefore expressed as a property of the corresponding random mapping.

**Definition 1 (Differential privacy)** *A random mapping* $A\colon \mathcal{U}^n \to \mathcal{Y}$ *is said to be* $(\epsilon, \delta)$*-differentially private if for any measurable subset* $Y \subset \mathcal{Y}$ *and any adjacent databases* $D_1, D_2$*, it holds true that*

$$Pr\left[A(D_1) \in Y\right] \leq \exp(\epsilon) Pr\left[A(D_2) \in Y\right] + \delta, \tag{1}$$

*where the probability is taken over the randomness of* $A$*.*

In this note, we will only consider $(\epsilon, 0)$-differential privacy which we also call $\epsilon$-differential privacy. This is less useful in practice, but it conveys much of the ideas and it is much easier to deal with. In this case, differential privacy amounts to control likelihood ratios.

**Example 1 (One boolean)** *Suppose that* $\mathcal{U} = \{0, 1\}$ *and consider a single row database. There are only two possibilities, either* $D = 0$ *or* $D = 1$*. Consider the following randomized algorithm* $A$*.*

- *Flip a coin.*
- *If heads, answer trustfully (the only entry of* $D$*).*
- *If tails, flip a second coin and responds* $1$ *if heads and* $0$ *otherwise.*

*We have* $\log(3)$ *differential privacy.*

$$\frac{Pr\left[A(1) = 1\right]}{Pr\left[A(0) = 1\right]} = \frac{Pr\left[A(0) = 0\right]}{Pr\left[A(1) = 0\right]} = \frac{3/4}{1/4} = 3.$$

**Example 2 (No differential privacy)** *Suppose that* $\mathcal{U} = \{0, 1\}$*,* $n \geq 2$ *and consider the algorithm* $A$ *that releases one database row at random. Let* $D_1$ *be a database full of* $0$ *and* $D_2$ *be a database with a single row equal to* $1$*, the rest being* $0$*. We have* $Pr\left[A(D_1) = 1\right] = 0$ *and* $Pr\left[A(D_2) = 1\right] = 1/n$ *so releasing random entries is not differentially private.*

## 2.4 Privacy guarantees

In this section, for simplicity (and because this is often the case in the litterature) we will consider finite universe $\mathcal{U}$ and algorithmic range $\mathcal{Y}$.

**Closure under postprocessing**  Since the released statistics, are considered to be in the public domain, the privacy guarantees should be robust to data post proccessing. This is illustrated by the following result which is similar to *data-processing inequalities* (if $X \to Y \to Z$ form a Markov chain, then, $I(X; Y) \geq I(X; Z)$ where $I$ is the mutual information: the KL divergence between joint distribution and product of marginals).

**Proposition 1 (Closure under postprocessing)** *Let $A$ be an $\epsilon$-differentially private random mapping which range is $\mathcal{Y}$. Let $M \colon \mathcal{Y} \to \mathcal{Z}$ be a second random mapping. Assume that they are independent[3]. Then $M \circ A$ is $\epsilon$-differentially private.*

**Proof :** For any $z \in \mathcal{Z}$ and adjacent databases $D_1, D_2 \in \mathcal{U}^n$, we have

$$
\begin{aligned}
Pr\left[M(A(D_1)) = z\right] &= \sum_{y \in \mathcal{Y}} Pr\left[M(A(D_1)) = z | A(D_1) = y\right] Pr\left[A(D_1) = y\right] \\
&= \sum_{y \in \mathcal{Y}} Pr\left[M(y) = z\right] Pr\left[A(D_1) = y\right] \\
&\leq \exp(\epsilon) \sum_{y \in \mathcal{Y}} Pr\left[M(y) = z\right] Pr\left[A(D_2) = y\right] \\
&= \exp(\epsilon) \sum_{y \in \mathcal{Y}} Pr\left[M(A(D_2)) = z | A(D_2) = y\right] Pr\left[A(D_2) = y\right] \\
&= \exp(\epsilon) Pr\left[M(A(D_2)) = z\right].
\end{aligned}
$$

$\square$

**Semantic privacy and robustness to linkage attacks**  Intuitively, since the output of a differentially private procedure does not depend too much on the value of a single row of the database, it should not allow an adversary to learn too much about a specific individual. This is illustrated in the following negative result.

**Proposition 2 (No test can identify a specific individual)** *Let $A : \mathcal{U}^n \to \mathcal{Y}$ be an $\epsilon$-differentially-private random mapping. Let $D_0$ and $D_1$ be two adjacent databases, e.g., that only differ in the first row. Suppose that after observing $A(D)$ we want to know whether $D = D_0$ or $D = D_1$, i.e., we want to solve the hypothesis testing problem*

$$
\begin{cases}
H_0 : & D = D_0 \\
H_1 : & D = D_1
\end{cases}
$$

*Then, the test affinity is lower bounded by $1 - \sqrt{\epsilon/2}$, i.e., for every test $\phi : \mathcal{Y} \to \{0,1\}$, the sum of the type I and type II risks is lower bounded as*

$$
Pr\left[\phi(A(D_0)) = 1\right] + Pr\left[\phi(A(D_1)) = 0\right] \geq 1 - \sqrt{\epsilon/2}.
$$

**Proof :** Set $R = \{\phi = 1\} \subset \mathcal{Y}$, and denote the distributions of $A(D_0)$ and $A(D_1)$ by $P_0$ and $P_1$ respectively. Then we have

$$
\begin{aligned}
Pr\left[\phi(A(D_0)) = 1\right] + Pr\left[\phi(A(D_1)) = 0\right] &= P_0(R) + P_1(R^c) = 1 - \left(P_0(R^c) - P_1(R^c)\right) \\
&\geq 1 - \mathrm{d}_{\mathrm{TV}}(P_0, P_1) \\
&\geq 1 - \sqrt{\mathrm{KL}(P_0, P_1)/2},
\end{aligned}
$$

---

[3]This should be properly defined and never appears in the proof arguments. We need that their "sources" of randomness are independent. In other words, for any $y \in \mathcal{Y}$, and $D \in \mathcal{U}^n$, $M(y)$ and $M(A(D))|A(D) = y$ have the same law.

where the first inequality above follows from the definition of the total variation distance $d_{TV}(P_0, P_1) = \sup_{B \subset \mathcal{Y}} |P_0(B) - P_1(B)|$, and where the second inequality follows from Pinsker's inequality. Now, note that by definition of $\epsilon$–differential privacy, the Kullback-Leibler divergence between $P_0$ and $P_1$ can be upper bounded as follows:

$$\mathrm{KL}(P_0, P_1) = \sum_{y \in \mathcal{Y}} P_0(y) \log \frac{P_0(y)}{P_1(y)} \leq \epsilon \ ,$$

which concludes the proof. $\qquad\square$

## 2.5 Composition of differentially private mappings

An important feature of differential privacy is the possibility to have simple and flexible composition laws. This is very important in practice because it means that one can "program differential privacy" by combining elementary differentially private operations in an arbitrary way, only controlling the privacy "level" $\epsilon$. This also offers the possibility to implement differential privacy in a sequential and adversarial context (*e.g.* counting querries). Again we suppose that the ouptut space is finite.

**Proposition 3 (Parallel composition)** *Let $A_1 \colon \mathcal{U}^n \to \mathcal{Y}_1$ (resp $A_2 \colon \mathcal{U}^n \to \mathcal{Y}_2$) be an $\epsilon_1$ (resp $\epsilon_2$) differentially private random mapping with finite range. Assume that for any $D$, $A_1(D)$ and $A_2(D)$ are independent[4]. For any database $D \in \mathcal{U}^n$, let $A_{1,2}(D) = (A_1(D), A_2(D))$. Then $A_{1,2}$ is $\epsilon_1 + \epsilon_2$ differentially private.*

**Proof :** For any $y_1 \in \mathcal{Y}_1$ and $y_2 \in \mathcal{Y}_2$ and adjacent databases $D, D' \in \mathcal{U}^n$, we have

$$\begin{aligned}
Pr\left[A_{1,2}(D) = (y_1, y_2)\right] &= Pr\left[A_1(D) = y_1\right] Pr\left[A_2(D) = y_2\right] \\
&\leq \exp(\epsilon_1) Pr\left[A_1(D') = y_1\right] \exp(\epsilon_2) Pr\left[A_2(D') = y_2\right] \\
&= \exp(\epsilon_1 + \epsilon_2) Pr\left[A_{1,2}(D') = (y_1, y_2)\right]
\end{aligned}$$

$\qquad\square$

**Proposition 4 (Sequential composition)** *Let $A_1 \colon \mathcal{U}^n \to \mathcal{Y}$ be an $\epsilon_1$ differentially private random mapping with finite range. Let $\mathcal{U}_2 = \mathcal{U}^n \times \mathcal{Y}_1$. Suppose that $A_2 \colon \mathcal{U}_2 \to \mathcal{Y}_2$ is a random mapping such that for any $y \in \mathcal{Y}_1$, $A_2(\cdot, y)$ is $\epsilon_2$ differentially private. Assume that $A_1$ and $A_2$ are independent[5]. Then, the random mapping $D \to (A_1(D), A_2(D, A_1(D)))$ is $\epsilon_1 + \epsilon_2$ differentially private.*

**Proof :** We only give details for the discrete case, they can be generalized to the continuous case. For any $(y_1, y_2) \in \mathcal{Y}_1 \times \mathcal{Y}_2$ and adjacent databases $D, D' \in \mathcal{U}^n$, we have

$$\begin{aligned}
&Pr\left[(A_1(D), A_2(D, A_1(D))) = (y_1, y_2)\right] \\
&= Pr\left[A_1(D) = y_1 \ \& \ A_2(D, A_1(D)) = y_2 | A_1(D) = y_1\right] Pr\left[A_1(D) = y_1\right] \\
&= Pr\left[A_2(D, y_1) = y_2 | A_1(D) = y_1\right] Pr\left[A_1(D) = y_1\right] \\
&= Pr\left[A_2(D, y_1) = y_2\right] Pr\left[A_1(D) = y_1\right] \\
&\leq \exp(\epsilon_1 + \epsilon_2) Pr\left[A_2(D', y_1) = y_2\right] Pr\left[A_1(D') = y_1\right] \\
&= \exp(\epsilon_1 + \epsilon_2) Pr\left[(A_1(D'), A_2(D', A_1(D'))) = (y_1, y_2)\right],
\end{aligned}$$

---

[4] Again, this important condition barely appears in proof arguments.

[5] For any $y \in \mathcal{Y}_1$, and $D \in \mathcal{U}^n$, $A_2(D, y)$ and $A_2(D, A_1(D))|A_1(D) = y$ have the same law.

where the last step is the same as the first four steps with $D'$ in place of $D$. □

# 3 Two important differentially private protocols

We present in this section two mappings which have differential privacy guarantees. They are important because they underpin many more advanced differentially private protocols. The first one involves a continuous output space, it corresponds to noise addition with Laplace distribution. It was introduced in [1] with the notion of differential privacy. The second one involves a utility function and is useful when addition of noise does not make sense (*e.g.* discrete spaces) or destroys utility (*e.g.* best auction).

## 3.1 Laplace mechanism

Laplace distribution with parameter $b$ is denoted by $Lap_b$ and corresponds to the probability distribution over $\mathbb{R}$ which density is $x \to \frac{1}{2b} \exp(-|x|/b)$.

**Definition 2 ($L_1$-sensitivity)** *Let $f \colon \mathcal{U}^n \to \mathbb{R}^k$, for some $k \in \mathbb{N}^*$. We take the following supremum over adjacent databases $D$ and $D'$.*

$$\Delta_f = \sup_{D,D'} \|f(D) - f(D')\|_1$$

The following proposition is one of the main results of [1] which purpose was "noise calibration".

**Proposition 5 (Laplace mechanism)** *Let $f \colon \mathcal{U}^n \to \mathbb{R}^k$, for some $k \in \mathbb{N}^*$. for some $k \in \mathbb{N}^*$ such that $\Delta_f$ is finite. For any $\epsilon > 0$, the random mapping $A_\epsilon \colon D \to f(D) + (Y_1, \ldots, Y_k)^T$ where $\{Y_i\}_{i=1}^k$ are iid $Lap_{\Delta_f/\epsilon}$, is $\epsilon$-differentially private.*

**Proof :** Fix $\epsilon > 0$. The subscript $i$ denotes the coordinate $i$. For any borelian subset $S \subset \mathbb{R}^k$ and adjacent databases $D, D' \in \mathcal{U}^n$, we have

$$
\begin{aligned}
Pr\left[A_\epsilon(D) \in S\right] &= \int_{z \in S} \left(\frac{\epsilon}{2\Delta_f}\right)^k \exp\left(-\frac{\epsilon}{\Delta_f}\|f(D) - z\|_1\right) dz \\
&= \int_{z \in S} \left(\frac{\epsilon}{2\Delta_f}\right)^k \exp\left(\frac{\epsilon}{\Delta_f}\left(-\|f(D') - z\|_1 - \|f(D) - z\|_1 + \|f(D') - z\|_1\right)\right) dz \\
&\leq \int_{z \in S} \left(\frac{\epsilon}{2\Delta_f}\right)^k \exp\left(-\frac{\epsilon}{\Delta_f}\|f(D') - z\|_1 + \frac{\epsilon}{\Delta_f}\|f(D) - f(D')\|_1\right) dz \\
&= \exp\left(\epsilon \frac{\|f(D) - f(D')\|_1}{\Delta_f}\right) \int_{z \in S} \left(\frac{\epsilon}{2\Delta_f}\right)^k \exp\left(-\frac{\epsilon}{\Delta_f}\|f(D') - z\|_1\right) dz \\
&\leq \exp(\epsilon) \int_{z \in S} \left(\frac{\epsilon}{2\Delta_f}\right)^k \exp\left(-\frac{\epsilon}{\Delta_f}\|f(D') - z\|\right) dz \\
&= \exp(\epsilon) Pr\left[A_\epsilon(D') \in S\right]
\end{aligned}
$$

□

We easily see from the proof that the only key properties are symmetry and triangular inequality and this result easily generalizes to other norms and the corresponding distributions.

The following result shows that Laplace mechanism preserves some utility.

**Proposition 6** $f \colon \mathcal{U}^n \to \mathbb{R}^k$, *for some* $k \in \mathbb{N}^*$, *such that* $\Delta_f$ *is finite. For any* $\epsilon > 0$, *the random mapping* $A_\epsilon \colon D \to f(D) + (Y_1, \ldots, Y_k)^T$ *where* $\{Y_i\}_{i=1}^k$ *are iid* $Lap_{\Delta_f/\epsilon}$, *satisfies, for any* $\delta > 0$ *and any database* $D \in \mathcal{U}^n$

$$Pr\left[\|f(D) - A(D)\|_\infty \geq \ln\left(\frac{k}{\delta}\right)\frac{\Delta_f}{\epsilon}\right] \leq \delta.$$

**Proof :** First note that if $Y$ is $Lap_b$, then $Pr[|Y| \geq bt] = \exp(-t)$ for all $t \geq 0$. We apply a union bound:

$$Pr\left[\|f(D) - A(D)\|_\infty \geq \ln\left(\frac{k}{\delta}\right)\frac{\Delta_f}{\epsilon}\right] = Pr\left[\max_{i=1,\ldots,k} |f_i(D) - A_i(D)| \geq \ln\left(\frac{k}{\delta}\right)\frac{\Delta_f}{\epsilon}\right]$$

$$\leq \sum_{i=1}^k Pr\left[|f_i(D) - A_i(D)| \geq \ln\left(\frac{k}{\delta}\right)\frac{\Delta_f}{\epsilon}\right]$$

$$= k\frac{\delta}{k} = \delta$$

$\square$

**Example 3 (First names)** *Consider releasing the counts of occurrences of first names (among 10000 possibilities) of a population of 300,000,000 people. This query has a sensitivity of 2. Laplace mechanism with* $\epsilon = 2$ *is 2-differentially private and ensures that, with probability 95%, no count will be off by more than an additive error of* $\ln(10000/0.05) \leq 12.21$. ($\epsilon = 1$ *yields an additive error at most of 24.42.)*

**Example 4 (Noisy max)** *Consider the problem of releasing the index which achieves the maximum of a set of counting queries. As shown in [1, Claim 3.9], the protocol which consists in adding iid* $Lap_{1/\epsilon}$ *noise to each count and reveal the index of the noisy maximum is* $\epsilon$-*differentially private. Note that this is not a direct application of Laplace Mechanism and requires additional technicalities to be shown.*

## 3.2 Exponential mechanism

This method was designed for situations where adding noise does not make sense (non numeric range) or destroys the structure of the problem (see next example). It was presented in [7].

**Example 5 (Auction pricing)** *An auctionneer has unlimited supply of some good and several bidders propose a price. A selling price is chosen and all the bidders who proposed higher price get one good at the selling price. Suppose that* $a, b$ *and* $c$ *bid respectively* $1, 1$ *and* $3$. *The auctionneer wants to find the optimal price. At* $1$ *and* $3$, *the revenue of the auctionneer is* $3$ *which is optimal. At* $1.01$, *it is* $1.01$ *and at* $3.01$ *it is* $0$. *Adding noise to the bids would destroy the utility as a tiny amount completely changes the revenue.*

We fix a discrete range $\mathcal{Y}$ and assume that we have a utility function which associates to a database $D$ and an element $y \in \mathcal{Y}$ a value $u(D, y) \in \mathbb{R}$. Given a database $D$, the objective is to release a value of $y \in \mathcal{Y}$ such that $u(D, y)$ is close to $\max_{y \in \mathcal{Y}} u(D, y)$. We need a notion of sensitivity.

**Definition 3 (Utility function sensitivity)** *Let $u$ be a utility function. We take the following supremum over adjacent databases $D$ and $D'$.*

$$\Delta_u = \sup_{y \in \mathcal{Y}} \sup_{D, D'} |u(D, y) - u(D', y)|.$$

**Proposition 7 (Exponential mechanism)** *Consider a finite range $\mathcal{Y}$ and a utility function $u$. For any $\epsilon > 0$, the random mapping $A_\epsilon$ which, given a database $D$ outputs $y \in \mathcal{Y}$ with probability proportional to $\exp\left(\epsilon \frac{u(D,y)}{2\Delta_u}\right)$ is $\epsilon$-differentially private.*

**Proof :** Fix $\epsilon > 0$. For any subset $y_0 \in \mathcal{Y}$ and adjacent databases $D, D' \in \mathcal{U}^n$, we have

$$
\begin{aligned}
Pr\left[A(D) = y_0\right] &= \frac{\exp\left(\epsilon \frac{u(D,y_0)}{2\Delta_u}\right)}{\sum_{y \in \mathcal{Y}} \exp\left(\epsilon \frac{u(D,y)}{2\Delta_u}\right)} \\
&= \frac{\exp\left(\epsilon \frac{u(D',y_0)+u(D,y_0)-u(D',y_0)}{2\Delta_u}\right)}{\sum_{y \in \mathcal{Y}} \exp\left(\epsilon \frac{u(D',y)+u(D,y)-u(D',y)}{2\Delta_u}\right)} \\
&\leq \frac{\exp\left(\frac{\epsilon}{2}\right) \exp\left(\epsilon \frac{u(D',y_0)}{2\Delta_u}\right)}{\exp\left(\frac{-\epsilon}{2}\right) \sum_{y \in \mathcal{Y}} \exp\left(\epsilon \frac{u(D',y)}{2\Delta_u}\right)} = \exp(\epsilon) Pr\left[A(D') = y_0\right]
\end{aligned}
$$

$\square$

**Example 6 (Auction pricing continued)** *The exponential mechanism preserves some utility. With the notation of the previous proposition, for any $c \in \mathbb{R}$, we have*

$$
\begin{aligned}
Pr\left[u(D, A(D)) \leq c\right] &= \frac{\sum_{u(D,y) \leq c} \exp\left(\epsilon \frac{u(D,y)}{2\Delta_u}\right)}{\sum_{y \in \mathcal{Y}} \exp\left(\epsilon \frac{u(D,y)}{2\Delta_u}\right)} \leq \frac{|\{y \in \mathcal{Y}, u(D,y) \leq c\}| \exp\left(\epsilon \frac{c}{2\Delta_u}\right)}{\exp\left(\epsilon \frac{\max_{y \in \mathcal{Y}} u(D,y)}{2\Delta_u}\right)} \\
&= |\{y \in \mathcal{Y}, u(D,y) \leq c\}| \exp\left(\epsilon \frac{c - \max_{y \in \mathcal{Y}} u(D,y)}{2\Delta_u}\right).
\end{aligned}
$$

*Choosing $c = \max_{y \in \mathcal{Y}} u(D, y) - \frac{2\Delta_u}{\epsilon}(t + \log(|\mathcal{Y}|))$, we get,*

$$Pr\left[u(D, A(D)) \leq \max_{y \in \mathcal{Y}} u(D, y) - \frac{2\Delta_u}{\epsilon}(t + \log(|\mathcal{Y}|))\right] \leq \exp(-t).$$

**Remark 1** *The problem addressed in Proposition 5 can somehow be seen as a particular case of the problem at hand in this section (though, strictly speaking, they correspond to two different settings: continuous versus discrete $\mathcal{Y}$). Indeed, estimating $f(D) \in \mathbb{R}^k$ is equivalent to producing a vector $y \in \mathbb{R}^k$ such that the utility*

$$u(D, y) = -\|f(D) - y\|_1$$

*is as large as possible. Note that in this case:*

- we have $\Delta_u = \Delta_f$;

- the Laplace mechanism is equivalent to the exponential mechanism if we used it with the Lebesgue measure on $\mathbb{R}^k$ (instead of the counting measure on $\mathcal{Y}$). The minor difference is that the scale parameter is $2\Delta_u/\epsilon$ instead of $\Delta_f/\epsilon$; this factor of 2 is due to the fact that the utility setting is in a sense more general.

# References

[1] C. Dwork and F. McSherry and K. Nissim, K and 1. Smith. *Calibrating noise to sensitivity in private data analysis.* Proceedings of Theory of Cryptography Conference, pp. 265-284. Springer Berlin Heidelberg (2006).

[2] C. Dwork. *Differential privacy: A survey of results.* Proceedings of the international conference on Theory and applications of models of computation, Springer Berlin Heidelberg 4078 (2008).

[3] C. Dwork and R. Aaron. *The algorithmic foundations of differential privacy.* Foundations and Trends in Theoretical Computer Science 9(3-4):211-407 (2014).

[4] M. Gymrek and A.L. McGuire and D. Golan E. Halperin and Y. Erlich. *Identifying personal genomes by surname inference.* Science, 339(6117), 321-324 (2013).

[5] N. Homer, and S. Szelinger, and M. Redman and D. Duggan and W. Tembe and J. Muehling and J.V. Pearson and D.A. Stephan S.F. Nelson and D.W. Craig *Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays.* PLoS Genet, 4(8), p.e1000167 (2008).

[6] J. Kleinberg and C. Papadimitriou R. Prabhakar Raghavan. *Auditing boolean attributes.* Journal of Computer and System Sciences 66(1):244-253 (2003).

[7] F. McSherry and K. Talwar. *Mechanism design via differential privacy.* Foundations of Computer Science, pages 94-103 (2007).

[8] Y. de Montjoye and C.A. Hidalgo and M. Verleysen and V.D. Blondel. *Unique in the crowd: The privacy bounds of human mobility.* Scientific reports, 3 2013.

[9] A. Narayanan and S. Vitaly Shmatikov. *Robust de-anonymization of large sparse datasets.* IEEE Symposium on Security and Privacy (2008).

[10] L. Sweeney. *k-anonymity: A model for protecting privacy.* International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems 10(5):557-570 (2002).

[11] H. Zang and J. Bolot. *Anonymization of location data does not work: A large-scale measurement study.* Proceedings of the 17th annual international conference on Mobile computing and networking (pp. 145-156). ACM (2011).