

# Statistical Quantification of Differential Privacy: A Local Approach

Önder Askin  
Ruhr-University Bochum  
oender.askin@rub.de

Tim Kutta  
Ruhr-University Bochum  
tim.kutta@rub.de

Holger Dette  
Ruhr-University Bochum  
holger.dette@rub.de

**Abstract**—In this work we introduce a new approach for statistical quantification of differential privacy in a black box setting. We present estimators and confidence intervals for the optimal privacy parameter of a randomized algorithm  $A$ , as well as other key variables (such as the novel "data-centric privacy level"). Our estimators are based on a local characterization of privacy and in contrast to the related literature avoid the process of "event selection" - a major obstacle to privacy validation. This makes our methods easy to implement and user-friendly. We show fast convergence rates of the estimators and asymptotic validity of the confidence intervals. An experimental study of various algorithms confirms the efficacy of our approach.

**Index Terms**—Differential privacy, data-centric privacy, local estimators, confidence intervals

## I. INTRODUCTION

Since its introduction in the seminal work of [1], the concept of *Differential Privacy* (DP) has become a standard tool to assess information leakage in data disseminating procedures. DP characterizes how strongly the output of a randomized algorithm is influenced by any one of its inputs, thus quantifying the difficulty of inferring arguments (i.e. user information) from algorithmic releases.

To formalize this situation, we consider a data base  $x = (x(1), \dots, x(m))$  where each data point  $x(i)$  takes values in a set  $\mathcal{D}$  and corresponds to the data provided by the  $i$ th individual among  $m$  users. Furthermore, we introduce the notion of *neighboring* or *adjacent* data bases, that is data bases that only differ in one component. Mathematically we can express neighborhood of  $x, x'$  by unit Hamming distance  $d_H(x, x') = 1$ , where the Hamming distance is defined as follows:

$$d_H(x, x') := |\{1 \leq i \leq m : x(i) \neq x'(i)\}|.$$

**Definition 1.** An Algorithm  $A$  is called  $\epsilon$ -differentially private for some  $\epsilon > 0$ , if for any two neighboring data bases  $x, x'$  and any measurable event  $E$  the inequality

$$\mathbb{P}(A(x) \in E) \leq e^\epsilon \mathbb{P}(A(x') \in E) \quad (1)$$

holds.

Definition 1 demands that (1) holds for all measurable events  $E$ , but what constitutes a measurable event depends on the output space  $\mathcal{Y}$  of the randomized algorithm  $A$ . If  $\mathcal{Y}$  is discrete (in particular if  $|\mathcal{Y}| < \infty$ ) we require that (1) holds for all events in the power set  $\mathcal{P}(\mathcal{Y})$ . If however  $A$  has outputs in

a continuum (e.g.  $\mathcal{Y} = \mathbb{R}^d$ ), then (1) has to hold for all Borel sets. In both cases, the collection of all measurable events is large and complex, which is an important obstacle in the practical validation of DP as we will discuss below.

The privacy parameter  $\epsilon$  in Definition 1 quantifies the information leakage of  $A$ , where small values correspond to small leakage (and thus high privacy). Hence, deploying differentially private algorithms with appropriate  $\epsilon$  provides users with strong privacy guarantees regarding their data. However, in practice it is often unclear whether an algorithm satisfies DP and if so, for which parameter  $\epsilon$ . It is therefore of crucial importance and the main objective of this work to develop procedures by which we can ascertain the level of privacy afforded by a given algorithm.

**Related work:** A number of languages and verification tools have been devised to validate differential privacy where possible and discard it where not (see [2]–[13] among others). Many of these approaches are designed specifically for developers and require knowledge of the inner structure of the algorithm in question. In contrast, in this paper, we want to investigate a black box scenario where we have little to no knowledge of the algorithm's design and have to rely solely on output samples. This scenario can occur naturally when a user relies on an algorithm provided by a third party that does not want to reveal its design (entirely). However, black box methods can also be valuable in settings where an algorithm is known but so complex, that focusing on its outputs is preferable. In any case, a procedure tailored to this scenario covers a wide range of algorithms with few requirements, which is a desirable feature in a validation scheme.

Relying solely on algorithmic outputs requires a statistical validation scheme of differential privacy. Such an approach is pursued in [14], built directly on Definition 1. For a fixed triplet  $(x, x', E)$  consisting of neighboring data bases  $x, x'$  and an event  $E$ , the privacy condition in (1) can be construed as a statistical hypothesis that needs to be checked. Given a preconceived privacy parameter  $\epsilon_0 > 0$ , candidate triplets are generated and a binomial statistical test is employed to find a counterexample  $(x_0, x'_0, E_0)$  that violates the privacy condition (1). These counterexamples expose faulty, non-private algorithms in a fast and practical manner and hint at potential deficiencies in the algorithm's design.

A related, but distinct approach is the examination of lower

bounds for differential privacy [15]. Here, privacy violations are determined with the help of the "privacy loss", which is defined for any triplet  $(x, x', E)$  as

$$L_{x,x'}(E) := \left| \ln(\mathbb{P}(A(x) \in E)) - \ln(\mathbb{P}(A(x') \in E)) \right|. \quad (2)$$

We interpret  $\infty - \infty := 0$  to account for events with 0 probability. In line with Definition 1, an algorithm  $A$  satisfies  $\epsilon$ -DP if and only if  $L_{x,x'}(E) \leq \epsilon$  for all permissible triplets. Thus, computing privacy violations  $L_{x,x'}(E)$  for different triplets naturally provides lower bounds for  $\epsilon$ . Note that in this context privacy violations and loss are used constructively to gather information about the privacy parameter. We also want to point out that this approach can be adapted to counterexample generation, if for some predetermined  $\epsilon_0$  a triplet  $(x_0, x'_0, E_0)$  is found s.t.  $L_{x_0,x'_0}(E_0) > \epsilon_0$ . However, lower bounds are somewhat more flexible, because they do not require some hypothesized  $\epsilon_0$  in the first place.

Even though [14] and [15] provide effective tools for privacy validation, they are not entirely compatible with our black box assumption. While the binomial test in [14] by itself requires little knowledge of  $A$ , the larger scheme, within which it is embedded, is designed to also consider the algorithm's program code. A symbolic execution of that code can be performed to facilitate the detection of counterexamples. Therefore, this approach is also labeled *semi-black-box* by the authors [14]. Even less compatible with the black box regime, the approach in [15] requires access to the program code of algorithm  $A$  in order to alter it in ways that produce a differentiable surrogate function for  $L_{x,x'}$ . Numerical optimizers can then be deployed to find triplets that yield high privacy violations.

A method that is based on the loss function  $L_{x,x'}$  and adheres to the black box setting is introduced in [16]. In pursuit of high privacy violations and given data bases  $x$  and  $x'$ , the goal is to maximize the difference in (2) by constructing events  $E$  containing output values that have been more likely generated by  $A(x)$  as opposed to  $A(x')$ . A machine learning classifier with that objective is employed to approximate posterior probabilities for output values and  $E$  is pieced together with output values whose posterior probability surpasses a given threshold (which is done to account for computational instabilities).

**The problem of event selection:** The task of finding a triplet  $(x, x', E)$  that provokes a high privacy violation is typically split into two separate parts: First finding a data base  $(x, x')$  such that the loss  $L_{x,x'}(E)$  is large for some event  $E$  and second finding this very event. Even though both problems are intractable, the greater challenge lies in the latter one, the *event selection*. This is due to the above mentioned intricacy of the space of measurable events (for a discussion see [16] and references therein). To illustrate the resulting statistical problem, suppose we want to estimate for fixed data bases  $x, x'$  the maximum of  $L_{x,x'}(E)$  over all measurable events  $E$ . Typically this is done by maximizing an empirical version of the loss, say  $\hat{L}_{x,x'}$  over a number of events  $E_1, \dots, E_M$ , yielding  $\max_{m=1, \dots, M} \hat{L}_{x,x'}(E_m)$ . However this

estimate depends completely on the choice of  $E_1, \dots, E_M$ , which raises practical questions, such as how many events are enough to maximize over. As we have seen above, the space of events is not practically exhaustible and a higher number of  $M$  requires larger samples for statistical estimation (and thus increased run-time). Furthermore it is not a priori clear which kinds of events work best: There is no one-size fits all event collection that works for any algorithm and so (in principle) the sets  $E_1, \dots, E_M$  have to be adjusted every time. Such an adjustment is not only troublesome for non-expert users, but also potentially runs counter to a black-box regime, because it may presuppose prior knowledge. So, despite recent attempts to streamline event selection via randomization [16], it remains the hard nut of statistical privacy validation.

In this work, we discuss an alternative route to the assessment of DP without the intermediary step of maximizing the empirical loss  $\hat{L}_{x,x'}$ . Instead, our estimators rely on the *local loss function* to approximate the maximum of  $L_{x,x'}(\cdot)$  directly (see Theorem 1) and thus circumvent event selection.

**This work:** In this work a central object of interest is the quantity

$$\epsilon_{x,x'} := \sup_E L_{x,x'}(E) \quad (3)$$

which we call *data-specific privacy violation* in  $x$  and  $x'$ . Recalling (2), we observe that  $\epsilon_{x,x'}$  indicates to which extent the algorithm outputs are indistinguishable for a fixed pair of data bases  $x$  and  $x'$ . Note that  $A$  satisfies  $\epsilon_0$ -DP if and only if  $\epsilon_{x,x'} \leq \epsilon_0$  for all pairs of adjacent data bases  $(x, x')$ . Thus, we define the smallest parameter  $\epsilon$ , for which  $\epsilon$ -DP still holds as

$$\epsilon := \sup_{x, x': d_H(x, x')=1} \epsilon_{x,x'}, \quad (4)$$

and refer to it as the *true privacy parameter* (privacy guarantees below  $\epsilon$  are not intractable while any  $\epsilon_0 > \epsilon$  underestimates the privacy level that is actually achievable).

We occasionally refer to  $\epsilon$  as the *global privacy parameter* which, in light of identity (4), only provides a "worst-case" guarantee for privacy leakage of any pair  $x, x'$ . In contrast the precise amount of privacy leakage associated with  $x$  and  $x'$  is captured by  $\epsilon_{x,x'}$ , which is potentially much smaller than  $\epsilon$ . The data-specific privacy violations comprise more granular and local information that we utilize to examine the following privacy aspects:

First, each  $\epsilon_{x,x'}$  constitutes a lower bound of  $\epsilon$ . Because  $L_{x,x'}(E) \leq \epsilon_{x,x'}$  holds for all events  $E$ , these lower bounds are more powerful than the ones derived in prior work. Lower bounds in themselves are useful as they can help expose faulty algorithms [15] and narrow down the extent to which a given algorithm can be private at all [16]. This ultimately provides us with a better understanding of the true privacy parameter  $\epsilon$ .

Secondly, data-specific privacy violations can be used to infer the *data-centric privacy level* for select data bases. More precisely, suppose that a curator has gathered a data base  $x$  and is interested in the amount of privacy conceded

specifically to  $x$ . The maximum privacy violation associated with  $x$  is obtained by forming the supremum over all data-specific privacy violations in its neighborhood, that is

$$\epsilon_x := \sup_{x': d_H(x, x')=1} \epsilon_{x, x'}. \quad (5)$$

Evidently, any procedure determining  $\epsilon_{x, x'}$  can be used to gather information about the data-centric privacy guarantee for specific data bases.

Statistically our approach is based on novel estimators  $\hat{\epsilon}_{x, x'}$  for the data-specific privacy violation  $\epsilon_{x, x'}$ . In view of the identities (4) and (5), such estimates are natural building blocks for the assessment of the true privacy parameter  $\epsilon$  or its data-centric version  $\epsilon_x$ . Contrary to the related literature, our estimators do not maximize an empirical version of the loss  $L_{x, x'}$ , but approximate the supremum  $\epsilon_{x, x'}$  directly, thus avoiding the pitfalls of event selection (see previous part). Mathematically, these estimates rest on a local version of the privacy loss discussed in Section III. Besides estimators we present new tools of statistical inference: In Section IV we devise the MPL algorithm, which generates one-sided confidence intervals  $[LB, \infty)$  for the privacy parameters  $\epsilon$  and  $\epsilon_x$  respectively. In this situation  $LB$  is a statistical lower bound (i.e. it holds with a high degree of certainty) and approximates the true parameter with increasing sample size. In particular, if MPL is applied to the quantification of  $\epsilon$  and outputs  $LB$ , the user can be confident that algorithm  $A$  is at best  $LB$ -differentially private. In Section V we confirm these findings via experiments.

**Main contributions:** We give a brief summary of our main contributions:

- A fully statistical black box procedure for the quantification of DP.
- A flexible approach based on data-specific privacy violations  $\epsilon_{x, x'}$ , that facilitates the study of global and local aspects of privacy.
- New estimators  $\hat{\epsilon}_{x, x'}$  for the data-specific privacy violation that circumvent the problem of event selection and are proved to converge at a fast rate.
- The MPL algorithm (Maximum Privacy Loss) that outputs a confidence interval for  $\epsilon$  (or  $\epsilon_x$ ), which demonstrably includes the parameter of interest with approximate level of confidence.
- A practical evaluation and validation of our methods.

## II. STATISTICAL PRELIMINARIES

In this section we review the statistical concepts of *confidence intervals* and *kernel density estimation*, which serve as technical background for the remainder of this paper. Readers who are only interested in discrete algorithms can omit Section II-B.

### A. Confidence Intervals

A confidence interval is a statistical method to localize a parameter of a probability distribution with a prescribed level of certainty. More concretely, consider a sample of

$n$  observations  $X_1, \dots, X_n$  (random variables), following an unknown distribution  $P$ . If a user is interested in a parameter  $\theta = \theta(P)$  derived from  $P$  (e.g. the expectation  $\theta := \mathbb{E}_P X_1$ ), the sample of observations can be used to approximately locate  $\theta$  in an interval  $\hat{I}(X_1, \dots, X_n) \subset \mathbb{R}$ . Notice that the term *confidence interval* usually refers to both the output  $\hat{I}(X_1, \dots, X_n)$ , which is an interval determined by the data, as well as the underlying algorithm  $\hat{I}(\cdot)$  itself. Given the randomness in the data, there is always a risk of mislocating  $\theta$ , i.e. that  $\theta \notin \hat{I}(X_1, \dots, X_n)$ . However, confidence intervals are constructed to guarantee  $\theta \in \hat{I}(X_1, \dots, X_n)$  with a prescribed probability (level of confidence). To be more precise,  $\hat{I}(\cdot)$  has an additional input parameter  $\alpha \in (0, 1)$ , such that the *confidence level*  $1 - \alpha$  holds:

$$\mathbb{P}(\theta \in \hat{I}_\alpha(X_1, \dots, X_n)) = 1 - \alpha, \quad (6)$$

where typically  $\alpha \in \{0.1, 0.05, 0.01\}$ . Notice that the choice of  $\alpha$  entails a trade-off: On the one hand a smaller  $\alpha$  provides the user with higher certainty that actually  $\theta \in \hat{I}_\alpha(X_1, \dots, X_n)$ , but on the other hand it translates into a wider confidence interval, which means less precision with regard to the location of  $\theta$ . Besides the choice of  $\alpha$ , the sample size  $n$  affects the width of the confidence interval, with larger  $n$  leading to narrower intervals.

In order to construct a confidence interval  $\hat{I}_\alpha$  s.t. (6) holds, it is necessary to have prior knowledge about the underlying distribution of the data sample  $X_1, \dots, X_n$ . For instance, it may be known that the sample comes from a normal distribution, with unknown mean and variance, and we want to give a confidence interval for the mean. In this situation parametric statistical theory equips the user with standard tools to construct  $\hat{I}_\alpha$  (see [17]).

Yet in many cases such prior knowledge about the data is not feasible and therefore a weaker requirement than (6) is reformulated: It states that the confidence confidence level  $1 - \alpha$  is approximated with increasing precision, as  $n$  grows larger, or mathematically speaking

$$\lim_{n \rightarrow \infty} \mathbb{P}(\theta \in \hat{I}_\alpha(X_1, \dots, X_n)) = 1 - \alpha. \quad (7)$$

If (7) is satisfied, we call  $\hat{I}_\alpha$  an *asymptotic confidence interval with confidence level*  $1 - \alpha$ . The advantages of asymptotic confidence intervals are their flexibility and robustness against deviations from a presumed distribution. Common approaches to prove asymptotic confidence levels include asymptotically normal estimators as well as the Delta-method for differentiable statistics. For details on asymptotic statistical theory we refer the interested reader to the monograph of [18].

### B. Kernel density estimation

Kernel density estimation is a method to estimate the unknown distribution of a data sample  $X_1, \dots, X_n$  on  $\mathbb{R}^d$ . It can be thought of as the creation of a smoothed, normalized histogram, where the jumps between the bins are interpolated continuously (for an introduction see [19]). This procedure is often preferred to a traditional histogram, particularly if the

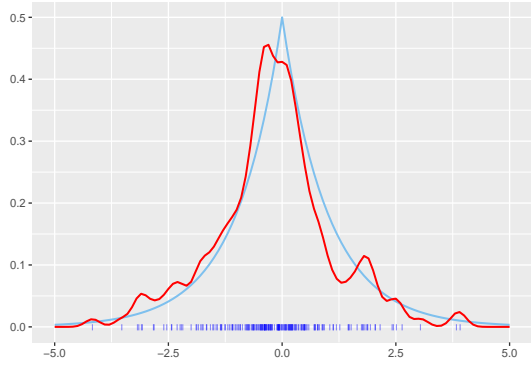


Fig. 1: Centered Laplace density with  $\lambda = 1$  (light blue) and kernel density estimate (red) for  $N = 200$ , with Gaussian kernel. On the  $x$ -axis we have plotted the observations  $X_1, \dots, X_{200}$  (dark blue).

data sample is distributed according to a continuous density  $f$  on  $\mathbb{R}^d$  (we write  $X_1, \dots, X_n \sim f$ ).

More precisely, let  $K : \mathbb{R}^d \rightarrow \mathbb{R}$  be a continuous, non-negative function, such that  $\int_{\mathbb{R}^d} K(u) du = 1$ . We call  $K$  a kernel and define the *kernel density estimator* (KDE)  $\tilde{f}$  for  $f$  pointwise as

$$\tilde{f}(t) := \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{t - X_i}{h}\right), \quad t \in \mathbb{R}^d, \quad (8)$$

where  $h > 0$  is the *bandwidth*, analogue to the bin-width in a histogram. For details on kernel density estimators as well as generalizations such as multi-dimensional bandwidths, we refer to [20]. As the number of observations  $n$  increases, the convergence speed of  $\tilde{f}$  to  $f$  depends on three distinct factors: First the smoothness of the true density  $f$ , secondly an adequate choice of the kernel  $K$  and thirdly the bandwidth  $h$ .

To quantify smoothness we require  $f$  to be *Hölder continuous*, i.e. for some  $\beta \in (0, 1]$  and  $C > 0$  it holds that

$$|f(t) - f(s)| \leq C|t - s|^\beta, \quad \forall t, s \in \mathbb{R}^d, \quad (9)$$

where  $|\cdot|$  denotes the Euclidean norm. Notice that  $\beta = 1$  corresponds to the well known *Lipschitz continuity*, which is satisfied by the densities corresponding to the Laplace, Gaussian and versions of the Exponential Mechanism. We also point out that a density which satisfies Hölder continuity for one  $\beta > 0$  is Hölder continuous for any other  $\beta' \in (0, \beta]$ .

The choice of the kernel  $K$  is a relatively simple task: To attain optimal convergence speed,  $K$  has to fulfill certain regularity properties (K1) and (K2), that we make precise in Appendix B. From now on we will always assume that  $K$  conforms to these assumptions. We point out that both of them are satisfied by all commonly used kernels (in particular by the Gaussian kernel, that we use in our experiments).

Finally, the choice of the bandwidth  $h$  should depend on the smoothness level  $\beta$  of  $f$ , as well as the sample size  $n$ . More

precisely it can be shown that

$$\sup_{t \in \mathbb{R}^d} |\tilde{f}(t) - f(t)| = \mathcal{O}_P\left(h^\beta + \sqrt{\frac{\ln(n)}{h^{d+n}}}\right), \quad (10)$$

which implies for the specific choice  $h = \mathcal{O}(n^{-\frac{1}{2\beta+d}})$

$$\sup_{t \in \mathbb{R}^d} |\tilde{f}(t) - f(t)| = \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+d}}\right). \quad (11)$$

Notice that this  $h$  minimizes the error rate (except for log-terms). For details on convergence rates in density estimation see [21] and for a definition of the stochastic Landau symbol  $\mathcal{O}_P$  we refer to the Appendix A.

In practical applications the true smoothness  $\beta$  and hence the optimal bandwidth is unknown and therefore data-driven procedures, such as cross validation are used to determine it. For details on bandwidth selection see [20].

In the subsequent discussion we consider log-transformed density estimators. These objects are potentially unstable for arguments where the true density  $f$  is close to 0, because small errors in the estimate of  $f$  translate into great errors in the logarithm. For this reason we define the truncated KDE pointwise in  $t$  as

$$\hat{f}(t) := \tilde{f}(t) \vee \delta,$$

where " $a \vee b$ " denotes the maximum of two numbers  $a, b \in \mathbb{R}$  and  $\delta > 0$  is a user-determined floor. In Section IV we discuss how to choose  $\delta$  dependent on  $n$  and  $\beta$ . The construction of the truncated KDE is described in Algorithm 1.

---

#### Algorithm 1 Truncated kernel density estimator

---

**Input:** data sample  $X = (X_1, \dots, X_n)$ , evaluation point  $t$ , bandwidth  $h$ , kernel function  $K$ , floor  $\delta$

```

1: function TKDE( $X, t, h, K, \delta$ )
2:  $out = 0$ 
3: for  $i = 1, 2, \dots, n$  do
4:    $out = out + K((t - X_i)/h)$ 
5: end for
6:  $out = out/(nh^d)$ 
7: return  $out \vee \delta$ 
8: end function
```

---

### III. DIFFERENTIAL PRIVACY AS A LOCAL PROPERTY

As we have seen in our Introduction,  $\epsilon$ -DP means that for any neighboring data bases  $x, x'$  the bound

$$\epsilon_{x, x'} = \sup_E L_{x, x'}(E) \leq \epsilon \quad (12)$$

holds, where the loss  $L_{x, x'}$  is defined in (2). Thus, in principle, validating DP requires the calculation of  $L_{x, x'}(E)$  for any measurable event  $E$ , a problem that is intractable from a practical perspective given the complexity of the space of measurable events (see Introduction). We can, however, drastically reduce the effort of *event selection* in the supremum by exploiting that differential privacy is an inherently *local property*, i.e. that the level of privacy is determined by the loss on small events. To get an intuition of this point, consider an

event  $E$  that can be decomposed into the disjoint subsets  $E_1$  and  $E_2$ . It is a simple exercise to show that

$$L_{x,x'}(E) \leq \max\{L_{x,x'}(E_1), L_{x,x'}(E_2)\}.$$

In this sense going from larger to smaller events increases the privacy loss and thus gets us closer to  $\epsilon_{x,x'}$ . Iterating this process suggests that we should look at "the smallest events possible", which are single points. So we expect that ultimately

$$\epsilon_{x,x'} \approx \sup_{t \in \mathcal{Y}} |L_{x,x'}(\{t\})|. \quad (13)$$

Admittedly, this statement is not formally correct for all algorithms, but we will make it rigorous for certain classes of algorithms in the course of this section. Compared with the supremum over all measurable events in (12), the expression in (13) is more convenient, because single points are easy to handle. We will explore this advantage in detail at the end of this section.

We now begin our formal discussion by specifying two classes of algorithms that are considered throughout this work: Discrete and continuous ones.

We call an algorithm  $A$  that maps a data base  $x$  to random values in either a finite or a countably infinite set  $\mathcal{Y}$  a *discrete algorithm*. Without loss of generality we will assume that  $\mathcal{Y} \subset \mathbb{N}$ . Moreover, we call the corresponding probability function  $f_x : \mathcal{Y} \rightarrow [0, 1]$  defined as

$$f_x(t) := \mathbb{P}(A(x) = t), \quad \forall t \in \mathcal{Y} \quad (14)$$

the *discrete density* of  $A$  in  $x$ . With this notation we can write for any  $E \subset \mathcal{Y}$

$$\mathbb{P}(A(x) \in E) = \sum_{t \in E} f_x(t). \quad (15)$$

Examples of discrete algorithms include Randomized Response [22], Report Noisy Max [23] and the Sparse Vector Technique [24].

Next suppose that  $\mathcal{Y} = \mathbb{R}^d$ . We say that  $A$  is a *continuous algorithm*, if for any data base  $x$ ,  $A(x)$  has a continuous density  $f_x : \mathbb{R}^d \rightarrow \mathbb{R}$ , such that for any Borel measurable event  $E$

$$\mathbb{P}(A(x) \in E) = \int_E f_x(t) dt.$$

Typical examples of continuous algorithms are, as mentioned before, the Laplace [23], the Gaussian [23] and versions of the Exponential Mechanism [25]. We want to highlight that in this definition the requirement of continuous densities on the whole space  $\mathbb{R}^d$  is only made for convenience of presentation and can be relaxed to densities on subsets, e.g.  $[0, \infty) \subset \mathbb{R}$  in the case  $d = 1$ . Notice that for continuous algorithms (13) is technically invalid because  $L_{x,x'}(\{t\}) = 0$  for any point  $t$ . However, it is possible to preserve the idea of (13) by reformulating it in terms of continuous densities (see Theorem 1).

Given the above definitions, the distribution of an algorithm  $A$  can be thoroughly characterized by its densities and we use the notation  $A(x) \sim f_x$  throughout this paper. In the following theorem we give a mathematically rigorous version of (13), proving that DP is a local property. Variants of this theorem, particularly the inequality " $\leq$ " in (16), are frequently used in the literature. However, the precise characterization is not trivial and therefore shown here explicitly.

**Theorem 1.** *Given a discrete or continuous algorithm  $A$  with  $A(x) \sim f_x$  and  $A(x') \sim f_{x'}$  we have*

$$\epsilon_{x,x'} = \sup_{t \in \mathcal{Y}} |\ln(f_x(t)) - \ln(f_{x'}(t))|, \quad (16)$$

where  $\infty - \infty := 0$ .

*Proof:* We first consider the discrete setting: In order to show " $\geq$ " we notice that for all  $t \in \mathcal{Y}$

$$L_{x,x'}(\{t\}) = |\ln(f_x(t)) - \ln(f_{x'}(t))|.$$

Recall that  $\epsilon_{x,x'} = \sup_E |L_{x,x'}(E)|$ . Here the supremum is taken over all elements  $E$  of the power set  $\mathcal{P}(\mathcal{Y})$  (which includes in particular sets with only one element) and this directly implies " $\geq$ ".

The proof of " $\leq$ " follows by standard techniques. We fix a set  $E \subset \mathcal{Y}$  and rewrite  $L_{x,x'}(E)$  using (15), s.t.

$$L_{x,x'}(E) = \left| \ln \left( \frac{\sum_{t \in E} f_x(t)}{\sum_{t \in E} f_{x'}(t)} \right) \right|. \quad (17)$$

Without loss of generality we assume that the numerator is greater than the denominator and we can therefore drop the absolute value. Now the inner fraction can be upper bounded as follows:

$$\frac{\sum_{t \in E} f_x(t)}{\sum_{t \in E} f_{x'}(t)} \leq \frac{\sum_{t \in E} f_{x'}(t) [f_x(t)/f_{x'}(t)]}{\sum_{t \in E} f_{x'}(t)} \leq \sup_{t \in \mathcal{Y}} \frac{f_x(t)}{f_{x'}(t)}.$$

Taking the logarithm on both sides and the supremum over all  $E$  on the left maintains the inequality, showing " $\leq$ ".

Moving to continuous algorithms we notice that the proof of " $\leq$ " follows along the same lines as for the discrete case and is therefore omitted (one simply has to replace all the sums by integrals).

To prove " $\geq$ " we first observe that a probability density in  $t$  gives the probability of a very small region around  $t$ . More precisely it can be expressed as follows

$$f_x(t) = \lim_{\delta \rightarrow 0} \frac{\mathbb{P}(A(x) \in U_\delta(t))}{\text{vol}(U_\delta(t))},$$

where  $U_\delta(t) := \{s \in \mathcal{Y} : |t - s| \leq \delta\}$  and  $\text{vol}()$  denotes the  $d$ -dimensional volume. The identity is a special case of Theorem 6.20 (c) in [26]. The same statement holds for  $x'$  instead of  $x$  and we can use that to get

$$\frac{f_x(t)}{f_{x'}(t)} = \lim_{\delta \rightarrow 0} \frac{\mathbb{P}(A(x) \in U_\delta(t))}{\mathbb{P}(A(x') \in U_\delta(t))} \leq \sup_E \frac{\mathbb{P}(A(x) \in E)}{\mathbb{P}(A(x') \in E)}$$

for any  $t \in \mathcal{Y}$ . Taking the logarithm on both sides and the supremum over  $t$  on the left preserves the inequality. Recalling

(3), this implies  $\sup_{t \in \mathcal{Y}} |\ln(f_x(t)) - \ln(f_{x'}(t))| \leq \epsilon_{x,x'}$ , which proves the theorem. ■

Theorem 1 allows us to characterize DP of an algorithm  $A$  by the absolute log-difference of the algorithm's densities. For ease of reference we define this difference, the *loss function*, explicitly as

$$\ell_{x,x'}(t) := |\ln(f_x(t)) - \ln(f_{x'}(t))|. \quad (18)$$

This definition admits the restatement of Theorem 1 as  $\epsilon_{x,x'} = \sup_{t \in \mathcal{Y}} \ell_{x,x'}(t)$ .

Figure 2 provides an illustration of the loss function for some standard examples of randomized algorithms (see e.g. [22], [23]). The plots help discern the amount of privacy leakage and where it occurs. For example, we observe that for Randomized Response (left) only two outputs elicit any privacy leakage at all, while the maximum loss associated with the Laplace Mechanism (middle panel) is assumed everywhere, except for the area enclosed by the density modes. For the Gaussian Mechanism (right panel) no single  $t$  exists that maximizes the loss. Instead  $\ell_{x,x'}(t)$  tends to infinity for growing  $|t|$ , which implies decreasing privacy for tail events.

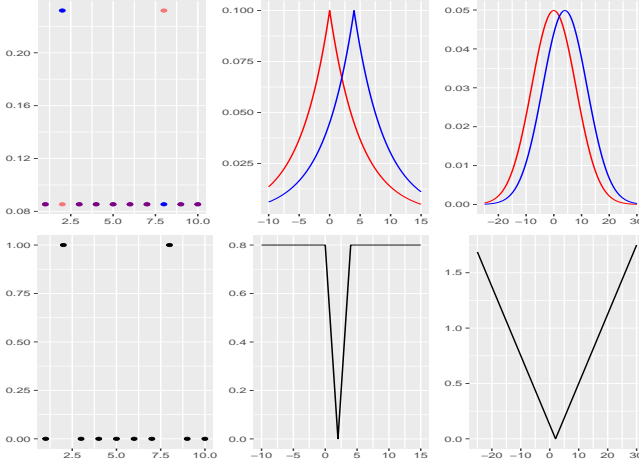


Fig. 2: The top row depicts the densities  $f_x \sim A(x)$ ,  $f_{x'} \sim A(x')$  for two neighboring data bases  $x, x'$  and algorithm  $A$  chosen (from left to right) as Randomized Response, the Laplace Mechanism and Gaussian Mechanism. The bottom row captures the corresponding loss functions  $\ell_{x,x'}$  from (18).

We briefly summarize the **key insights of this section**:

Instead of examining large and complex sets in order to quantify  $\epsilon_{x,x'}$ , Theorem 1 shows that it suffices to consider single output values  $t \in \mathcal{Y}$ . In fact, larger events  $E$  potentially dilute the observed privacy violation and lead to an underestimation of privacy leakage. Numerically, the task of maximizing  $L_{x,x'}$  (a function with sets as arguments), is much more difficult than to maximize  $\ell_{x,x'}$  (which has arguments in

$\mathbb{R}^d$  or  $\mathbb{N}$ ), where standard solutions exist (see [27]). Finally, the loss function  $\ell_{x,x'}$  is far more amenable to interpretation than  $L_{x,x'}$ . In fact,  $\ell_{x,x'}$  can be plotted and thus problematic areas with respect to privacy can be easily displayed and understood (e.g. we see at one glance, that for the Gaussian Mechanism the problem lies in extreme values of  $t$ ; see Figure 2, right).

We conclude this section with a non-trivial example, where we utilize the loss function to derive the true privacy parameter.

**Example 1.** We consider a data base  $x$  containing the information of only one individual ( $m = 1$ ). Assuming that said individual's data is a vector  $v = (v_1, \dots, v_k) \in [0, 1]^k$ , i.e.  $\mathcal{D} = [0, 1]^k$ , we can identify our data base as  $x = v$ . It is our intention to publish the maximum entry of  $v$  in a differentially private manner. We can do this by employing a version of the Noisy Max algorithm (Algorithm 7 in [14]) where we add independent Laplace noise  $L_i \sim \text{Lap}(0, \frac{1}{\lambda})$  to each component  $v_i$  and publish the maximum  $\max_i(v_i + L_i)$ . We demonstrate how  $\ell_{x,x'}$  can be used to determine the true privacy parameter  $\epsilon$  of this algorithm.

On the one hand, releasing a noisy component  $v_i + L_i$  by itself satisfies  $\lambda$ -DP by virtue of the Laplace Mechanism. The maximum can then be understood as a function over the vector of noisy components and the composition theorem of DP yields  $k\lambda$  as an upper bound of  $\epsilon$ . On the other hand, define  $F_i$  as the distribution function of  $v_i + L_i$  and  $f_i = F'_i$  as the corresponding density. Then the density  $f_v$  of the random variable  $\max_i(v_i + L_i)$  is of the form

$$f_v(t) = \left( \sum_{i=1}^k \frac{f_i(t)}{F_i(t)} \right) \left( \prod_{i=1}^k F_i(t) \right).$$

In the case where  $v_1 = \dots = v_k$ , this can be simplified to  $f_v(t) = k f_1(t) [F_1(t)]^{k-1}$ . Using this formula, it is a straightforward calculation to show that for  $v = (0, \dots, 0)$ ,  $w = (1, \dots, 1)$  and sufficiently large  $t \in \mathbb{R}$

$$\ell_{v,w}(t) = |\ln(f_v(t)) - \ln(f_w(t))| = k\lambda.$$

Theorem 1 especially implies that  $k\lambda$  is also a lower bound of  $\epsilon$  and thus the equality  $\epsilon = k\lambda$  holds.

#### IV. QUANTIFYING THE MAXIMUM PRIVACY VIOLATION

In this section we proceed to the statistical aspects of our discussion. According to Theorem 1 the data-specific privacy violation  $\epsilon_{x,x'}$  defined in (3) can be attained by maximizing the loss function  $\ell_{x,x'}$  defined in (18). We devise an estimator  $\hat{\epsilon}_{x,x'}$  for  $\epsilon_{x,x'}$ , by maximizing an empirical version  $\hat{\ell}_{x,x'}$  of the loss function, specified in Section IV-A. In Proposition 1, we demonstrate mathematically that such estimators are consistent with fast convergence rates. Besides estimation, we consider confidence intervals for the pointwise privacy loss  $\ell_{x,x'}(t)$  in Section IV-B. If applied to a  $t^*$  close to the argmax of  $\ell_{x,x'}$ , these can be used to statistically locate  $\epsilon_{x,x'} \approx \ell_{x,x'}(t^*)$ .

Next recall that the true privacy parameter  $\epsilon$  as well as the data-centric privacy level  $\epsilon_x$ , defined in (4) and (5) respectively, can be attained by maximizing  $\epsilon_{x,x'}$  over a (sub)space of data

bases. It therefore makes sense to approximate them (from below) by a finite maximum, s.t. for instance

$$\epsilon \approx \max(\epsilon_{x_1, x'_1}, \dots, \epsilon_{x_B, x'_B}), \quad (19)$$

where  $(x_1, x'_1), \dots, (x_B, x'_B)$  are  $B$  pairs of adjacent data bases (approximating  $\epsilon_x$  works by setting  $x = x_1 = \dots = x_B$ ). If the data bases are chosen appropriately, the maximum on the right side of (19) comes arbitrarily close to  $\epsilon$ . Prior work suggests that oftentimes simple heuristics already yield data bases that point to the true privacy parameter  $\epsilon$  [14]. Furthermore, the structure of the data space  $\mathcal{D}$  can naturally motivate search patterns (typically choosing  $x_b$  and  $x'_b$  to be "far apart" in some sense).

Combined with our statistical tools for the data-specific privacy violations, we apply the approximation in (19) to estimation and statistical inference for the parameters  $\epsilon$  and  $\epsilon_x$ . We integrate these methods into the MPL algorithm presented in Section IV-C and demonstrate that its output  $[LB, \infty)$  is a one-sided, asymptotic confidence interval (Theorem 2).

#### A. Estimating data-specific privacy violations

We now consider the problem of estimating the data-specific privacy violation  $\epsilon_{x, x'}$  for two adjacent data bases  $x, x'$  defined in (3). According to Theorem 1 we can express  $\epsilon_{x, x'}$  as the maximum of the loss function, i.e.

$$\epsilon_{x, x'} = \sup_{t \in \mathcal{Y}} \ell_{x, x'}(t),$$

where  $\ell_{x, x'}$  is defined in (18). It stands to reason to first estimate the privacy loss  $\ell_{x, x'}(\cdot)$  by an empirical version  $\hat{\ell}_{x, x'}(\cdot)$ , which is then maximized to obtain an estimate for  $\epsilon_{x, x'}$ . Suppose that  $A$  is either discrete or continuous, s.t. a realization of  $A(x)$  has density  $f_x$ . By running that algorithm  $n$  times on data bases  $x$  and  $x'$  respectively, we can generate two independent samples of i.i.d observations  $X_1, \dots, X_n \sim f_x$  and  $Y_1, \dots, Y_n \sim f_{x'}$ . Recalling the definition of the loss function in (18) we can naturally define the *empirical loss function* as

$$\hat{\ell}_{x, x'}(t) := |\ln(\hat{f}_x(t)) - \ln(\hat{f}_{x'}(t))|, \quad (20)$$

where  $\hat{f}_x, \hat{f}_{x'}$  are density estimators for  $f_x, f_{x'}$ . In the case of continuous densities we can obtain such estimators via the TKDE algorithm (see Section II-B). For discrete densities we can simply use the relative amplitudes estimator which is described in the DDE (discrete density estimator) algorithm and mathematically defined as follows:

$$\hat{f}_x(t) := \frac{|\{X_i : X_i = t\}|}{n}.$$

We notice that by some basic calculations (union bound and Chebyshev inequality) it can be shown that for any discrete algorithm the uniform approximation rate

$$\sup_{t \in \mathcal{Y}} |\hat{f}_x(t) - f_x(t)| = \mathcal{O}_P(n^{-1/2}) \quad (21)$$

holds.

In the case of discrete algorithms with finite range we could now in principle just maximize the empirical loss  $\hat{\ell}_{x, x'}(\cdot)$

---

#### Algorithm 2 Discrete density estimator

---

**Input:**  $X = (X_1, \dots, X_n)$ : data sample,  $t$ : evaluation point

**Output:**  $\hat{f}(t)$ : density estimate at point  $t$

---

```

1: function DDE( $X, t$ )
2:  $out := 0$ 
3: for  $i = 1, 2, \dots, n$  do
4:   if  $X_i = t$  then
5:      $out = out + 1$ 
6:   end if
7: end for
8:  $out = out/n$ 
9: return  $out$ 
10: end function

```

---

to approximate  $\epsilon_{x, x'}$ . However, in the case of continuous algorithms and discrete algorithms with infinite output space, a straightforward maximization leads to computational instabilities, as the density estimates are potentially unreliable for extreme arguments, where (almost) no observations are sampled. We therefore restrict maximization to a closed, bounded set  $C \subset \mathcal{Y}$ , usually an interval (or hypercube in the multivariate case). Notice that

$$\epsilon_{x, x', C} := \sup_{t \in C} \ell_{x, x'}(t) \approx \sup_{t \in \mathcal{Y}} \ell_{x, x'}(t) = \epsilon_{x, x'} \quad (22)$$

in the sense that the difference between  $\epsilon_{x, x', C}$  and  $\epsilon_{x, x'}$  can be made arbitrarily small for sufficiently large  $C$ . For most standard algorithms even strict equality holds for some fixed  $C$  (this is true for all algorithms investigated in Section V). This is in particular true for discrete algorithms with finite range where we can always choose  $C = \mathcal{Y}$ .

We now state two regularity conditions that pertain to continuous algorithms and guarantee reliable inference:

- (C1) There exists a constant  $\beta \in (0, 1]$ , such that for all  $x$  the density  $f_x$  corresponding to  $A(x)$  is  $\beta$ -Hölder continuous.
- (C2) For any  $x, x'$  and any sequence  $(t_n)_{n \in \mathbb{N}}$  in  $C$ , which satisfies

$$\lim_{n \rightarrow \infty} \ell_{x, x'}(t_n) = \sup_{t \in C} \ell_{x, x'}(t),$$

it holds that  $(t_n)_{n \in \mathbb{N}}$  has a limit point in  $\arg \max_{t \in C} \ell_{x, x'}(t)$ .

We briefly comment on these assumptions: Condition (C1) demands that our algorithm is not only continuous in the sense that it has probability densities everywhere, but additionally that these satisfy a weak regularity condition of  $\beta$ -smoothness (see Section II-B). This guarantees reliable kernel density estimators and thus a good approximation of  $\ell_{x, x'}$  by  $\hat{\ell}_{x, x'}$ . Condition (C2) is a technical requirement that appears more complicated than it is: It prohibits that the maximum privacy violation (of  $A$  on  $C$ ) occurs in locations where both densities are 0, thus excluding pathological cases. Many continuous algorithms satisfy both of these conditions (among them all those discussed in this paper).



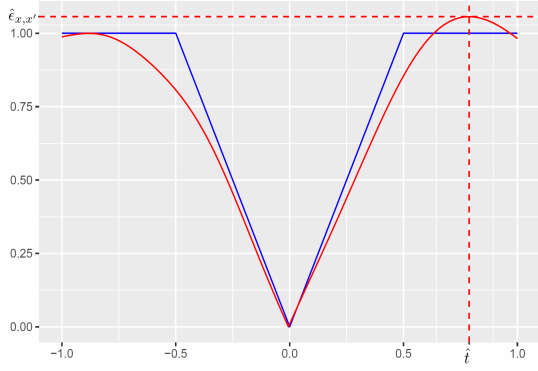


Fig. 3: Loss function  $\ell_{x,x'}$  (blue) and empirical loss  $\hat{\ell}_{x,x'}$  (red) for the Laplace algorithm. The vertical line indicates the location of the argmax  $\hat{t}$  and the horizontal line the maximum  $\epsilon_{x,x'}$  of the empirical loss function.

We now define the location  $\hat{t}$  of maximum privacy violation:

$$\hat{t} \in \arg \max_{t \in C} \hat{\ell}_{x,x'}(t). \quad (23)$$

In the following we demonstrate that the maximum of the empirical loss function, i.e.

$$\hat{\epsilon}_{x,x'} := \hat{\ell}_{x,x'}(\hat{t}) \quad (24)$$

is close to the the maximum of the true loss function.

To derive asymptotic convergence rates in the continuous case, the bandwidths  $h$  and  $h'$  of the truncated kernel density estimators  $\hat{f}_x$  and  $\hat{f}_{x'}$  in (20) have to be chosen appropriately. In addition, the floor  $\delta$  must not be smaller than the precision level of the density estimators (see Section II-B). We specify the proper choice of parameters in the following condition:

(C3) The parameters  $h, h'$  and  $\delta$  are adapted to  $n$  and satisfy

$$h, h' = \mathcal{O}(n^{-\frac{1}{2\beta+d}}) \quad \delta = \mathcal{O}(n^{-\frac{\beta}{2\beta+d}} \ln(n)).$$

**Proposition 1.** Suppose that  $C$  is a closed, bounded set and  $\epsilon_{x,x',C} \in (0, \infty)$ . If  $A$  is a discrete algorithm it follows that

$$|\hat{\epsilon}_{x,x'} - \epsilon_{x,x',C}| = \mathcal{O}_P(n^{-1/2})$$

and  $|\ell_{x,x'}(\hat{t}) - \epsilon_{x,x',C}| = \mathcal{O}_P(n^{-1/2}).$

If  $A$  is a continuous algorithm such that conditions (C1) – (C3) are satisfied, it follows that

$$|\hat{\epsilon}_{x,x'} - \epsilon_{x,x',C}| = \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+d}}\right)$$

and  $|\ell_{x,x'}(\hat{t}) - \epsilon_{x,x',C}| = \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+d}}\right).$

Furthermore if  $\epsilon_{x,x',C} \in \{0, \infty\}$  it holds that

$$\hat{\epsilon}_{x,x'} \rightarrow_P \epsilon_{x,x',C}$$

where " $\rightarrow_P$ " denotes convergence in probability (see Appendix A for a definition).

The first part of Proposition 1 suggests that the maximum privacy violation for  $x, x'$  is approximated by its empirical

counterpart at the same rate as the densities  $f_x, f_{x'}$  by their estimators (which is different in the discrete and the continuous case; see (21) and (11) respectively). This rate - specifically in the continuous case- should not be taken for granted: Admittedly, if the two continuous densities  $f_x, f_{x'}$  are bounded away from 0 on  $C$ , it is not difficult to show that

$$\sup_{t \in C} |\hat{\ell}_{x,x'}(t) - \ell_{x,x'}(t)| = \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+d}}\right),$$

which implies the Proposition. However if the densities are not bounded away from 0, it may not be true that  $\ell_{x,x'}$  is uniformly approximated by  $\hat{\ell}_{x,x'}$ . Still the approximation of the maxima holds and is not slowed down in this case (even though the mathematical proof gets substantially more involved). The second part of the Proposition states that  $\hat{t}$  is close to the argmax of  $\ell_{x,x'}$  in the sense that the true loss function evaluated at  $\hat{t}$  is close to its maximum on  $C$ . This fact will be used in the two subsequent sections, where we argue that a confidence interval for  $\ell_{x,x'}(\hat{t})$  automatically contains  $\epsilon_{x,x',C}$ .

We conclude this section by stating the DPL algorithm which, given  $x$  and  $x'$ , calculates the maximum empirical privacy loss, as well as  $\hat{t}$ . In the algorithm, the binary variable *discr* indicates whether a discrete (1) or continuous (0) setting is on hand and the set  $C$  encloses the area of interest.

---

### Algorithm 3 Data-specific privacy loss

---

**Input:** neighboring data bases  $x$  and  $x'$ , closed and bounded set  $C$ , sample size  $n$ , specification variable *discr*

**Output:** estimated loss  $\hat{\epsilon}_{x,x'}$ , location of loss  $\hat{t}$

```

1: function DPL( $x, x', n, C, discr$ )
2:   Generate  $X = (X_1, \dots, X_n)$  with  $X_i \sim A(x)$ 
3:   Generate  $Y = (Y_1, \dots, Y_n)$  with  $Y_i \sim A(x')$ 
4:   Set  $h, h'$  and  $\delta$  in accordance with (C3)
5:   Choose appropriate kernel  $K$ 
6:   if  $discr = 1$  then
7:      $\hat{f}_x(\cdot) = \text{DDE}(X, \cdot)$ 
8:      $\hat{f}_{x'}(\cdot) = \text{DDE}(Y, \cdot)$ 
9:   else
10:     $\hat{f}_x(\cdot) = \text{TKDE}(X, \cdot, h, K, \delta)$ 
11:     $\hat{f}_{x'}(\cdot) = \text{TKDE}(Y, \cdot, h', K, \delta)$ 
12:   end if
13:    $\ell_{x,x'}(\cdot) = |\ln(\hat{f}_x(\cdot)) - \ln(\hat{f}_{x'}(\cdot))|$ 
14:    $\hat{t} = \arg \max\{\ell_{x,x'}(t) : t \in C\}$ 
15:    $\hat{\epsilon}_{x,x'} = \hat{\ell}_{x,x'}(\hat{t})$ 
16:   return ( $\hat{t}, \hat{\epsilon}_{x,x'}$ )
17: end function
```

---

### B. Statistical bounds for pointwise privacy loss

In the previous section we have considered the problem of estimating data-specific privacy violations. We now move to the related topic of statistical inference in the sense of Section II-A: Finding a confidence interval for  $\epsilon_{x,x',C}$ .

More precisely we show in this section how to construct an asymptotic confidence interval for the pointwise privacy loss  $\ell_{x,x'}(t)$  for an arbitrary  $t \in C$ , which we apply later to the choice  $t = \hat{t}$  (recall that according to Proposition 1 we have  $\ell_{x,x'}(\hat{t}) \approx \epsilon_{x,x',C}$ ).



Suppose that  $\ell_{x,x'}(t) \in (0, \infty)$ . In this situation it can be shown by asymptotic normality of the density estimators and the Delta method (see [28]), that for all  $t \in \mathbb{R}$

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\frac{c_n}{\sigma}(\hat{\ell}_{x,x'}(t) - \ell_{x,x'}(t)) \leq t\right) = \Phi(t). \quad (25)$$

Here  $\Phi(\cdot)$  is the distribution function of a standard normal random variable and  $c_n = \sqrt{n}$  if the algorithm  $A$  is discrete and  $c_n = \sqrt{nh^d}$  if it is continuous. In the latter case  $h$  denotes the bandwidth of both  $\hat{f}_x, \hat{f}_{x'}$  and is assumed to be adapted to the sample size  $n$  as  $h = \mathcal{O}(n^{-\frac{1}{2\beta+d}-\gamma})$  for some  $\gamma > 0$ . This bandwidth is smaller than the one suggested in (C3) and leads to a slower uniform convergence of the corresponding density estimators (see Section II-B, (11)). Such a bandwidth choice, which makes the variance of the density estimator larger than its bias, is referred to as "undersmoothing". Undersmoothing is a standard tool in the statistical analysis of continuous densities, where the two tasks of estimation and inference require different degrees of smoothing (see [29] p.3999).

The variance  $\sigma^2$  on the right side of (25) can be expressed as follows

$$\sigma^2 := \begin{cases} \frac{1}{f_x(t)} + \frac{1}{f_{x'}(t)} - 2, & A \text{ discrete} \\ \int K^2(s) ds \left( \frac{1}{f_x(t)} + \frac{1}{f_{x'}(t)} \right), & A \text{ continuous.} \end{cases}$$

Note that  $\sigma^2$  is well defined in both cases (in particular in the discrete case  $1/f_x(t), 1/f_{x'}(t) > 1$ , s.t. the variance is indeed positive). Also notice that  $\sigma^2$  is unknown, but easy to estimate in practice, replacing the true densities by their estimators  $\hat{f}_x, \hat{f}_{x'}$ , which yields

$$\hat{\sigma}^2 := \begin{cases} \frac{1}{\hat{f}_x(t)} + \frac{1}{\hat{f}_{x'}(t)} - 2, & A \text{ discrete} \\ \int K^2(s) ds \left( \frac{1}{\hat{f}_x(t)} + \frac{1}{\hat{f}_{x'}(t)} \right), & A \text{ continuous.} \end{cases}$$

It is straightforward to show that  $\hat{\sigma}^2 = \sigma^2 + o_P(1)$ . We can now use this fact, together with the convergence in (25), to see that for any  $\alpha \in (0, 1)$

$$\begin{aligned} 1 - \alpha &\approx \mathbb{P}\left(\frac{c_n}{\hat{\sigma}}(\hat{\ell}_{x,x'}(t) - \ell_{x,x'}(t)) \leq \Phi^{-1}(1 - \alpha)\right) \quad (26) \\ &= \mathbb{P}\left(\hat{\ell}_{x,x'}(t) + \frac{\Phi^{-1}(\alpha)\hat{\sigma}}{c_n} \leq \ell_{x,x'}(t)\right). \end{aligned}$$

Here  $\Phi^{-1}$  denotes the quantile function of the standard normal distribution and we have used the identity  $\Phi^{-1}(1 - \alpha) = -\Phi^{-1}(\alpha)$ . The approximation of  $1 - \alpha$  by the probability gets more accurate as the sample size  $n$  increases and we see that

$$\hat{I}_\alpha := [\hat{\ell}_{x,x'}(t) + \hat{\sigma}c_n^{-1}\Phi^{-1}(\alpha), \infty)$$

is an asymptotic confidence interval for  $\ell_{x,x'}(t)$  (in the sense of Section II-A).

### C. A statistical procedure for the maximum privacy violation

Recall the definition of  $\epsilon_{x,x',C}$  in (22). In this section we construct an algorithm called MPL (Maximum Privacy Loss) whose output  $LB$  lower bounds the maximum of  $\epsilon_{x_1,x'_1,C}, \dots, \epsilon_{x_B,x'_B,C}$  with prescribed probability  $1 - \alpha$ . The choice of  $\alpha$  is determined by the user, but, guided by

common practice in hypothesis testing, we recommend  $\alpha \in \{0.1, 0.05, 0.01\}$ . By construction the inequality

$$\max\{\epsilon_{x_1,x'_1}, \dots, \epsilon_{x_B,x'_B}\} \geq \max\{\epsilon_{x_1,x'_1,C}, \dots, \epsilon_{x_B,x'_B,C}\}$$

holds and both sides are arbitrarily close for large enough  $C$ . Hence,  $LB$  will also constitute a tight lower bound for the maximum on the left and thus of the true privacy parameter  $\epsilon$  (see (19)).

We now outline the structure of the algorithm MPL, which calculates  $LB$  for a given set

$$\mathcal{X} = \{(x_1, x'_1), \dots, (x_B, x'_B)\}$$

of  $B$  adjacent pairs and is composed of two parts. The first part of the algorithm is dedicated to finding the pair of data bases  $(x_{max}, x'_{max}) \in \mathcal{X}$  along with the corresponding location  $\hat{t}_{max}$  that maximize the empirical privacy violation. For that purpose, MPL runs the DPL algorithm for each pair  $(x_b, x'_b)$  to approximate the data-specific privacy violation  $\epsilon_{x_b, x'_b}$  by an estimate  $\hat{\epsilon}_{x_b, x'_b}$ . Based on the empirical violations  $\hat{\epsilon}_{x_1, x'_1}, \dots, \hat{\epsilon}_{x_B, x'_B}$ , the pair of data bases  $(x_{max}, x'_{max})$  with the highest privacy loss is chosen. The location, where the empirical privacy loss  $\hat{\ell}_{x_{max}, x'_{max}}$  is maximized is called  $\hat{t}_{max}$  (which is an output of  $\text{DPL}(x_{max}, x'_{max})$ ). Structurally, this part of the algorithm resembles counterexample generation [14] and the tuple  $(\hat{\epsilon}_{x_{max}, x'_{max}}, x_{max}, x'_{max}, \hat{t}_{max})$  already yields useful information concerning the location and magnitude of the maximum privacy violation.

The second part of the MPL algorithm is designed to establish a confidence region for the privacy loss at  $(x_{max}, x'_{max}, \hat{t}_{max})$ . Notice that by construction  $\ell_{x_{max}, x'_{max}}(\hat{t}_{max}) \approx \epsilon_{x_{max}, x'_{max}}$  holds (see Theorem 1) and that therefore said confidence region captures the maximum privacy violation. The methods for deriving  $LB$  are borrowed from Section IV-B and are performed independently from the first part of the algorithm. MPL creates two fresh samples  $X_1^*, \dots, X_N^* \sim A(x_{max})$  and  $Y_1^*, \dots, Y_N^* \sim A(x'_{max})$  with sample size  $N > n$ . These are used to approximate the loss  $\ell_{x_{max}, x'_{max}}(\hat{t}_{max})$  by its empirical version  $\hat{\ell}_{x_{max}, x'_{max}}^*(\hat{t}_{max})$ . The density estimators  $\hat{f}_x^*, \hat{f}_{x'}^*$  underlying this empirical loss function are constructed with parameters  $h_{max}$  and  $\delta$  tailored to the construction of confidence intervals. This choice is expressed in the following condition:

$$(C4) \text{ Let } \nu \geq 0. \text{ With } N = \mathcal{O}(n^{1+\nu}) \text{ and } \gamma > \nu/((1+\nu)6) \text{ we choose } h_{max} = \mathcal{O}(N^{-\frac{1}{2\beta+d}-\gamma}) \text{ and } \delta = o(1).$$

As already indicated in Section IV-C, bandwidths for confidence intervals have to be chosen smaller than for estimation (this explains why  $\gamma > 0$ ). The trade-off between  $\gamma$  and  $\nu$  expresses that in the second part of the algorithm MPL, a larger sample size  $N$  compared to  $n$  requires more undersmoothing to control the bias. Yet as  $\nu$  is usually small in practice (in our experiments about 0.1), the undersmoothing requirement is rather weak. The fact that  $\delta$  can decay at any rate shows that

$\hat{t}_{max}$  (selected by truncated estimators in the first step) locates automatically in regions where the densities are not too close to 0 and thus a second truncation by  $\delta$  is not important. In applications one can simply put  $\delta = 0$  in this step.

Recalling Section IV-B and particularly (26), we can now give a confidence interval  $[LB, \infty)$  for  $\epsilon_{x_{max}, x'_{max}, C}$ , where the statistical lower bound  $LB$  is defined as follows:

$$LB := \hat{\ell}_{x, x'}^*(\hat{t}_{max}) + \frac{\Phi^{-1}(\alpha)\hat{\sigma}_N}{c_N}. \quad (27)$$

Here  $\Phi^{-1}$  is, again, the quantile function of the standard normal distribution and  $1 - \alpha$  is the confidence level. The normalizing constants  $c_N$  and  $\hat{\sigma}_N$  are described in Section IV-B. An outline of MPL is given in algorithm 4.

---

**Algorithm 4** Maximum Privacy Loss

---

**Input:** set of data pairs  $\mathcal{X}$ , sample sizes  $n$  and  $N$ , region of investigation  $C$ , specification variable  $discr$ , level  $\alpha$

**Output:** Statistical lower bound for privacy violation  $LB$

```

1: function MPL( $\mathcal{X}, n, N, C, discr, \alpha$ )
2:   for  $b = 1, \dots, B$  do
3:      $(\hat{t}_{x_b, x'_b}, \hat{\epsilon}_{x_b, x'_b}) = \text{DPL}(x_b, x'_b, n, C, discr)$ 
4:   end for
5:    $(x_{max}, x'_{max}) \in \arg \max\{\hat{\epsilon}_{x_b, x'_b} : (x_b, x'_b) \in \mathcal{X}\}$ 
6:    $\hat{t}_{max} := \hat{t}_{x_{max}, x'_{max}}$ 
7:   Generate  $X^* = (X_1^*, \dots, X_N^*)$  with  $X_i^* \sim A(x_{max})$ 
8:   Generate  $Y^* = (Y_1^*, \dots, Y_N^*)$  with  $Y_i^* \sim A(x'_{max})$ 
9:   Choose  $h_{max}, \delta$  in accordance with (C4)
10:  Choose appropriate kernel  $K$ 
11:  if  $discr = 1$  then
12:     $\hat{f}_{x_{max}}^*(\hat{t}_{max}) = \text{DDE}(X^*, \hat{t}_{max})$ 
13:     $\hat{f}_{x'_{max}}^*(\hat{t}_{max}) = \text{DDE}(Y^*, \hat{t}_{max})$ 
14:  else
15:     $\hat{f}_{x_{max}}^*(\hat{t}_{max}) = \text{TKDE}(X^*, \hat{t}_{max}, h_{max}, K, \delta)$ 
16:     $\hat{f}_{x'_{max}}^*(\hat{t}_{max}) = \text{TKDE}(Y^*, \hat{t}_{max}, h_{max}, K, \delta)$ 
17:  end if
18:   $\hat{\ell}_{x_{max}, x'_{max}}^*(\hat{t}_{max}) = |\ln(\hat{f}_{x_{max}}^*(\hat{t}_{max})) - \ln(\hat{f}_{x'_{max}}^*(\hat{t}_{max}))|$ 
19:  Calculate  $\hat{\sigma}^2$  and  $c_N$  based on  $X^*, Y^*$  and  $discr$ 
20:  Define  $LB := \hat{\ell}_{x_{max}, x'_{max}}^*(\hat{t}_{max}) + \frac{\Phi^{-1}(\alpha)\hat{\sigma}}{c_N}$ 
21:  return  $LB$ 
22: end function

```

---

The following theorem validates theoretically the lower bound  $LB$  produced by the MPL algorithm.

**Theorem 2.** Suppose that  $A$  is either a discrete algorithm or a continuous one such that conditions (C1)-(C4) are satisfied with regard to  $A$  and the MPL algorithm.

i) If

$$\epsilon_C^* := \max(\epsilon_{x_1, x'_1, C}, \dots, \epsilon_{x_B, x'_B, C}) \in (0, \infty)$$

it holds that

$$\lim_{n \rightarrow \infty} \mathbb{P}(LB \leq \epsilon_C^*) = 1 - \alpha. \quad (28)$$

ii) If  $\epsilon_C^* = \infty$ , then  $LB \rightarrow_P \infty$ . If  $\epsilon_C^* = 0$ , then  $LB \rightarrow_P 0$

The proof of the theorem is technical and therefore deferred to the Appendix.

## V. EXPERIMENTS

In this section we analyze the performance of our methodology by applying it to some standard algorithms in DP validation where the true privacy parameter is known. Our method is implemented in R and since it employs only common statistical tools, we can rely on standard packages and functions. For kernel density estimation we use the "kdensity" package, which also provides automatic bandwidth selection. In the following we give a short outline of the algorithms and experiment settings before discussing our empirical findings.

**Query model:** We briefly discuss the query model used in [14]. Many discrete algorithms do not operate on data bases  $x$  directly, but instead process query outputs  $q(x)$ . Thus, the search and selection of data bases  $x = (x(1), \dots, x(m))$  translates into a choice of query outputs

$$q = (q_1, \dots, q_d) = (q_1(x), \dots, q_d(x)).$$

Here counting queries, which check how many data points  $x(i)$  in  $x$  satisfy a given property, are of particular interest. A change in a single data point can impair the output of each counting query by at most 1. Hence, query answers on neighboring data bases are captured by vectors of natural numbers  $q, q'$  where  $q_i$  and  $q'_i$  can differ by at most 1. Simple query answers that are created following patterns displayed in Table I are sufficient to deduce the privacy parameter [14] and we will draw on these to evaluate discrete algorithms.

Pattern	Query $q$	Query $q'$
One Above	(1, 1, 1, 1, 1)	(2, 1, 1, 1, 1)
One Below	(1, 1, 1, 1, 1)	(0, 1, 1, 1, 1)
One Above Rest Below	(1, 1, 1, 1, 1)	(2, 0, 0, 0, 0)
One Below Rest Above	(1, 1, 1, 1, 1)	(0, 2, 2, 2, 2)
Half Half	(1, 1, 1, 1, 1)	(0, 0, 0, 1, 1)
All Above All Below	(1, 1, 1, 1, 1)	(2, 2, 2, 2, 2)
X Shape	(1, 1, 1, 0, 0)	(0, 0, 0, 1, 1)

TABLE I: Input patterns used in [14]

Similar to the discrete case, continuous algorithms are usually applied to aggregate statistics  $S$  of the data and not to the raw data itself. We therefore consider algorithmic inputs of the form  $s = S(x)$  and  $s' = S(x')$ , that lie in a continuous domain (in the following examples intervals and cubes).

**Algorithms:** We test our approach on 4 algorithms in total. The **Report Noisy Max** algorithm [23] publishes the query with the largest value within a vector of noisy query answers. More precisely, the index  $\arg \max\{q_i + L_i : 1 \leq i \leq d\}$  with  $L_i \sim \text{Lap}(\frac{2}{\epsilon})$  is calculated and returned (see [14], Algorithm 5). We implement Report Noisy Max and our procedure on vectors that entail 6 query answers and choose data bases  $q_b$  and  $q'_b$ ,  $b = 1, \dots, 10$ , that are similar to the patterns described in Table I.

Given a query vector  $q$  and a threshold  $T$ , the **Sparse Vector Technique (SVT)** goes through each query answer  $q_i$  and reports whether said query lies above or below  $T$  [23]. The maximum number of positive responses  $M$  is an adjustable feature of the algorithm that forces it to abort after

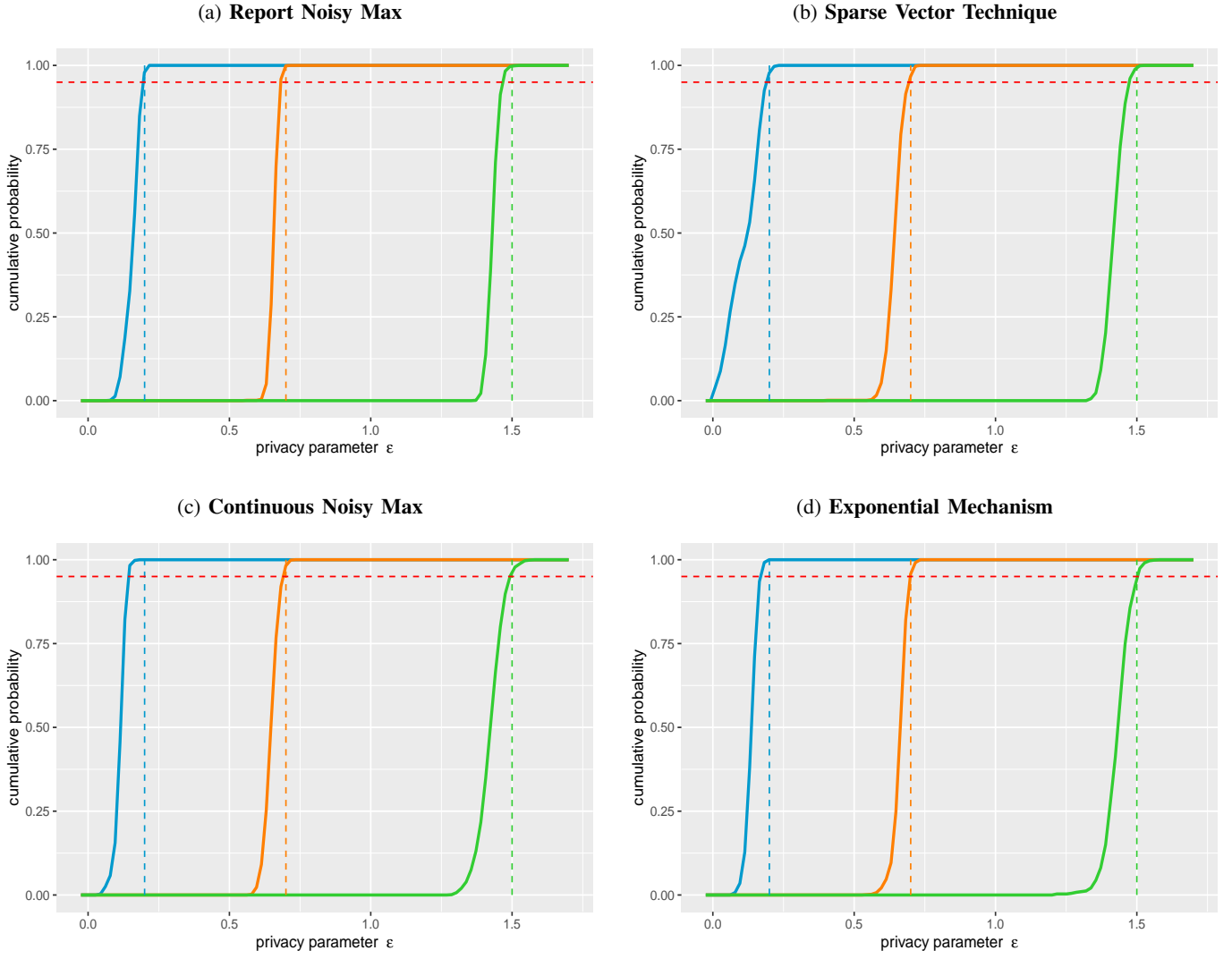


Fig. 4: Empirical distribution functions of the lower bound  $LB$ , generated by the MPL algorithm.

$M$  query answers above  $T$  have been reported. We consider query vectors  $q_b$  and  $q'_b$ ,  $b = 1, \dots, 10$ , with 10 entries and again similar to the patterns in Table I. This choice emulates the one in prior work (see [14], [16]) and we do the same for the tuning parameters with  $T = 1$  and  $M = 1$  [16].

The **continuous Noisy Max** algorithm (see Algorithm 7, [14]) has been discussed in Example 1. Here we use it to publish the maximum entry of a statistic  $s \in [0, 1]^k$ . We consider the case  $k = 3$  and input statistics  $s_b = (0, 0, 0)$  and  $s'_b = (b/10, b/10, b/10)$  for  $b = 1, \dots, 10$ . The set  $C$  in MPL is chosen as the symmetric interval  $[-1, 1]$ .

The **Exponential Mechanism** provides a general principle for the construction of private algorithms. We consider a version where we privatize real numbers from the interval  $[1, 2]$ , with non-negative outputs. More precisely, for a number  $s \in [1, 2]$  the output is sampled according to a continuous density proportional to  $\exp(-\lambda|s - t|)$  for  $t \geq 0$ . Here  $\lambda > 0$  is a parameter determining the privacy level. Recall that

this setup fits our (relaxed) notion of continuous algorithms discussed in Section III (continuous density on the half-line). It is well known that using this construction, the exponential mechanism affords (at least)  $2\lambda$ -DP. We can however employ Theorem 1 to derive the true privacy parameter  $\epsilon$  precisely:

$$\epsilon = \lambda + \ln(2 - \exp(-2\lambda)) - \ln(2 - \exp(-\lambda)).$$

Notice that  $\epsilon \approx 2\lambda$  for small  $\lambda$ . In the following simulations we consider input statistics  $s_b = 1$  and  $s'_b = 1 + b/10$  for  $b = 1, \dots, 10$  and choose  $C = [0, 2]$ .

**Experiment settings:** To study privacy violations we employ the MPL algorithm described in Section IV-C. The sample sizes in MPL are chosen as  $n = 20000$  and  $N = 50000$  for the discrete and continuous version of Noisy Max (which is substantially smaller than those used in [14]), as well as for the Exponential Mechanism. For SVT we use the larger samples  $n = 100000$  and  $N = 500000$  (matching the choice in [14]). This increase is necessary as SVT allows for extreme

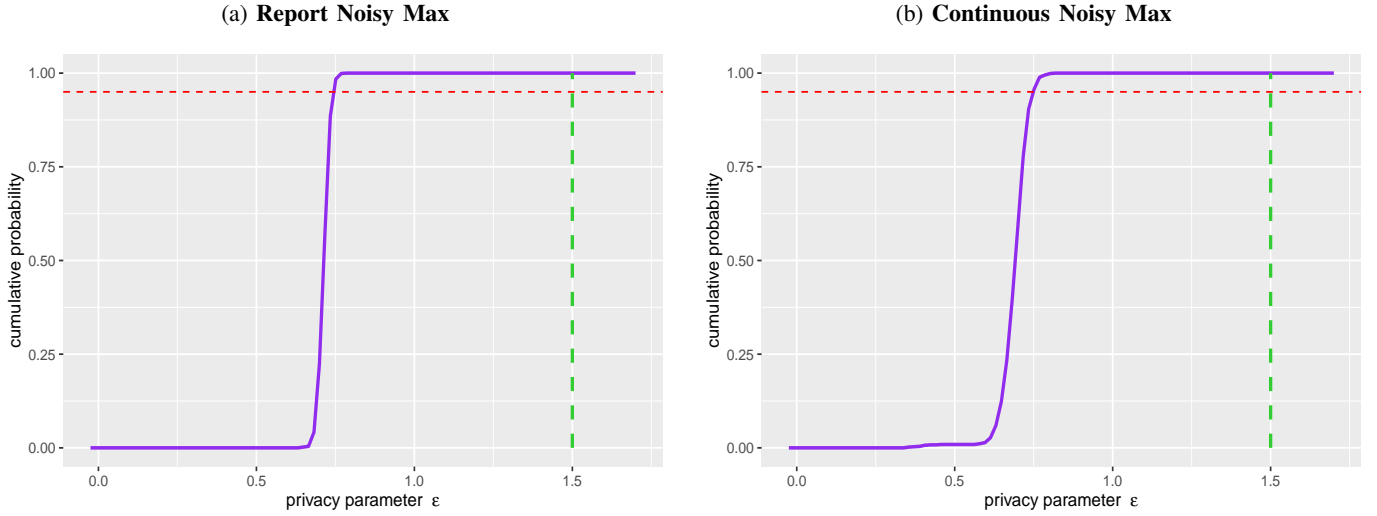


Fig. 5: Empirical distribution function of  $LB$  for fixed data bases.

events (with low probability) that cause instabilities. As a result insufficient sampling can lead to distorted estimation. For the continuous algorithms, the kernel in KDE is the Gaussian Kernel (described in Appendix B), the bandwidths in the first step of MPL are chosen by a pre-implemented selection rule in the "kdensity" package (both are the default options) and the floor  $\delta$  is 0.001.

We examine each algorithm in different privacy regimes  $\epsilon \in \{0.2, 0.7, 1.5\}$  [14] (we adjust the level of privacy by tuning the Laplace Noise or changing  $\lambda$  in the Exponential Mechanism).

**Results:** In order to evaluate MPL we consider the cumulative distribution function (cdf) of the lower bound  $LB$  defined in (27). In Figure 4 we display a panel where each plot corresponds to one algorithm under investigation and each curve to the empirical cdf for a different choice of  $\epsilon$  (each based on 1000 simulation runs). The vertical lines (in the same color as the corresponding cdfs) indicate the true privacy parameters and the horizontal, purple line the prescribed confidence level  $1 - \alpha$ , where we have chosen  $\alpha = 0.05$ . Notice that evaluated in the true privacy parameter  $\epsilon$ , the cdf describes the confidence level  $\mathbb{P}(LB \leq \epsilon)$ , which according to our theory should approximately equal  $1 - \alpha$  (see Theorem 2). In most scenarios we observe that the prescribed confidence level is indeed well approximated, while sometimes it is slightly too large (corresponding to small values of  $LB$ ). This tendency is inherent in the empirical study of DP and should not surprise us: To approximate  $\epsilon$ , one has to first select the right data pair out of  $B$  and then empirically maximize the privacy loss. Poor performance in either step biases estimates away from  $\epsilon$  towards smaller values - a trend that has been observed in other empirical studies (see e.g. [14], where the  $p$ -values are in each instance much higher than the prescribed level).

A second performance measure of our procedure is the ascent

of the cdf in a neighborhood of  $\epsilon$ : In all of our simulations (particularly for the two Noisy Max versions and the Exponential Mechanism) we observe a rapid increase close to  $\epsilon$ , suggesting that  $LB$  is a tight and reliable bound for  $\epsilon$  (even though our sample sizes are relatively small compared to [14]). In the case of SVT the ascent is slightly slower, which hints at higher variance in  $LB$  caused by smaller values of the discrete densities.

**The data-centric privacy level for fixed data bases:** As pointed out in Section IV, we can use the MPL algorithm to determine the data-centric privacy guarantee for select data bases defined in (5). We demonstrate this on both versions (discrete and continuous) of the Noisy Max algorithm.

Regarding the discrete case, suppose we have a data base  $x$  that, given 6 counting queries, evaluates to 0 for each query, that is  $q = q(x) = (0, 0, 0, 0, 0, 0)$ . Recalling our discussion of the query model, we know that any data base  $x'$  in the neighborhood of  $x$  evaluates to a binary vector  $q' \in \{0, 1\}^6$ . This means that the entire neighborhood of  $x$  can be exhausted by the collection of all such query pairs  $(q, q')$ . We set the privacy parameter  $\epsilon = 1.5$  and run the MPL algorithm for Report Noisy Max on that collection of query pairs 1000 times. In Figure 5 (left panel) we plot the empirical cdf of  $LB$  (purple), which exhibits a sharp rise, long before the true privacy parameter  $\epsilon$  (vertical green line). In view of our earlier results and given the exhaustive search of query pairs, we can be confident that the empirical cdf captures the data-centric privacy leakage  $\epsilon_x$ . The plot suggests that the data-centric privacy parameter is only about half the size of  $\epsilon$ , confirming that the amount of privacy afforded to this specific data base outstrips the worst case guarantee by far.

For the continuous case, we consider a data base  $x$  that produces the statistic  $s = S(x) = (1/2, 1/2, 1/2)$  and assume that  $S$  maps neighboring data bases  $x'$  anywhere on the unit cube  $[0, 1]^3$ . Let  $s' \in \{0, 1/2, 1\}^3$  (which forms an even grid of

27 points on the unit cube). We can run MPL on the collection of queries thus obtained. It can be shown by similar methods as employed in Example 1, that  $\epsilon_{x,x'} = \epsilon_x$  is attained for data bases  $x'$  with  $S(x') = s' = (0, 0, 0)$  or  $S(x') = s' = (1, 1, 1)$ , both of which are covered by our grid. Similarly as for the discrete case we observe that  $\epsilon_x$  is about half the size of  $\epsilon$  (see Figure 5 (b)). In conclusion, the amount of privacy ceded to our specific data bases  $x$  in both examples is about twice as high as the true privacy parameter suggests (i.e.  $\epsilon_x \approx \epsilon/2$ ).

**Estimation of data-specific privacy violations:** Up to this point we have focused on the lower bound  $LB$ , produced by the MPL algorithm. We now want to consider the estimation of data-specific privacy violations defined in (3), which is the key novelty of our local approach and, as an integral part of MPL, has an outsize effect on the quality of  $LB$ . We especially focus on the two continuous algorithms (Noisy Max and the Exponential Mechanism), where our estimator  $\hat{\epsilon}_{x,x'}$  differs most noticeably from prior approaches by virtue of kernel density estimation.

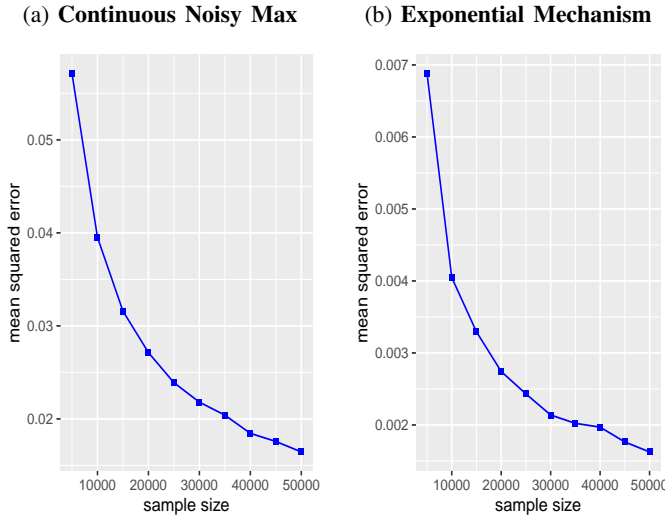


Fig. 6: Mean squared error  $\mathbb{E}(\hat{\epsilon}_{x,x'} - \epsilon_{x,x'})^2$  for different sample sizes  $n$  and  $\epsilon_{x,x'} = 1.5$ .

Regarding the Noisy Max algorithm, suppose we choose data bases  $x$  and  $x'$  that produce statistics  $s = S(x) = (0, 0, 0)$  and  $s' = S(x') = (1, 1, 1)$ , and similarly for the Exponential Mechanism data bases  $x$  and  $x'$  that result in  $s = 1$  and  $s' = 2$ . In both situations the choice of these data bases provokes a privacy violation  $\epsilon_{x,x'} = \epsilon$  that is equal to the true privacy parameter, which we fix at 1.5.

To study the quality of the estimator  $\hat{\epsilon}_{x,x'}$  based on  $n$  observations we consider the mean squared error  $\mathbb{E}(\hat{\epsilon}_{x,x'} - \epsilon_{x,x'})^2$  (approximated by 1000 simulation runs) for both algorithms. In Figure 6 we display the simulated errors for the two algorithms and different sizes of  $n$ . In both cases we observe for a sample size as moderate as 5000 only small estimation errors (less than 4% of the true  $\epsilon$  for Noisy Max and less than 0.5% for the Exponential Mechanism) and the errors are less

than half of this for  $n = 20000$  (which is used in our previous experiments). This shows that the strong performance of MPL can also be attributed to the precision of our estimators for the data-specific privacy violations.

## CONCLUSION

In this work we have discussed a way to assess privacy with statistical guarantees in a black box scenario. In contrast to prior works, our approach relies on a local conception of DP that facilitates the estimation and interpretation of privacy violations by circumventing the problem of event selection. Besides quantification of the global privacy parameter, our methods can be used for a more refined analysis, measuring the amount of privacy ceded to a specific data base. The findings of this analysis might not only help to understand existing algorithms better, but also aid the design of new privacy preserving mechanisms. This can, for instance, be algorithms that are tailored to provide greater privacy to data bases that require more protection.

## ACKNOWLEDGMENT

This work was partially funded by the DFG under Germany's Excellence Strategy - EXC 2092 CASA - 390781972.

# REFERENCES

- [1] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," ser. TCC'06. Berlin, Heidelberg: Springer-Verlag, 2006, p. 265–284.
- [2] J. Reed and B. C. Pierce, "Distance makes the types grow stronger: A calculus for differential privacy," ser. ICFP '10. New York, NY, USA: Association for Computing Machinery, 2010, p. 157–168.
- [3] M. Gaboardi, A. Haeberlen, J. Hsu, A. Narayan, and B. C. Pierce, "Linear dependent types for differential privacy," ser. POPL '13. New York, NY, USA: Association for Computing Machinery, 2013, p. 357–370.
- [4] G. Barthe, G. Danezis, B. Grégoire, C. Kunz, and S. Z. Béguelin, "Verified computational differential privacy with applications to smart metering," 2013 IEEE 26th Computer Security Foundations Symposium, pp. 287–301, 2013.
- [5] G. Barthe, M. Gaboardi, E. G. Arias, J. Hsu, C. Kunz, and P. Strub, "Proving differential privacy in hoare logic," Los Alamitos, CA, USA: IEEE Computer Society, jul 2014, pp. 411–424.
- [6] G. Barthe, N. Fong, M. Gaboardi, B. Grégoire, J. Hsu, and P.-Y. Strub, "Advanced probabilistic couplings for differential privacy," ser. CCS '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 55–67.
- [7] G. Barthe, M. Gaboardi, B. Grégoire, J. Hsu, and P.-Y. Strub, "Proving differential privacy via probabilistic couplings," ser. LICS '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 749–758.
- [8] A. Albarghouthi and J. Hsu, "Synthesizing coupling proofs of differential privacy," vol. 2, no. POPL, Dec. 2017.
- [9] D. Zhang and D. Kifer, "Lightdp: Towards automating differential privacy proofs," ser. POPL 2017. New York, NY, USA: Association for Computing Machinery, 2017, p. 888–901.
- [10] Y. Wang, Z. Ding, G. Wang, D. Kifer, and D. Zhang, "Proving differential privacy with shadow execution," ser. PLDI 2019. New York, NY, USA: Association for Computing Machinery, 2019, p. 655–669.
- [11] H. Zhang, E. Roth, A. Haeberlen, B. C. Pierce, and A. Roth, "Testing differential privacy with dual interpreters," vol. 4, no. OOPSLA, Nov. 2020.
- [12] G. Barthe, R. Chadha, V. Jagannath, A. P. Sistla, and M. Viswanathan, "Deciding differential privacy for programs with finite inputs and outputs," ser. LICS '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 141–154.
- [13] Y. Wang, Z. Ding, D. Kifer, and D. Zhang, "Checkdp: An automated and integrated approach for proving differential privacy or finding precise counterexamples," ser. CCS '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 919–938.
- [14] Z. Ding, Y. Wang, G. Wang, D. Zhang, and D. Kifer, "Detecting violations of differential privacy," ser. CCS '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 475–489.
- [15] B. Bichsel, T. Gehr, D. Drachler-Cohen, P. Tsankov, and M. Vechev, "Dp-finder: Finding differential privacy violations by sampling and optimization," ser. CCS '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 508–524.
- [16] B. Bichsel, S. Steffen, I. Bogunovich, and M. Vechev, "Dp-sniper: Black-box discovery of differential privacy violations using classifiers," 2021.
- [17] P. J. Bickel and K. A. Doksum, "Mathematical statistics." CRC Press, 2015.
- [18] A. van der Vaart and J. Wellner, "Weak convergence and empirical processes. with applications to statistics." Springer Series in Statistics., 1996.
- [19] D. W. S. Scott, "Multivariate density estimation: theory, practice, and visualization." Wiley, 1992.
- [20] A. Gramacki, Nonparametric Kernel Density Estimation and Its Computational Aspects. Cham, Switzerland: Springer International Publishing AG, 2018.
- [21] H. Jiang, "Uniform convergence rates for kernel density estimation," in Proceedings of the 34th International Conference on Machine Learning, ser. Proceedings of Machine Learning Research, D. Precup and Y. W. Teh, Eds., vol. 70. PMLR, 06–11 Aug 2017, pp. 1694–1703.
- [22] P. Kairouz, S. Oh, and P. Viswanath, "Extremal mechanisms for local differential privacy," J. Mach. Learn. Res., vol. 17, no. 1, p. 492–542, Jan. 2016.
- [23] C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," Found. Trends Theor. Comput. Sci., vol. 9, no. 3–4, p. 211–407, Aug. 2014.
- [24] M. Lyu, D. Su, and N. Li, "Understanding the sparse vector technique for differential privacy," Proc. VLDB Endow., vol. 10, no. 6, p. 637–648, Feb. 2017.
- [25] F. McSherry and K. Talwar, "Mechanism design via differential privacy," 11 2007, pp. 94–103.
- [26] A. W. Knap, "Basic real analysis." Birkhäuser, 2005.
- [27] W. Forst and D. Hoffmann, "Optimization—theory and practice." Springer-Verlag New York, 2010.
- [28] A. W. van der Vaart, "Asymptotic statistics." Cambridge University Press, 1998.
- [29] J. J. Heckman and E. Leamer, "Handbook of econometrics, volume 5." Elsevier Science B.V., 2001.
- [30] Y. M. Bishop, S. E. Fienberg, and P. W. Holland, "Discrete multivariate analysis: Theory and practice." Springer, 2007.
- [31] T. T. D. Team, R. Al-Rfou, G. Alain, A. Almahairi, C. Angermueller, D. Bahdanau, N. Ballas, F. Bastien, J. Bayer, A. Belikov, A. Belopolsky, Y. Bengio, A. Bergeron, J. Bergstra, V. Bisson, J. B. Snyder, N. Bouchard, N. Boulanger-Lewandowski, X. Bouthillier, A. de Brébisson, O. Breuleux, P.-L. Carrier, K. Cho, J. Chorowski, P. Christiano, T. Cooijmans, M.-A. Côté, M. Côté, A. Courville, Y. N. Dauphin, O. Delalleau, J. Demouth, G. Desjardins, S. Dieleman, L. Dinh, M. Ducoffe, V. Dumoulin, Samira Ebrahimi Kahou, Dumitru Erhan, Z. Fan, O. Firat, M. Germain, X. Glorot, I. Goodfellow, M. Graham, C. Gulcehre, P. Hamel, I. Harlouchet, J.-P. Heng, B. Hidasi, S. Honari, A. Jain, S. Jean, K. Jia, M. Korobov, V. Kulkarni, A. Lamb, P. Lamblin, E. Larsen, C. Laurent, S. Lee, S. Lefrançois, S. Lemieux, N. Léonard, Z. Lin, J. A. Livezey, C. Lorenz, J. Lowin, Q. Ma, P.-A. Manzagol, O. Mastropietro, R. T. McGibbon, R. Memisevic, B. van Merriënboer, V. Michalski, M. Mirza, A. Orlandi, C. Pal, R. Pascanu, M. Pezeshki, C. Raffel, D. Renshaw, M. Rocklin, A. Romero, M. Roth, P. Sadowski, J. Salvatier, F. Savard, J. Schlüter, J. Schulman, G. Schwartz, I. V. Serban, D. Serdyuk, Samira Shabanian, Étienne Simon, S. Spieckermann, S. R. Subramanyam, J. Synowski, J. Tanguay, G. van Tulder, J. Turian, S. Urban, P. Vincent, F. Visin, H. de Vries, David Warde-Farley, D. J. Webb, M. Willson, K. Xu, L. Xue, L. Yao, S. Zhang, and Y. Zhang, "Theano: A python framework for fast computation of mathematical expressions," ArXiv:1605.02688v1, 2016.
- [32] A. Ratner, D. Alistarh, G. Alonso, D. G. Andersen, P. Bailis, S. Bird, N. Carlini, B. Catanzaro, J. Chayes, E. Chung, B. Dally, J. Dean, I. S. Dhillon, A. Dimakis, P. Dubey, Charles Elkan, G. Fursin, G. R. Ganger, L. Getoor, P. B. Gibbons, G. A. Gibson, J. E. Gonzalez, J. Gottschlich, S. Han, K. Hazelwood, F. Huang, M. Jaggi, Kevin Jamieson, M. I. Jordan, G. Joshi, R. Khalaf, J. Knight, J. Koenig, T. Kraska, A. Kumar, A. Kyriillidis, A. Lakshmiratan, J. Li, S. Madden, H. B. McMahan, E. Meijer, I. Mitliagkas, R. Monga, D. Murray, K. Olukotun, D. Papailiopoulos, G. Pekhimenko, T. Rekatsinas, A. Rostamizadeh, C. Ré, C. D. Sa, H. Sedghi, Siddhartha Sen, V. Smith, A. Smola, D. Song, E. Sparks, I. Stoica, V. Sze, M. Udell, J. Vanschoren, S. Venkataraman, R. Vinayak, M. Weimer, A. G. Wilson, E. Xing, M. Zaharia, C. Zhang, and A. Talwalkar, "Mlsys: The new frontier of machine learning systems," ArXiv:1904.03257v3, 2019.
- [33] Y. Mirsky, Ambra Demontis, J. Kotak, R. Shankar, Deng Gelei, L. Yang, X. Zhang, W. Lee, Y. Elovici, and B. Biggio, "The threat of offensive ai to organizations," ArXiv:2106.15764v1, 2021.
- [34] Y. Zhang, M. Humbert, T. A. Rahman, C.-T. Li, J. Pang, and M. Backes, "Tagvisor: A privacy advisor for sharing hashtags," in Proc. of International World Wide Web Conference (WWW), 2018, pp. 287–296.
- [35] C.-T. Huang, "Ringcnn: Exploiting algebraically-sparse ring tensors for energy-efficient nonmem 5d i9 cnn-based computational imaging," ArXiv:2104.09056v1, 2021.
- [36] D. Mudigere, Y. Hao, J. Huang, A. Tulloch, S. Sridharan, X. Liu, M. Ozdal, J. Nie, J. Park, L. Luo, J. A. Yang, L. Gao, D. Ivchenko, A. Basant, Y. Hu, J. Yang, E. K. Ardestani, Xiaodong Wang, R. Komuravelli, C.-H. Chu, S. Yilmaz, H. Li, J. Qian, Z. Feng, Y. Ma, J. Yang, E. Wen, H. Li, L. Yang, C. Sun, W. Zhao, Dmitry Melts, K. Dhulipala, K. Kishore, T. Graf, A. Eisenman, K. K. Matam, A. Gangidi, G. J. Chen, M. Krishnan, A. Nayak, K. Nair, B. Muthiah, M. khorashadi, P. Bhattacharya, P. Lapukhov, M. Naumov, L. Qiao, M. Smelyanskiy, B. Jia, and V. Rao, "High-performance, distributed training of large-scale deep learning recommendation models," ArXiv:2104.05158v3, 2021.
- [37] M. Antonakakis, T. April, M. Bailey, M. Bernhard, E. Bursztein, Jaime Cochran, Z. Durumeric, J. A. Halderman, L. Invernizzi, M. Kallitsis, D. Kumar, C. Lever, Z. Ma, J. Mason, D. Menscher, C. Seaman,

- N. Sullivan, K. Thomas, and Y. Zhou, “Understanding the mirai botnet.” in *Proc. of USENIX Security Symposium*, 2017, pp. 1093–1110.
- [38] H. Xu, H. Wang, and Angelos Stavrou, “Privacy risk assessment on online photos.” in *Proc. of International Symposium on Research in Attacks, Intrusions and Defenses (RAID)*, 2015, pp. 427–447.
- [39] R. Thiago, R. Souza, L. Azevedo, E. Soares, R. Santos, W. Santos, M. D. Bayser, M. Cardoso, M. Moreno, and R. Cerqueira, “Managing stainer stingray 345m data lineage of o&g machine learning models: The sweet spot for shale use case,” *ArXiv:2003.04915v1*, 2020.
- [40] B. Delaware, S. Suriyakarn, C. Pit-Claudel, Qianchuan Ye, and A. Chlipala, “Narcissus: Deriving correct-by-construction decoders and encoders from binary formats,” *ArXiv:1803.04870v3*, 2018.
- [41] A. Mariakakis, S. Chen, B. Nguyen, K. Bray, M. Blank, J. Lester, Lauren Ryan, P. Johns, Gonzalo Ramos, and A. Roseway, “Project calico: Wearable chemical sensors for environmental monitoring,” *ArXiv:2006.15292v2*, 2020.
- [42] Yuanshun Yao, Bimal Viswanath, J. Cryan, H. Zheng, and B. Y. Zhao, “Automated crowdturfing attacks and defenses in online review systems.” in *Proc. of ACM Conference on Computer and Communications Security (CCS)*, 2017, pp. 1143–1158.
- [43] P. Parys, “Compositionality maxpool cho unperturb of the mso+u logic,” *ArXiv:2005.02384v1*, 2020.
- [44] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P. Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A. Nichol, A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr, J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati, K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and Wojciech Zaremba, “Evaluating large language models trained on code,” *ArXiv:2107.03374v2*, 2021.
- [45] A. Demontis, M. Melis, M. Pintor, M. Jagielski, B. Biggio, A. Oprea, Cristina Nita-Rotaru, and F. Roli, “Why do adversarial attacks transfer? explaining transferability of evasion and poisoning attacks.” in *Proc. of USENIX Security Symposium*, 2019, pp. 321–338.
- [46] M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter, “Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition.” in *Proc. of ACM Conference on Computer and Communications Security (CCS)*, 2016, pp. 1528–1540.
- [47] S. Akbarzadeh, S. Lee, and C.-T. Tan, “The spatial selective auditory attention of cochlear implant users xia indiscrimin nullifi in different conversational sound levels,” *ArXiv:2103.02703v1*, 2021.
- [48] J. Rahme, S. Jelassi, J. Bruna, and S. M. Weinberg, “A permutation-equivariant neural network architecture for auction design,” *ArXiv:2003.01497v2*, 2020.
- [49] Z. Huang, M. Liang, S. Liang, and W. He, “Altersgd: Finding flat minima for continual learning by alternative training,” *ArXiv:2107.05804v1*, 2021.
- [50] Jinyuan Jia, A. Salem, M. Backes, Y. Zhang, and Neil Zhenqiang Gong, “Memguard: Defending against black-box membership inference attacks via adversarial examples.” in *Proc. of ACM Conference on Computer and Communications Security (CCS)*, 2019, pp. 259–274.
- [51] S. R. Karingula and S. Lovett, “Codes over integers, and stingray tde nokia the singularity of random matrices with large entries,” *ArXiv:2010.12081v1*, 2020.

## APPENDIX A PROOFS AND TECHNICAL DETAILS

The appendix is dedicated to the mathematical details of our analysis: The definition of stochastic convergence, additional facts on the kernel  $K$  in KDE, as well as the proofs of Proposition 1 and Theorem 2.

### A. Stochastic Landau symbols and convergence in probability

Let  $(Z_n)_{n \in \mathbb{N}}$  be a sequence of random variables and  $(a_n)_{n \in \mathbb{N}}$  a sequence of positive, real numbers. We now say that

$Z_n = \mathcal{O}_P(a_n)$ , if for every  $\varepsilon > 0$  there exists a (sufficiently large)  $C > 0$  s.t.

$$\limsup_{n \rightarrow \infty} \mathbb{P}(|Z_n|/a_n \geq C) < \varepsilon.$$

Notice that analogous rules hold for the stochastic as for the deterministic Landau notation, such as  $\mathcal{O}_P(a_n) = a_n \mathcal{O}_P(1)$  or, for another positive sequence  $(b_n)_{n \in \mathbb{N}}$ , that  $\mathcal{O}_P(a_n) + \mathcal{O}_P(b_n) = \mathcal{O}_P(a_n + b_n)$ . Next we say that  $Z_n = o_P(a_n)$ , if for every (arbitrarily small)  $c > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|Z_n|/a_n \geq c) = 0.$$

Finally we say that for a constant  $a \in \mathbb{R}$  it holds that  $Z_n \rightarrow_P a$  if  $|Z_n - a| = o_P(1)$ . We say that  $Z_n \rightarrow_P \infty$ , if for any  $C > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}(Z_n \geq C) = 1.$$

For an extensive explanation of Landau symbols and convergence see [30]

### B. Kernel density estimation

Recall the definition of a kernel  $K$  as a continuous function  $K : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$  with  $\int_{\mathbb{R}^d} K(u) du = 1$ . In our discussion, we make the following two regularity assumptions, which are taken from [21] (Assumptions 2 and 3):

(K1)  $K$  satisfies *spherical symmetry*, i.e. there exists a non-increasing function  $k : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ , s.t.

$$K(u) = k(|u|), \quad \forall u \in \mathbb{R}^d.$$

(K2)  $k$  has *exponentially decaying tails*, i.e. there exist  $\rho, C_\rho, t_0$ , s.t.

$$k(t) \leq C_\rho \exp(-t^\rho), \quad \forall t > t_0.$$

A typical example of a kernel satisfying (K1) and (K2) is the *Gaussian kernel*, which corresponds to the density function of a standard normal and is given for  $d = 1$  as follows:

$$K(t) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{t^2}{2}\right).$$

We use this kernel in our experiments to study continuous algorithms.

### C. Proof of Proposition 1

We only show the Proposition for the case of a continuous algorithm  $A$  and only for  $d = 1$ . The discrete case works by similar, but simpler techniques and the continuous case for  $d > 1$  is a straightforward generalization. Furthermore we restrict ourselves to the case where  $\epsilon_{x,x',C} \in (0, \infty)$ . Proving consistency in the remaining cases  $\epsilon_{x,x',C} \in \{0, \infty\}$  is easier and therefore omitted.

We begin by defining two sets, that will be used extensively in our subsequent discussion: The argmax of the loss function

$$\mathcal{M} := \arg \max_{t \in C} \ell_{x,x'}(t).$$

and the closed  $\zeta$ -environment of  $\mathcal{M}$

$$U_\zeta(\mathcal{M}) := \{t \in C : \min_{t' \in \mathcal{M}} |t - t'| \leq \zeta\}.$$



Notice that  $\mathcal{M}$  is non-empty and closed. To see this, consider a sequence  $(t_n)_{n \in \mathbb{N}} \subset C$ , such that  $\ell_{x,x'}(t_n) \rightarrow \sup_{t \in C} \ell_{x,x'}(t)$ . Condition C5) implies that there exists a limit point in  $C$ , where the maximum is attained. In particular  $\mathcal{M} \neq \emptyset$ . Similarly we can show that  $\mathcal{M}$  is closed: If  $t$  is in the closure of  $\mathcal{M}$ , we can construct a sequence  $(t_n)_{n \in \mathbb{N}} \subset \mathcal{M}$  with  $t_n \rightarrow t$  and by Condition C5) it follows that  $t \in \mathcal{M}$ .

We now formulate an auxiliary result, that is the main stepping stone in the proof of Proposition 1.

**Lemma 1.** *Suppose that the Assumptions of Proposition 1 hold and  $\epsilon_{x,x',C} \in (0, \infty)$ . Then the following statements hold:*

i) *For any sufficiently small  $\zeta > 0$*

$$\sup_{t \in U_\zeta(\mathcal{M})} |\hat{\ell}_{x,x'}(t) - \ell_{x,x'}(t)| = \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+1}}\right).$$

ii) *There exists a  $\kappa = \kappa(\zeta) > 0$  s.t.*

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\sup_{t \notin U_\zeta(\mathcal{M})} \hat{\ell}_{x,x'}(t) > \sup_{t \in C} \ell_{x,x'}(t) - \kappa\right) = 0.$$

Let us verify that the Lemma indeed entails Proposition 1. We first show that for a small enough  $\zeta > 0$  it holds that

$$\lim_{n \rightarrow \infty} \mathbb{P}(\hat{t} \in U_\zeta(\mathcal{M})) = 1. \quad (29)$$

To see this we notice that according to Lemma 1, part ii) there exists a  $\kappa > 0$ , s.t.

$$\sup_{t \notin U_\zeta(\mathcal{M})} \hat{\ell}_{x,x'}(t) \leq \sup_{t \in \mathcal{M}} \ell_{x,x'}(t) - \kappa + o_P(1).$$

Here we have used  $\sup_{t \in \mathcal{M}} \ell_{x,x'}(t) = \sup_{t \in C} \ell_{x,x'}(t)$ . Combining this with part i) of the Lemma we have

$$\sup_{t \notin U_\zeta(\mathcal{M})} \hat{\ell}_{x,x'}(t) \leq \sup_{t \in \mathcal{M}} \hat{\ell}_{x,x'}(t) - \kappa + o_P(1).$$

As a consequence it holds with probability converging to 1, that  $\hat{\ell}_{x,x'}$  does not attain its maximum in  $C \setminus U_\zeta(\mathcal{M})$  and conversely that (29) holds. We now have for any  $t^* \in \mathcal{M}$ , that

$$\begin{aligned} |\hat{\ell}_{x,x'}(t^*) - \ell_{x,x'}(t^*)| &= \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+1}}\right) \\ |\hat{\ell}_{x,x'}(\hat{t}) - \ell_{x,x'}(\hat{t})| &= \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+1}}\right), \end{aligned} \quad (30)$$

where we have used part i) of the Lemma and for the second rate additionally (29). Recalling the definition  $\epsilon_{x,x',C} := \sup_{t \in C} \ell_{x,x'}(t)$ , (30) already entails part i) of Proposition 1. Part ii) of the proposition also follows from (30), as

$$\begin{aligned} |\epsilon_{x,x',C} - \ell_{x,x'}(\hat{t})| &= \ell_{x,x'}(t^*) - \ell_{x,x'}(\hat{t}) \\ &= [\ell_{x,x'}(t^*) - \hat{\ell}_{x,x'}(\hat{t})] + [\hat{\ell}_{x,x'}(\hat{t}) - \ell_{x,x'}(\hat{t})]. \end{aligned}$$

In the first step we have used that  $\epsilon_{x,x',C} = \ell_{x,x'}(t^*) \geq \ell_{x,x'}(\hat{t})$  because  $t^* \in \mathcal{M}$ . We can now treat the two terms on the right separately. The first term in the square brackets decays at the desired rate according to Proposition 1 part i) and the second part according to the second identity in (30). This shows Proposition 1 in the continuous case.

We now show that Lemma 1 holds. We begin with two technical observations: For any, sufficiently small  $\zeta > 0$  there exist positive constants  $\kappa, \rho > 0$ , such that simultaneously

$$\min_{t \in U_\zeta(\mathcal{M})} f_x(t) \wedge f_{x'}(t) \geq \rho > 0 \quad (31)$$

$$\sup_{t \in C \setminus U_\zeta(\mathcal{M})} \ell(x, x', t) < \sup_{t \in C} \ell(x, x', t) - \kappa, \quad (32)$$

where " $a \wedge b$ " denotes the minimum of two numbers  $a$  and  $b$ . We begin by proving (31): For all  $t \in \mathcal{M}$  it holds that  $f_x(t) \wedge f_{x'}(t) > 0$  (otherwise the assumption  $\sup_{t \in C} \ell_{x,x'}(t) \in (0, \infty)$  would be violated). Now  $f_x \wedge f_{x'}$  is a continuous function on the closed (thus compact) set  $\mathcal{M}$  and it therefore attains its (positive) minimum. Therefore for some  $\tilde{\rho} > 0$

$$\min_{t \in \mathcal{M}} f_x(t) \wedge f_{x'}(t) \geq \tilde{\rho}.$$

Now let  $t \in U_\zeta(\mathcal{M})$  and  $\tilde{t} \in \mathcal{M}$ , s.t.  $|t - \tilde{t}| \leq \zeta$ . According to Assumption C2) it holds that

$$\begin{aligned} &f_x(t) \wedge f_{x'}(t) \\ &\geq f_x(\tilde{t}) \wedge f_{x'}(\tilde{t}) - |f_x(t) \wedge f_{x'}(t) - f_x(\tilde{t}) \wedge f_{x'}(\tilde{t})| \\ &\geq \tilde{\rho} - a|t - \tilde{t}|^\beta \geq \tilde{\rho} - a\zeta^\beta. \end{aligned}$$

Here we have used for the second inequality that the minimum of two  $\beta$ -Hölder continuous functions is again  $\beta$ -Hölder (where we have called the constant  $a$ ). In the last step we have used that  $|t - \tilde{t}| \leq \zeta$ . It is now obvious that with sufficiently small  $\zeta$ , say  $\zeta < (\tilde{\rho}/(2a))^{1/\beta}$ , it follows (31) with  $\rho := \tilde{\rho}/2$ . Next we show (32). Suppose (32) was wrong. Then there must exist a sequence  $(t_n)_{n \in \mathbb{N}} \subset C \setminus U_\zeta(\mathcal{M})$  s.t.  $\ell_{x,x'}(t_n) \rightarrow \sup_{t \in C} \ell_{x,x'}(t)$ . According to C5) there exists a limit point  $t^*$ , where the maximum is attained. By definition  $t^* \in \mathcal{M}$ . This however is a contradiction to the fact, that  $|t_n - t^*| > \zeta$  for all  $n \in \mathbb{N}$ , showing (32). In the following we assume that  $\kappa, \rho, \zeta$  are chosen such that (31) and (32) hold.

We now prove part i) of Lemma 1. To show this, we first notice that for any fixed  $\rho' \in (0, \rho)$  it holds that

$$\lim_{n \rightarrow \infty} \mathbb{P}(\tilde{f}_x(t) \wedge \tilde{f}_{x'}(t) > \rho' : \forall t \in U_\zeta(\mathcal{M})) = 1, \quad (33)$$

where  $\tilde{f}_x(t), \tilde{f}_{x'}(t)$  are the KDEs defined in (8), Section II-B. (33) is a direct consequence of the uniform consistency of KDEs (see (11)). Now recall the definition of the truncated KDE  $\hat{f}_x := \tilde{f}_x \vee \delta$ . Since  $\delta \rightarrow 0$  and (33) holds, it follows for all  $t \in U_\zeta(\mathcal{M})$  simultaneously that  $\hat{f}_x(t) = \tilde{f}_x(t)$ , with probability converging to 1. Consequently the definition of the empirical loss in (20) implies that with probability converging to 1

$$\hat{\ell}_{x,x'}(t) = |\ln(\tilde{f}_x(t)) - \ln(\tilde{f}_{x'}(t))|, \quad \forall t \in U_\zeta(\mathcal{M}).$$

This means that to establish part i) of the Lemma, it suffices to show

$$\begin{aligned} &||\ln(\tilde{f}_x(t)) - \ln(\tilde{f}_{x'}(t))| \\ &- |\ln(f_x(t)) - \ln(f_{x'}(t))|| = \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+1}}\right). \end{aligned}$$

By the triangle inequality we can show the desired rate separately for  $|\ln(\tilde{f}_x(t)) - \ln(f_x(t))|$  and  $|\ln(\tilde{f}_{x'}(t)) - \ln(f_{x'}(t))|$ . We restrict ourselves to the first term (the second one follows by analogous arguments). By the mean value theorem it follows that

$$|\ln(\tilde{f}_x(t)) - \ln(f_x(t))| = \frac{|\tilde{f}_x(t) - f_x(t)|}{\xi(t)}, \quad (34)$$

where  $\xi(t)$  is a number between  $\tilde{f}_x(t)$  and  $f_x(t)$ . The numerator is of order

$$\sup_t |\tilde{f}_x(t) - f_x(t)| = \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+1}}\right), \quad (35)$$

where we have used the uniform approximation of kernel density estimators, from (11). The denominator is bounded away from 0, with probability converging to 1, as the bound

$$\xi(t) \geq f_x(t) - |\tilde{f}_x(t) - f_x(t)| \geq \rho - o_P(1), \quad (36)$$

holds uniformly for  $t \in U_\zeta(\mathcal{M})$ . Here we have used the lower bound (31) of the density  $f_x$  on  $U_\zeta(\mathcal{M})$ . Together (35) and (36) imply the desired rate for the right side of (34). By our above arguments this shows part i) of Lemma 1.

Next we prove part ii) of the Lemma 1. Let us therefore define pointwise in  $t$  the truncated density

$$f_x^{(\delta)}(t) := \begin{cases} f_x(t), & \text{if } \hat{f}_x(t) > \delta, \\ \delta, & \text{else} \end{cases}$$

and analogously the function  $f_{x'}^{(\delta)}$ . Therewith define the truncated loss

$$\ell_{x,x'}^{(\delta)}(t) := |\ln(f_x^{(\delta)}(t)) - \ln(f_{x'}^{(\delta)}(t))|. \quad (37)$$

By definition it holds for any  $\delta > 0$  and any  $t$ , that  $\ell_{x,x'}(t) \geq \ell_{x,x'}^{(\delta)}(t)$  ("=" if  $\hat{f}_x(t), \hat{f}_{x'}(t) > \delta$  and " $\geq$ " else). Now for any  $t \in C \setminus U_\zeta(\mathcal{M})$  we consider the following decomposition

$$\sup_{s \in C} \ell_{x,x'}(s) - \hat{\ell}_{x,x'}(t) = A_1 + A_2 + A_3 + A_4, \quad (38)$$

where

$$\begin{aligned} A_1 &:= \sup_{s \in C} \ell_{x,x'}(s) - \sup_{s \in C \setminus U_\zeta(\mathcal{M})} \ell_{x,x'}(s) \\ A_2 &:= \sup_{s \in C \setminus U_\zeta(\mathcal{M})} \ell_{x,x'}(s) - \sup_{s \in C \setminus U_\zeta(\mathcal{M})} \ell_{x,x'}^{(\delta)}(s) \\ A_3 &:= \sup_{s \in C \setminus U_\zeta(\mathcal{M})} \ell_{x,x'}^{(\delta)}(s) - \ell_{x,x'}^{(\delta)}(t) \\ A_4 &:= \ell_{x,x'}^{(\delta)}(t) - \hat{\ell}_{x,x'}(t). \end{aligned}$$

Now  $A_1 \geq \kappa$  holds according to (32). Furthermore  $A_2 \geq 0$  due to the inequality  $\ell_{x,x'}(s) \geq \ell_{x,x'}^{(\delta)}(s)$  and  $A_3 \geq 0$  because  $t \in C \setminus U_\zeta(\mathcal{M})$ . Finally we turn to  $A_4$  and show that it is uniformly in  $t$  of order  $o_P(1)$ . Using the triangle inequality, we can upper bound  $A_4$  by

$$|\ln(f_x^{(\delta)}(t)) - \ln(\hat{f}_x(t))| + |\ln(f_{x'}^{(\delta)}(t)) - \ln(\hat{f}_{x'}(t))|.$$

Both terms on the right can be treated analogously and so we focus on the first one. If  $\hat{f}_x(t) \leq \delta$  it is equal to 0 and thus

we consider the case where  $\hat{f}_x(t) > \delta$ . According to the mean value theorem

$$|\ln(f_x^{(\delta)}(t)) - \ln(\hat{f}_x(t))| = \frac{|f_x^{(\delta)}(t) - \hat{f}_x(t)|}{\xi'(t)}, \quad (39)$$

where  $\xi'(t)$  lies between  $f_x^{(\delta)}(t)$  and  $\hat{f}_x(t)$ . Just as before, the numerator is uniformly of order

$$\sup_{t \in C} |f_x(t) - \tilde{f}_x(t)| = \mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+1}}\right),$$

and the denominator is (asymptotically) bounded away from 0, as

$$\xi'(t) = \hat{f}_x(t) + \mathcal{O}_P(\sup_{t \in C} |f_x(t) - \tilde{f}_x(t)|) \geq \delta + o_P(\delta).$$

In both cases we have used that if  $\hat{f}_x(t) > \delta$  we have  $f_x^\delta(t) - \hat{f}_x(t) = f_x(t) - \tilde{f}_x(t)$ . Furthermore we have used for the denominator the approximation rate (11) and that according to C3)

$$\mathcal{O}_P\left(\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+1}}\right) = o_P(\delta).$$

These arguments imply that the right side of (39) is uniformly in  $t$  of order  $o_P(\delta)/[\delta + o_P(\delta)] = o_P(1)$ . By our above arguments we now have

$$A_1 + A_2 + A_3 + A_4 \geq \kappa + o_P(1),$$

which implies by (38) part ii) of Lemma 1 (if we replace  $\kappa$  by  $2\kappa$  in the above calculations).

#### D. Proof of Theorem 2

As for Proposition 1, we only show Theorem 2 for continuous algorithms and  $d = 1$ . Furthermore we confine ourselves to part i) of this theorem (as the convergence in part ii) follows by similar but simpler techniques). Finally, for clarity of presentation we will assume that there exists a unique  $b^* \in \{1, \dots, B\}$ , s.t.

$$\epsilon_{x_{b^*}, x'_{b^*}, C} = \max(\epsilon_{x_1, x'_1, C}, \dots, \epsilon_{x_B, x'_B, C}). \quad (40)$$

Recall that the algorithm MPL consists of two parts: In the first part the algorithm creates  $B$  pairs of samples with  $n$  elements each, to approximate  $\epsilon_{x_b, x'_b, C}$  by  $\hat{\epsilon}_{x_b, x'_b}$ . According to Proposition 1 these estimates are consistent and therefore with probability converging to 1 it holds that  $b_{max} = b^*$  (where  $b_{max}$  is an estimator, defined in Algorithm MPL and  $b^*$  in (40)). For simplicity we will subsequently assume that  $(x_{max}, x'_{max}) = (x_{b^*}, x'_{b^*})$  (formally we can do this by conditioning of the event  $\{b_{max} = b^*\}$ ). Next recall that from the first step of MPL we get empirical estimates  $\hat{\ell}_{x_{max}, x'_{max}}(\cdot)$  of the loss function and  $\hat{t}_{max}$  of the location of maximum privacy violation. These estimates are based on samples  $X_1, \dots, X_n \sim f_{x_{max}}$ ,  $Y_1, \dots, Y_n \sim f_{x'_{max}}$ . We will use these estimators in our subsequent discussion and it is important to keep them distinct from the randomness in the second part of the algorithm.

In the second step MPL generates fresh samples of size  $N$

$X_1^*, \dots, X_N^* \sim f_{x_{max}}, Y_1^*, \dots, Y_N^* \sim f_{x'_{max}}$ . The corresponding density estimates, generated by the TKDE algorithm are denoted by  $\hat{f}_{x_{max}}^*$  and  $\hat{f}_{x'_{max}}^*$  (to distinguish them from the estimators from the first step of the algorithm). Notice that these density estimators use the same kernel  $K$  as in the first step, but bandwidth  $h_{max}$  of a smaller size (the asymptotic rate is described in Condition C4)). Correspondingly we define the loss based on the  $*$ -samples

$$\hat{\ell}_{x_{max}, x'_{max}}^*(t) := |\hat{f}_{x_{max}}^*(t) - \hat{f}_{x'_{max}}^*(t)|.$$

We point out that by the choices of  $n, N$  and the bandwidth  $h_{max}$  (see Condition and C4)) it holds that

$$\sqrt{\ln(n)} n^{-\frac{\beta}{2\beta+1}} = o\left(\frac{1}{\sqrt{N h_{max}}}\right). \quad (41)$$

Now consider the decomposition

$$\sqrt{N h_{max}} \left( \sup_{t \in C} \ell_{x_{max}, x'_{max}}(t) - \hat{\ell}_{x_{max}, x'_{max}}^*(\hat{t}_{max}) \right) \quad (42)$$

$$=: B_1 + B_2 + B_3$$

where

$$B_1 := \sqrt{N h_{max}} \left( \sup_{t \in C} \ell_{x_{max}, x'_{max}}(t) - \hat{\ell}_{x_{max}, x'_{max}}^*(\hat{t}_{max}) \right)$$

$$B_2 := \sqrt{N h_{max}} \left( \hat{\ell}_{x_{max}, x'_{max}}(\hat{t}_{max}) - \ell_{x_{max}, x'_{max}}(\hat{t}_{max}) \right)$$

$$B_3 := \sqrt{N h_{max}} \left( \ell_{x_{max}, x'_{max}}(\hat{t}_{max}) - \hat{\ell}_{x_{max}, x'_{max}}^*(\hat{t}_{max}) \right).$$

According to Proposition 1 together with (41) it follows that  $B_1, B_2 = o_P(1)$ . Thus to show weak convergence of the left side of (42) (which is key to our asymptotic result) we can show weak convergence of  $B_3$ .

In order to study  $B_3$  we consider the more general object

$$G(t) := \sqrt{N h_{max}} \left( \ell_{x_{max}, x'_{max}}(t) - \ell_{x_{max}, x'_{max}}^*(t) \right)$$

which is defined for any  $t \in U_\zeta(\mathcal{M})$  (for some small enough, fixed  $\zeta$  s.t. (31) and (32) hold), where from now on

$$\mathcal{M} := \arg \max_{t \in C} \ell_{x_{max}, x'_{max}}(t).$$

We now notice that with probability converging to 1 it holds for all  $t \in U_\zeta(\mathcal{M})$  that

$$\begin{aligned} & \text{sign}(\ln(\hat{f}_{x_{max}}^*(t)) - \ln(\hat{f}_{x'_{max}}^*(t))) \\ &= \text{sign}(\ln(f_{x_{max}}(t)) - \ln(f_{x'_{max}}(t))). \end{aligned} \quad (43)$$

This follows because the density estimators are uniformly consistent (see Section II-B, equation (10)), together with boundedness away from 0 on  $U_\zeta(\mathcal{M})$  (see (31)).

For simplicity of presentation we subsequently assume that the signum on the right side of (43) is always 1. This means that with probability converging to 1

$$\begin{aligned} G(t) &= \sqrt{N h_{max}} \left( [\ln(\hat{f}_{x_{max}}^*(t)) - \ln(f_{x_{max}}(t))] \right. \\ &\quad \left. - [\ln(\hat{f}_{x'_{max}}^*(t)) - \ln(f_{x'_{max}}(t))] \right). \end{aligned}$$

By the mean value theorem we can transform the right side to

$$\sqrt{N h_{max}} \left( \frac{\hat{f}_{x_{max}}^*(t) - f_{x_{max}}(t)}{\xi_1(t)} - \frac{\hat{f}_{x'_{max}}^*(t) - f_{x'_{max}}(t)}{\xi_2(t)} \right).$$

Here  $\xi_1(t)$  lies between  $\hat{f}_{x_{max}}^*(t)$  and  $f_{x_{max}}(t)$ , and  $\xi_2(t)$  between  $\hat{f}_{x'_{max}}^*(t)$  and  $f_{x'_{max}}(t)$ . We now focus on the fraction of densities in  $x_{max}$  (the other one is analyzed step by step in the same fashion). Using (31) and the uniform consistency of the density estimates it is a simple calculation to show that

$$\frac{\hat{f}_{x_{max}}^*(t) - f_{x_{max}}(t)}{\xi_1(t)} = \frac{\hat{f}_{x_{max}}^*(t) - f_{x_{max}}(t)}{f_{x_{max}}(t)} + \text{Rem},$$

where  $\text{Rem}$  is a (negligible) remainder of size  $o_P(1/\sqrt{N h_{max}})$  (here we have applied the same techniques as in the discussion of (34)). We can rewrite the fraction on right side as follows

$$\begin{aligned} & \frac{\hat{f}_{x_{max}}^*(t) - f_{x_{max}}(t)}{f_{x_{max}}(t)} \\ &= \frac{1}{N f_{x_{max}}(t)} \sum_{i=1}^N \left[ h_{max}^{-1} K\left(\frac{t - X_i^*}{h_{max}}\right) - f_{x_{max}}(t) \right]. \end{aligned}$$

By standard arguments it is now possible to replace  $f_{x_{max}}(t)$  in the sum by  $\mathbb{E} h_{max}^{-1} K\left(\frac{t - X_i^*}{h_{max}}\right)$ , while only incurring a (uniformly in  $t$ ) negligible error. More precisely:

$$\begin{aligned} & \mathbb{E} h_{max}^{-1} K\left(\frac{t - X}{h_{max}}\right) = \int h_{max}^{-1} K\left(\frac{t - s}{h_{max}}\right) f_{x_{max}}(s) ds \\ &= \int K(s) f_{x_{max}}(s h_{max} + t) ds \\ &= f_{x_{max}}(t) + \int K(s) |f_{x_{max}}(s h_{max} + t) - f_{x_{max}}(t)| ds \\ &= f_{x_{max}}(t) + \mathcal{O}(|h_{max}|^\beta) \end{aligned}$$

Here we have used symmetry of the kernel (K1) in Appendix B) in the second and Hölder continuity of order  $\beta$  in the last equality (see Assumption C3); for a definition of Hölder continuity recall (9)). We also notice that  $\mathcal{O}(|h_{max}|^\beta) = o_P(1/\sqrt{N h_{max}})$ , which makes the remainder asymptotically negligible. By similar calculations we can show that

$$\begin{aligned} & \text{Var}\left(h_{max}^{-1} K\left(\frac{t - X_i^*}{h_{max}}\right)\right) \\ &= h_{max}^{-1} f_{x_{max}}(t) \int K^2(y) dy + \text{Rem}_2, \end{aligned} \quad (44)$$

where  $\text{Rem}_2$  is a remainder of negligible order. We can use the same considerations for  $f_{x'_{max}}$  to rewrite

$$G(t) = \frac{1}{\sqrt{N}} \sum_{i=1}^N \{Z_i(t) - \mathbb{E} Z_i(t)\} + o_P(1),$$

where

$$Z_i(t) = h_{max}^{-1/2} \left[ K\left(\frac{t - X_i^*}{h_{max}}\right) + K\left(\frac{t - Y_i^*}{h_{max}}\right) \right].$$

All variables  $Z_i$  are i.i.d. and according to (44) (and analogous calculations for  $f_{x'_{max}}$ ) asymptotically have variance

$$\sigma^2(t) := \int K^2(y) dy ([f_{x_{max}}(t)]^{-1} + [f_{x'_{max}}(t)]^{-1}),$$

Now define the estimator

$$\hat{\sigma}^2(t) := \int K^2(y) dy ([\hat{f}_{x_{max}}^*(t)]^{-1} + [\hat{f}_{x'_{max}}^*(t)]^{-1}),$$

which is identical to  $\hat{\sigma}^2$  in MPL for  $t = \hat{t}_{max}$ . By similar techniques as before we can show that  $\hat{\sigma}^2(t)$  is uniformly (for  $t \in U_\zeta(\mathcal{M})$ ) consistent for  $\sigma^2(t)$ . As a consequence we have  $G(t)/\hat{\sigma}(t) = S(t) + o_P(1)$ , where

$$S(t) := \frac{1}{\sqrt{N}} \sum_{i=1}^N \tilde{Z}_i(t) \quad (45)$$

and  $\tilde{Z}_i(t) := \{Z_i(t) - \mathbb{E}Z_i(t)\}/\sqrt{\text{Var}(Z_i)}$ . We can now prove the identity (28): First notice that

$$\begin{aligned} \mathbb{P}(LB \leq \epsilon_C^*) &= \mathbb{P}(LB \leq \epsilon_{x_{max}, x'_{max}, C}) \\ &= \mathbb{P}\left(\hat{\ell}^*(x, x', \hat{t}_{max}) + \frac{\Phi^{-1}(\alpha)\hat{\sigma}}{c_N} \leq \sup_{t \in C} \ell_{x_{max}, x'_{max}}(t)\right) \\ &= \mathbb{P}\left(\frac{c_N}{\hat{\sigma}}(\ell_{x_{max}, x'_{max}}(\hat{t}_{max}) - \hat{\ell}_{x_{max}, x'_{max}}^*(\hat{t}_{max})) \leq \Phi^{-1}(\alpha)\right) \\ &\quad + o(1). \end{aligned} \quad (46)$$

In the second equality we have used the decomposition (42), together with the fact, that  $B_1, B_2 = o_P(1)$ . We can plug in the definition of the process  $G$  into the probability on the right of (46), which gives us

$$\begin{aligned} &\mathbb{P}\left(\frac{G(\hat{t}_{max})}{\hat{\sigma}} \leq \Phi^{-1}(\alpha)\right) \\ &= \mathbb{P}\left(S(\hat{t}_{max}) \leq \Phi^{-1}(\alpha)\right) + o(1). \end{aligned} \quad (47)$$

Here we have used the definition of  $S$  in (45), as well as the (above mentioned) identity  $G(t)/\hat{\sigma} = S(t) + o_P(1)$ , which holds uniformly in  $t \in U_\zeta(\mathcal{M})$  (recall that  $\hat{t}_{max} \in \mathcal{M}$  with probability converging to 1 according to (29)). Moreover we have strictly speaking used that  $S$  has (asymptotically) a continuous distribution function (see below). Now recall that  $\hat{t}_{max}$  (which is based on the samples  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$  from the first step of the algorithm) is independent of all  $X_1^*, \dots, X_N^*, Y_1, \dots, Y_N^*$  (and so loosely speaking of the randomness in  $Z_i(\cdot)$ ). Thus we can express

$$\begin{aligned} &\mathbb{P}\left(S(\hat{t}_{max}) \leq \Phi^{-1}(\alpha)\right) \\ &= \int \mathbb{P}\left(S(t) \leq \Phi^{-1}(\alpha)\right) dP^{\hat{t}_{max}}(t), \end{aligned} \quad (48)$$

where  $P^{\hat{t}_{max}}$  is the image measure of  $\hat{t}_{max}$ . Again we use that asymptotically the probability that  $\hat{t}_{max} \notin U_\zeta(\mathcal{M})$  converges to 0 (see (29)). Now adding and subtracting  $\alpha$  yields

$$\begin{aligned} &\alpha + o(1) \\ &\quad + \int_{U_\zeta(\mathcal{M})} \mathbb{P}\left(S(t) \leq \Phi^{-1}(\alpha)\right) - \alpha \, dP^{\hat{t}_{max}}(t) \\ &= \alpha + o(1) \\ &\quad + \mathcal{O}\left(\sup_{t \in U_\zeta(\mathcal{M})} \left|\mathbb{P}\left(S(t) \leq \Phi^{-1}(\alpha)\right) - \Phi(\Phi^{-1}(\alpha))\right|\right). \end{aligned} \quad (49)$$

Given some fixed  $t$ , the sum  $S$  consists of i.i.d. random variables with unit variance and expectation 0. We can therefore apply the Berry-Esseen theorem to see that

$$\sup_{t \in U_\zeta(\mathcal{M})} \left|\mathbb{P}\left(S(t) \leq \Phi^{-1}(\alpha)\right) - \Phi(\Phi^{-1}(\alpha))\right| = o(1),$$

if we can show that (uniformly in  $t$ )

$$\frac{\mathbb{E}|\tilde{Z}_1(t) - \mathbb{E}\tilde{Z}_1(t)|^3}{\sqrt{N}} = o(1).$$

Similar calculations as before show that

$$\mathbb{E}|\tilde{Z}_1(t) - \mathbb{E}\tilde{Z}_1(t)|^3 = \mathcal{O}(h_{max}^{-1/2}),$$

which proves the approximation and thus entails that (49) equals  $\alpha + o(1)$ . This again implies by (46), (47) and (48) that the weak convergence in (28) holds and thus Theorem 1.3 part i).