

---

# Instance-optimality in differential privacy via approximate inverse sensitivity mechanisms

---

**Hilal Asi**  
Stanford University  
asi@stanford.edu

**John C. Duchi**  
Stanford University  
jduchi@stanford.edu

## Abstract

We study and provide instance-optimal algorithms in differential privacy by extending and approximating the inverse sensitivity mechanism. We provide two approximation frameworks, one which only requires knowledge of local sensitivities, and a gradient-based approximation for optimization problems, which are efficiently computable for a broad class of functions. We complement our analysis with instance-specific lower bounds for vector-valued functions, which demonstrate that our mechanisms are (nearly) instance-optimal under certain assumptions and that minimax lower bounds may not provide an accurate estimate of the hardness of a problem in general: our algorithms can significantly outperform minimax bounds for well behaved instances. Finally, we use our approximation framework to develop private mechanisms for unbounded-range mean estimation, principal component analysis, and linear regression. For PCA, our mechanisms give an efficient (pure) differentially private algorithm with near-optimal rates.

## 1 Introduction

We study the estimation of a function (statistic) of interest under differential privacy, where strong privacy protections usually decrease utility relative to non-private data analysis. In an effort to improve the utility of private algorithms, it is of utmost importance to design mechanisms that adapt to the hardness of the underlying data. Such mechanisms are of growing prevalence in the privacy literature, with prominent examples including the smooth sensitivity [25] and propose-test-release [12] frameworks.

To further investigate adaptivity to underlying instance, Asi and Duchi [4] recently study instance-optimal mechanisms—which, in a sense, achieve optimal utility for every possible data instance—in differentially private release of 1-dimensional quantities, moving beyond the more standard (worst case) minimax optimality. Inspired by classical statistical theory, Asi and Duchi develop local-minimax optimality and optimality against unbiased mechanisms, both of which aim to capture the hardness of the underlying data. By developing instance-specific lower bounds, they show that classical frameworks such as smooth sensitivity and propose-test-release may not be instance-optimal in general. To overcome this challenge, they investigate what they term the *inverse sensitivity mechanism*, showing it is instance-optimal for a wide range of functions.

Yet instance-optimality in private statistical estimation remains widely unexplored. First, the implementation of the inverse sensitivity mechanism requires a calculation of a particular sample distance (see Section 1.1.1), which may be intractable. Moreover, the current instance-optimality guarantees are not sharp for vector-valued functions. This is in part because the paper [4] tailors its instance-optimality notions for 1-dimensional functions by leveraging Stein’s “hardest one-dimensional alternative” approach to lower bounds [31, 9], which gives tight lower bounds for 1-dimensional functions but fails to yield correct bounds in higher dimensions.

To address these challenges, in this work we develop extensions and approximations to the inverse sensitivity mechanism with efficient implementations for a broad class of functions, which allows us to (for example) develop efficient algorithms for private PCA with near-optimal sample complexity. We also establish complementary instance-optimality results for vector-valued functions by proposing two approaches for instance-specific lower bounds: (i) a local minimax approach that measures the risk of an instance through the loss that an algorithm must incur on instances in a small neighborhood around it, and (ii) lower bounds against families of appropriately unbiased mechanisms, which includes many standard mechanisms. These instance-specific bounds suggest the limitations of more prevalent minimax (worst-case) bounds in privacy [19, 11]: they do not always give the correct limits on the performance of algorithms, and algorithms exist that achieve lower error on many instances.

## 1.1 Preliminaries

Given a function  $f : \mathcal{X}^n \rightarrow \mathcal{T}$  and instance  $\mathbf{x} \in \mathcal{X}^n$ , we wish to design differentially private mechanisms that accurately estimates the value  $f(\mathbf{x})$ . We usually take  $\mathcal{X}, \mathcal{T} \subset \mathbb{R}^d$  for a dimension  $d$ .

We begin by recalling the standard definition of differential privacy [16, 15]. We say that two instances  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$  are *neighboring* if they differ in at most one example, that is,  $d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq 1$ .

**Definition 1.1.** A randomized algorithm  $M : \mathcal{X}^n \rightarrow \mathcal{T}$  is  $(\varepsilon, \delta)$ -differentially private if for all neighboring datasets  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$  and all measurable  $S \subseteq \mathcal{T}$ ,

$$\mathbb{P}(M(\mathbf{x}) \in S) \leq e^\varepsilon \mathbb{P}(M(\mathbf{x}') \in S) + \delta.$$

If  $\delta = 0$ , then  $M$  is  $\varepsilon$ -differentially private.

Given a loss function  $L : \mathcal{T} \times \mathcal{T} \rightarrow \mathbb{R}_+$ , we quantify the utility of a mechanism  $M$  on instance  $\mathbf{x}$  through its expected loss  $\mathbb{E}[L(M(\mathbf{x}), f(\mathbf{x}))]$ . A mechanism is instance-optimal if it achieves the best utility for every instance. We formalize this through instance-specific lower bounds in Section 3.

For a function  $f : \mathcal{X}^n \rightarrow \mathbb{R}$ , the standard method to preserve privacy is the Laplace mechanism [16]. Defining the global sensitivity of  $f$  to be  $\text{GS}_f := \sup_{\mathbf{x}, \mathbf{x}' : d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq 1} |f(\mathbf{x}) - f(\mathbf{x}')|$ , it adds Laplace noise,  $M_{\text{Lap}}(\mathbf{x}) := f(\mathbf{x}) + \frac{\text{GS}_f}{\varepsilon} \text{Lap}(1)$ . This can be conservative, therefore Nissim et al. [25] consider the local sensitivity at instance  $\mathbf{x}$  at hand  $\text{LS}_f(\mathbf{x}) := \sup_{\mathbf{x}' : d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq 1} |f(\mathbf{x}) - f(\mathbf{x}')|$ . Directly using the local sensitivities may compromise privacy, hence the smooth sensitivity framework adds noise that is proportional to a smooth upper bound  $S^\beta(\mathbf{x})$  on the local sensitivity, that is,  $M_{\text{sm}}(\mathbf{x}) := f(\mathbf{x}) + \frac{2S^\beta(\mathbf{x})}{\varepsilon} Z$ , where  $Z$  is sampled from an admissible noise distribution and  $S^\beta(\mathbf{x})$  is the smooth sensitivity satisfying  $\text{LS}(\mathbf{x}) \leq S^\beta(\mathbf{x})$  and  $S^\beta(\mathbf{x}) \leq e^\beta S^\beta(\mathbf{x}')$  for neighboring instances  $\mathbf{x}, \mathbf{x}'$ , and  $\beta$  is chosen appropriately to guarantee the desired privacy level.

### 1.1.1 The inverse sensitivity mechanism

Our work builds on the inverse sensitivity mechanism [4], which we review. Key to the mechanism is the path-length (inverse sensitivity), which, for a target  $t$ , measures how many users we must change in  $\mathbf{x}$  to reach  $\mathbf{x}'$  with a target value  $t$ :

$$\text{len}_f(\mathbf{x}; t) := \inf_{\mathbf{x}'} \{d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \mid f(\mathbf{x}') = t\}. \quad (1)$$

The basic inverse sensitivity mechanism then instantiates the exponential mechanism [24] with the path-length function (1), yielding the density

$$\pi_{M_{\text{inv}}(\mathbf{x})}(t) = \frac{e^{-\text{len}_f(\mathbf{x}; t)\varepsilon/2}}{\int_{\mathcal{T}} e^{-\text{len}_f(\mathbf{x}; s)\varepsilon/2} ds}. \quad (\text{M.1})$$

A smoother variant of mechanism (M.1) is sometimes necessary to achieve instance-optimality, where one instead uses

$$\text{len}_f^\rho(\mathbf{x}; t) = \inf_{s \in \mathcal{T} : \|s - t\| \leq \rho} \text{len}_f(\mathbf{x}; s),$$

with a smoothing parameter  $\rho > 0$  [4]. Different variations of these mechanisms are instance-optimal for a range of real-valued functions. Yet while examples exist, it is often unclear how to compute the length (1).

Instance-specific bounds depend on the *modulus of continuity*, which (focusing in this work on a vector space  $\mathcal{T}$  with norm  $\|\cdot\|_p$ ) measures the sensitivity of a function when changing  $k$  users:

$$\omega_f^p(\mathbf{x}; k) = \sup_{\mathbf{x}' \in \mathcal{X}^n} \left\{ \|f(\mathbf{x}) - f(\mathbf{x}')\|_p : d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq k \right\}. \quad (2)$$

Instance-specific lower bounds show that the risk we expect for an  $\varepsilon$ -differentially private algorithm on instance  $\mathbf{x}$  is in general roughly  $\omega_f(\mathbf{x}; 1/\varepsilon)$  for 1-dimensional functions (with  $p = 1$ ) and loss  $L(s, t) = |s - t|$  [4]. Unfortunately, this is not tight for  $d$ -dimensional functions.

**Notation** We denote samples using bold symbols  $\mathbf{x} \in \mathcal{X}^n$  and individual examples using non-bold symbol  $x \in \mathcal{X}$ . We let  $d_{\text{ham}}(\mathbf{x}, \mathbf{x}')$  denote the Hamming distance of instance  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$ . The local sensitivity of  $f : \mathcal{X}^n \rightarrow \mathcal{T}$  at instance  $\mathbf{x}$  is  $\text{LS}_f^p(\mathbf{x}) = \omega_f^p(\mathbf{x}; 1)$ , and the global sensitivity of  $f$  is  $\text{GS}_f^p = \sup_{\mathbf{x} \in \mathcal{X}^n} \omega_f^p(\mathbf{x}; 1)$ . To facilitate notation, we sometimes remove the superscript  $p$  if  $p = 2$ . We let  $\text{diam}_p(\mathcal{T}) = \sup_{s, t \in \mathcal{T}} \|s - t\|_p$ , and  $\mathbb{B}_p^{d-1} = \{x \in \mathbb{R}^d : \|x\|_p \leq 1\}$  denote the  $\ell_p$ -ball.

## 1.2 Contributions

**Approximate inverse sensitivity mechanisms** We develop two approximation methods for the inverse sensitivity mechanism: (i) using local sensitivities in Section 2.1 and (ii) a gradient-based method for minimization problems in Section 2.2. These methods have efficient implementations for a wide range of problems and can outperform smooth sensitivity mechanisms for pure differential privacy. In contrast to Cauchy and Student’s T distributions used in such instantiations [7]—which have infinite first and third moment respectively—our mechanisms add noise with bounded  $p$ ’th moments for all finite  $p$ , resulting in improved high-probability bounds for utility analysis which is especially important for high dimensional functions as our examples demonstrate.

**Instance-optimality and lower bounds** We propose two notions of instance optimality for vector-valued functions and prove tight lower bounds for both notions in Section 3. Similarly to the 1-dimensional setting, our results give a characterization of the risk through the modulus of continuity. Combined with our instance-specific upper bounds, these bounds establish that approximate and exact inverse mechanisms are (nearly) instance-optimal for vector-valued functions under some assumptions.

**Applications** We study three problems that illustrate the methodological possibilities of the inverse sensitivity framework and its approximations in Section 4: mean estimation, PCA and linear regression. The utility improvements in these examples demonstrate the advantages of our mechanisms over standard frameworks and the importance of these notions of instance-optimality. Here we highlight the PCA example where smooth sensitivity algorithms require sample complexity (for dimension  $d$  and ignoring other parameters)  $O(d^{3/2})$  [18], whereas our mechanisms require  $O(d)$  samples, which is the optimal dependence on the dimension  $d$  according to PCA lower bounds [10, 21].

## 1.3 Related work

The most widely used frameworks for instance-dependent noise are smooth sensitivity [25] and propose-test-release [12]. The former adds noise that scaling with a smooth upper bound on the local sensitivity, and the latter adds noise scaling with a prespecified upper bound on the local sensitivity—whose validity the algorithm tests—in a neighborhood of the instance. Applications are numerous: Smith and Thakurta [30] develop an algorithm based on propose-test-release for high-dimensional regression problems, and Bun and Steinke [7] design noise distributions for smooth sensitivity and use them to estimate the mean of distributions with unbounded range. Other applications include principal component analysis [18], outlier analysis [26], and graph data [22, 32]. The inverse sensitivity framework is a distinct approach to instance-dependent noise that Asi and Duchi [4] investigate ([20, 29, 8] propose variants of the mechanism). Their results suggest that this framework, in contrast to smooth sensitivity and propose-test-release, is instance-optimal for a range of functions, and can have quadratically better sample complexity than smooth sensitivity mechanisms.

## 2 Approximate inverse sensitivity mechanisms

Having described the difficulty of sampling from the inverse sensitivity mechanism in general, in this section we develop two approximation frameworks that are applicable for a broader range of functions while maintaining some of the instance-optimality guarantees of the exact mechanism. First, in Section 2.1 we describe a method that uses the local sensitivities to approximate the path-length, and in Section 2.2 we describe an approximation for the specific setting of empirical risk minimization.

### 2.1 Approximation using local sensitivities

The mechanisms we develop in this section first construct an approximation  $\overline{\text{len}}_f(\mathbf{x}; t)$  for the path-length, then apply the exponential mechanism for a base measure  $\mu$  on  $\mathcal{T}$  with this approximation

$$\pi_{M_{\text{appr}}}(\mathbf{x})(t) = \frac{e^{-\overline{\text{len}}_f(\mathbf{x}; t)\varepsilon/2}}{\int_{\mathcal{T}} e^{-\overline{\text{len}}_f(\mathbf{x}; s)\varepsilon/2} d\mu(s)}. \quad (\text{M.2})$$

Our main tool for calculating  $\overline{\text{len}}_f(\mathbf{x}; t)$  are the local sensitivities of instance  $\mathbf{x}$  at distance  $\ell$

$$\text{LS}_\ell^p(\mathbf{x}) = \sup_{\mathbf{x}' : d_{\text{ham}}(\mathbf{x}, \mathbf{x}') = \ell} \text{LS}^p(\mathbf{x}').$$

This definition implies that changing  $k$  users can vary the function value by at most  $\sum_{\ell=1}^k \text{LS}_\ell^p(\mathbf{x})$ . As a consequence, we have the lower bound  $\text{len}_f(\mathbf{x}; t) \geq \min\{k : \sum_{i=1}^k \text{LS}_i^p(\mathbf{x}) \geq \|t - f(\mathbf{x})\|_p\}$ . Unfortunately, directly using this lower bound may result in mechanisms that are not private. The following theorem shows how to construct suitable approximations that preserve privacy.

**Theorem 1.** *Let  $f : \mathcal{X}^n \rightarrow \mathcal{T}$  and  $R_\ell : \mathcal{X}^n \rightarrow \mathbb{R}$  satisfy  $\text{LS}^p(\mathbf{x}) \leq R_1(\mathbf{x})$  and  $R_\ell(\mathbf{x}) \leq R_{\ell+1}(\mathbf{x}')$  for any neighboring instances  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$ . Then, using the approximation*

$$\overline{\text{len}}_f(\mathbf{x}; t) = \min \left\{ k : \sum_{i=1}^k R_i(\mathbf{x}) \geq \|t - f(\mathbf{x})\|_p \right\}, \quad (3)$$

*mechanism (M.2) is  $\varepsilon$ -differentially private.*

Algorithm 1 efficiently samples from the approximate inverse sensitivity mechanism for reasonable choices of  $p$ : the main bottleneck is step 3 but efficient algorithms exist for  $p \in \{1, 2\}$  using truncated Gamma distributions [23].

---

**Algorithm 1:** Sampling from approximate inverse sensitivity

---

**Input:**  $\mathbf{x} \in \mathcal{X}^n, p, \{R_i(\cdot)\}_{i=1}^n$

- 1 Denote  $S_k = \{t : \sum_{\ell=1}^{k-1} R_\ell(\mathbf{x}) \leq \|t\|_p \leq \sum_{\ell=1}^k R_\ell(\mathbf{x})\}$ ;
  - 2 Sample  $k \sim K$  from  $\mathbb{P}(K = k) \propto \text{Vol}(S_k) e^{-k\varepsilon/2}$  for  $1 \leq k \leq n$ ;
  - 3 Sample  $z \sim \text{Uni}(S_k)$ ;
  - 4 **return**  $f(\mathbf{x}) + z$
- 

Before proceeding to our utility analysis, we show that given the local sensitivities, we can always find an appropriate choice of  $R_i$  without calculating the smooth sensitivities.

**Proposition 2.1.** *Let  $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$  and assume  $\overline{\text{LS}}(\mathbf{x})$  is such that  $\text{LS}(\mathbf{x}) \leq \overline{\text{LS}}(\mathbf{x})$  for every  $\mathbf{x}$ . Then mechanism (M.2) using the approximation (3) with  $R_\ell(\mathbf{x}) = \sup_{\mathbf{x}' : d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq \ell} \overline{\text{LS}}(\mathbf{x}')$  is  $\varepsilon$ -differentially private.*

#### 2.1.1 Utility guarantees for vector-valued functions

In this section, we provide utility guarantees for the exact and approximate inverse sensitivity mechanisms for vector-valued functions. Combined with our lower bounds of Section 3, this establishes (near) instance optimality of these methods. Our guarantees hold with high probability, in contrast to those of the smooth sensitivity framework which uses distributions with heavy tails. We also show that our approximations can outperform smooth Laplace for real-valued functions.

We begin by analyzing the utility of the exact and approximate inverse sensitivity mechanisms.

**Theorem 2.** Let  $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$ ,  $\text{diam}_2(f(\mathcal{X}^n)) \leq D$ ,  $r > 0$  and  $1 \leq K \leq n$ . Then the (smooth) inverse sensitivity mechanism (M.1) with  $\rho = 1/n^r$  has

$$\mathbb{P}(\|M_{\text{inv}}(\mathbf{x}) - f(\mathbf{x})\|_2 \geq \omega_f(\mathbf{x}; K) + 1/n^r) \leq e^{-K\varepsilon/2}(n^r D)^d.$$

Moreover, if  $R_i(\mathbf{x}) \leq D$ , the approximate mechanism (M.2) using (3) with  $p = 2$  has

$$\mathbb{P}\left(\|M_{\text{appr}}(\mathbf{x}) - f(\mathbf{x})\|_2 \geq \sum_{i=1}^K R_i(\mathbf{x})\right) \leq e^{-K\varepsilon/2+1} \left(nD / \sum_{i=1}^K R_i(\mathbf{x})\right)^d.$$

We remark that using the smooth sensitivity framework to preserve pure differential privacy does not usually result in such high probability bounds due to using noise distributions with heavy tails such as Cauchy distribution [25]. Moreover, Theorem 2 implies that using  $k \approx \frac{Cd \log n}{\varepsilon}$  for large constant  $C$ , with high probability the inverse sensitivity mechanism roughly has

$$\|M_{\text{inv}}(\mathbf{x}) - f(\mathbf{x})\|_2 \leq O(\omega_f(\mathbf{x}; Cd \log n/\varepsilon)).$$

The approximate mechanism has similar loss whenever our approximate  $R_i$  are accurate such that  $\sum_{i=1}^K R_i(\mathbf{x}) = O(\omega_f(\mathbf{x}; K))$ . The lower bounds in Section 3 show this is (near) instance optimal.

We conclude this section with another choice of  $R_i$  that uses the smooth sensitivities instead. This guarantees that the approximate mechanism always outperforms the smooth Laplace mechanism [25].

**Proposition 2.2.** Let  $f : \mathcal{X}^n \rightarrow \mathbb{R}$ ,  $\varepsilon = O(1)$  and  $R_\ell(\mathbf{x}) = \sup_{\mathbf{x}': d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq \ell} S^\beta(\mathbf{x}')$ . Then mechanism (M.2) using (3) is  $\varepsilon$ -differentially private. If  $p = 1$  and  $\beta = \frac{\varepsilon}{8}$  then  $\mathbb{E}[\|M_{\text{appr}}(\mathbf{x}) - f(\mathbf{x})\|] \leq O\left(\frac{S^\beta(\mathbf{x})}{\varepsilon}\right)$ .

The smooth Laplace mechanism—which guarantees only approximate  $(\varepsilon, \delta > 0)$ -DP—has loss  $O\left(\frac{S^\beta(\mathbf{x})}{\varepsilon}\right)$  with a much smaller  $\beta = \frac{\varepsilon}{2 \log 2/\delta}$ , which can be  $\log 1/\delta$  worse in some settings.

## 2.2 Gradient-based approximations for empirical risk minimization

In this section, we describe our second approximation which applies to empirical risk minimization problems. Given data points  $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$  and  $L$ -Lipschitz loss function  $\ell(\theta; x_i)$  for  $\theta \in \Theta$ , we wish to solve the following minimization problem

$$\hat{\theta}_n = \underset{\theta \in \Theta}{\text{argmin}} L_n(\theta; \mathbf{x}, \mathbf{y}) := \frac{1}{n} \sum_{i=1}^n \ell(\theta; x_i, y_i).$$

It is possible to calculate the path-length using gradients for robust regression [4]. Here, we use similar techniques to approximate the inverse sensitivity mechanism in general settings. As  $\ell$  is  $L$ -Lipschitz, we need to change  $\text{len}(\mathbf{x}, \mathbf{y}; \theta) \geq \frac{n}{L} \|\nabla L_n(\theta; \mathbf{x}, \mathbf{y})\|_2$  users to make  $\theta$  a minimizer with  $\nabla L_n(\theta; \mathbf{x}', \mathbf{y}') = 0$ . The *gradient mechanism* uses this approximation of  $\text{len}$ , resulting in the density

$$\pi_{\text{Grad}}(\theta \mid \mathbf{x}, \mathbf{y}) \propto e^{-\frac{n\varepsilon}{2L} \|\nabla L_n(\theta; \mathbf{x}, \mathbf{y})\|_2}. \quad (4)$$

Sampling from this distribution can be hard in general, but we show an efficient implementation for linear regression in Section 4.3. For general twice differentiable functions, we propose an efficient heuristic of the gradient mechanism based on Taylor's expansion which gives  $\nabla L_n(\theta; \mathbf{x}, \mathbf{y}) \approx \nabla^2 L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y})(\theta - \hat{\theta}_n)$ . Letting  $\text{GS}_{\text{Hess}}$  denote the global sensitivity of  $\|\nabla^2 L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y})(\theta - \hat{\theta}_n)\|_2$ , we define the following *Hessian-based mechanism* for  $\theta \in \Theta$

$$\pi_{\text{Hess}}(\theta \mid \mathbf{x}, \mathbf{y}) \propto e^{-\frac{\varepsilon}{2\text{GS}_{\text{Hess}}} \|\nabla^2 L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y})(\theta - \hat{\theta}_n)\|_2}. \quad (5)$$

The main advantage of the Hessian mechanism (5) is that now we can design efficient and simple sampling procedures (see Section 4.3). It also provides an accurate approximation of the gradient mechanism with good utility whenever  $\nabla^2 L_n$  is  $H$ -Lipschitz with small  $H$ .

The privacy of these mechanisms follow immediately from the privacy of the exponential mechanism. For utility, we start with the following lemma which upper bound  $\text{GS}_{\text{Hess}}$ .

**Lemma 2.1.** Assume  $\ell(\cdot; x_i, y_i)$  is  $L$ -Lipschitz,  $\nabla^2 L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y}) \succeq \lambda I$ ,  $\hat{\theta}_n \in \text{int } \Theta$ , and  $\nabla^2 L_n(\cdot; \mathbf{x}, \mathbf{y})$  is  $H$ -Lipschitz. If  $H \leq \lambda = O(1)$  and  $n \geq 4L \text{diam}_2(\Theta) + 1$  then  $\text{GS}_{\text{Hess}} \leq O\left(\frac{L}{n}\right)$ .

We are now ready to analyze the utility of the Hessian-based mechanism.

**Proposition 2.3.** Let the set of instances  $(x_i, y_i)_{i=1}^n$  satisfy the assumptions of Lemma 2.1. Then the Hessian mechanism (5) is  $\varepsilon$ -DP. If  $\inf_{\theta \in \text{bd } \Theta} \|\theta - \hat{\theta}_n\|_2 \geq \Omega\left(\frac{dL^2}{n^2\varepsilon^2} \text{tr}(\nabla^2 L(\hat{\theta}_n; \mathbf{x}, \mathbf{y})^{-2})\right)$ , then

$$\mathbb{E}_{\theta \sim \pi_{\text{Hess}}(\cdot | \mathbf{x}, \mathbf{y})} \left[ \|\theta - \hat{\theta}_n\|_2^2 \right] \leq O\left( \frac{dL^2 \text{tr}(\nabla^2 L(\hat{\theta}_n; \mathbf{x}, \mathbf{y})^{-2})}{n^2\varepsilon^2} \right).$$

When  $H = 0$ , the gradient (4) and Hessian mechanisms (5) are identical and the gradient mechanism has the same utility. In Section 4.3, we use these mechanisms for solving regression problems and show the significant advantages of instance-specific bounds over standard minimax bounds.

Finally, we remark that our gradient-based approximations of the inverse sensitivity mechanism are closely related to the K-norm mechanism [27] where the authors use gradient norms as a score function for the exponential mechanism. However, their work only provides asymptotic utility analyses without finite-sample guarantees, and they propose an approximate implementation of their mechanisms using an MCMC procedure without providing privacy guarantees for the implementation.

### 3 Instance-specific lower bounds for vector-valued functions

Given a function  $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$ , in this section we prove instance-specific lower bounds on the loss that any private mechanism must incur. Unfortunately, the instance-specific notions in [4] were tailored for 1-dimensional functions, hence do not result in satisfactory lower bounds in our setting. To this end, we propose extensions that result in tight bounds. The first notion gives lower bounds by restricting to families of appropriately unbiased mechanisms. The second is a local-minimax approach that measures the performance in a small neighborhood around a given instance.

We begin with our optimality notion for unbiased mechanisms, which we define now.

**Definition 3.1.** We say that a randomized algorithm  $M$  is  $\|\cdot\|$ -unbiased if for any  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$  and  $\rho > 0$ ,

$$\mathbb{P}(\|M(\mathbf{x}) - f(\mathbf{x})\| \leq \rho) \geq \mathbb{P}(\|M(\mathbf{x}) - f(\mathbf{x}')\| \leq \rho).$$

Definition 3.1 says that when applying an unbiased mechanism  $M$  on instance  $\mathbf{x}$ , the output is more likely to be in a ball around the correct value  $f(\mathbf{x})$  rather than  $f(\mathbf{x}')$  for some other instance  $\mathbf{x}'$ . Anderson's theorem [2] implies that the Laplace mechanism, Gaussian mechanism, their smooth sensitivity instantiations, the approximate inverse sensitivity mechanism, and any instantiation of the exponential mechanism with a concave score function are  $\|\cdot\|$ -unbiased.

Our lower bounds require a growth condition on the set of values at distance at most  $k$ ,  $W_f(\mathbf{x}; k) = \{f(\mathbf{x}') : d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq k\}$ . We have the following instance-specific lower bound for unbiased mechanisms.

**Theorem 3.** Let  $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$  and assume  $W_f(\mathbf{x}; k) \supseteq c \cdot \omega_f(\mathbf{x}; k) \cdot \mathbb{B}_2^{d-1}$  for  $c > 0$ . If  $M$  is  $\varepsilon$ -DP,

$$\sup_{\mathbf{x} \in \mathcal{X}^n} \mathbb{E} [\|M(\mathbf{x}) - f(\mathbf{x})\|_2] \geq \frac{c}{8} \sup_{\mathbf{x} \in \mathcal{X}^n} \max_{1 \leq k \leq n} e^{-k\varepsilon/d} \omega_f(\mathbf{x}; k).$$

Moreover, if  $M$  is  $\|\cdot\|_2$ -unbiased, then for any  $\mathbf{x} \in \mathcal{X}^n$ ,

$$\mathbb{E} [\|M(\mathbf{x}) - f(\mathbf{x})\|_2] \geq \frac{c}{8} \max_{1 \leq k \leq n} e^{-2k\varepsilon/d} \omega_f(\mathbf{x}; k).$$

Theorem 3 suggests that worst-case lower bounds may be too pessimistic: while the minimax risk is roughly  $\sup_{\mathbf{x} \in \mathcal{X}^n} \omega_f(\mathbf{x}; d/\varepsilon)$ , we may hope to achieve a better risk for instance  $\mathbf{x}$ , that is,  $\omega_f(\mathbf{x}; d/\varepsilon)$ .

Now we define the local-minimax risk for an instance  $\mathbf{x}$  following similar ideas in statistical theory [cf. 33, Ch. 8]. Let  $\mathcal{M}_\varepsilon$  be the family of  $\varepsilon$ -differentially private mechanisms. For a radius  $r$ , we define the local-minimax risk of  $\mathbf{x}$  to be the worst-case risk in a small neighborhood around  $\mathbf{x}$ , that is,

$$\mathcal{R}(\mathbf{x}; r) := \inf_{M \in \mathcal{M}_\varepsilon} \sup_{\mathbf{x}' : d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq r} \mathbb{E} [\|M(\mathbf{x}') - f(\mathbf{x}')\|_2]. \quad (6)$$



The choice of  $r$  in this definition is important to exclude trivial mechanisms such as  $M(\mathbf{x}') = f(\mathbf{x})$ ; briefly, we choose the smallest radius that excludes such mechanisms, which is in this case  $r = \Theta(d/\varepsilon)$  (see Appendix C.2 for more details about this definition).

We have the following lower bound for the local-minimax risk.

**Theorem 4.** *Let  $f : \mathcal{X}^n \rightarrow \mathbb{R}^d$  and assume  $W_f(\mathbf{x}; k) \supseteq c \cdot \omega_f(\mathbf{x}; k) \cdot \mathbb{B}_2^{d-1}$  for  $c > 0$ . Then for any  $\mathbf{x} \in \mathcal{X}^n$ ,  $\mathcal{R}(\mathbf{x}; d/\varepsilon) \geq \Omega(\omega_f(\mathbf{x}; d/\varepsilon))$ .*

Similarly to our lower bounds for unbiased mechanisms, Theorem 4 shows that any mechanism must incur local-minimax risk roughly  $\omega_f(\mathbf{x}; d/\varepsilon)$  for instance  $\mathbf{x}$ . The upper bounds of Theorem 2 show that the exact inverse mechanism achieves this loss for every instance up to logarithmic factors, as well as the approximate version if the approximations  $R_i$  are accurate.

## 4 Applications

We investigate three examples that demonstrate different advantages and applications of the exact, approximate, and gradient inverse sensitivity mechanisms. Our examples include (i) mean estimation with unbounded range, (ii) principal component analysis and (iii) linear regression, and show that our techniques yield private algorithms with better noise distributions resulting in improved utility, which in some cases can significantly outperform existing minimax-optimal algorithms.

### 4.1 Unbounded-range mean estimation

Given  $x_i \stackrel{\text{iid}}{\sim} P$  with unbounded range, our goal is to privately estimate the mean  $\mu = \mathbb{E}_{x \sim P}[x]$ . The difficulty here is that the empirical mean has infinite global and even local sensitivity, leading Bun and Steinke [7] to use the trimmed mean which calculates the mean after removing the smallest and largest  $m$  samples. Letting  $x_{(1)} \leq \dots \leq x_{(n)}$  denote the order statistics, the trimmed mean is

$$\text{trim}_m(\mathbf{x}) = \frac{x_{(m+1)} + x_{(m+2)} + \dots + x_{(n-m)}}{n - 2m}. \quad (7)$$

This is useful as it leads to small local sensitivity under distributional assumptions. Bun and Steinke [7] use the smooth sensitivity to estimate the trimmed mean, resulting in strong utility but only with the weaker concentrated differential privacy [14, 6]. To preserve pure differential privacy, they use Student's T distribution which has infinite third moments and consequently heavy tails.

We use the exact inverse mechanism to estimate the mean with strong utility. This algorithm has finite  $p$ 'th moment for any finite  $p$ , therefore yields tight confidence intervals. We assume  $\mu \in [a, b]$  and let  $[c]_{[a,b]}$  denote projection to  $[a, b]$ . The following lemma enables exact calculation of the path-length.

**Lemma 4.1.** *Let  $f(\mathbf{x}) = [\text{trim}_m(\mathbf{x})]_{[a,b]}$ . Then, for any  $t \in [a, b]$ , if  $t \geq \text{trim}_m(\mathbf{x})$ , we have  $\text{len}_f(\mathbf{x}; t) = \min\{k : k \leq m, t - f(\mathbf{x}) \leq \frac{1}{n-2m} \sum_{i=1}^k (x_{(n-m+i)} - x_{(m+i)})\} \cup \{m+1\}$ .*

The calculation for  $t < \text{trim}_m(\mathbf{x})$  is similar and we present it in Appendix D. Using Lemma 4.1, we can efficiently sample (Algorithm 4 in Appendix D.1 which runs in  $O(n \log n)$  time) from the inverse sensitivity mechanism. To analyze the performance of this algorithm, we assume  $P$  is  $\sigma$ -subgaussian while noting that these results can be extended to other settings in [7]. The following proposition upper bounds the error of our algorithm, which resembles the bounds that the algorithms of [7] achieve with the weaker concentrated differential privacy.

**Proposition 4.1.** *Let  $a, b \in \mathbb{R}$  and  $x_i \stackrel{\text{iid}}{\sim} P$  where  $P$  is  $\sigma$ -subgaussian with mean  $\mu \in [a, b]$ . If  $P$  is symmetric about its mean and  $n \geq \frac{12 \log(n(b-a)/\sigma^2)}{\varepsilon}$ , the inverse sensitivity mechanism (Algorithm 4) with  $\rho = \frac{\sigma^2}{n^2}$  is  $\varepsilon$ -differentially private and has*

$$\mathbb{E}[(\hat{x} - \mu)^2] \leq \frac{\sigma^2}{n} + \frac{\sigma^2}{n^2} \cdot O\left(\frac{\log((b-a)/\sigma)}{\varepsilon} + \frac{\log n}{\varepsilon^2}\right).$$

## 4.2 Principal component analysis

In this section, we apply our approximations to calculate a rank  $k$  approximation of a matrix. Given  $x_1, \dots, x_n \in \mathbb{B}_2^{d-1}$  with covariance  $\Sigma(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ , we wish to find  $V \in \mathbb{R}^{d \times k}$  that solves

$$\hat{V}(\mathbf{x}) = \operatorname{argmax}_{V: V^T V = I_k} F(V) := \operatorname{tr}(V^T \Sigma(\mathbf{x}) V) \quad (8)$$

Gonen and Gilad-Bachrach [18] design  $\varepsilon$ -DP algorithms based on the smooth sensitivity framework with suboptimal error of roughly  $\frac{d^{3/2}}{n \operatorname{GAP}(\mathbf{x}) \varepsilon}$ . We show that our algorithms achieve a near-optimal rate  $\frac{d}{n \operatorname{GAP}(\mathbf{x}) \varepsilon}$ . Though there exist algorithms that achieve this rate using the exponential mechanism [10, 21], these algorithms require sampling from complex distributions and the only implementation with theoretical runtime analysis requires  $O(d^6)$  time. In contrast, given the eigenvectors, our algorithm (Algorithm 2) returns a private version of the leading eigenvector in time  $O(n + d)$  with high probability.

We only consider  $k = 1$  as extensions to larger  $k$  are straightforward using QR factorization [18]. Our algorithm builds on techniques from [18] and the approximate inverse mechanism. It requires a non-private PCA algorithm  $\mathcal{A}_1$  that calculates the first eigenvector  $\hat{v} \in \mathbb{R}^d$  (which maximizes (8)) and the gap between the two largest eigenvalues  $\operatorname{GAP}(\mathbf{x}) := \lambda_1(\Sigma(\mathbf{x})) - \lambda_2(\Sigma(\mathbf{x}))$ . Then it randomly flips the sign of  $\hat{v}$  as  $-\hat{v}$  is also a solution, and adds noise using the approximate inverse sensitivity. Algorithm 2 describes our private PCA procedure. Given the output of the non-private PCA algorithm, the main computational difficulty in Algorithm 2 is step 3 which requires sampling from the noise distribution of Algorithm 1. Using the rejection-sampling algorithms of Laud et al. [23] for sampling from truncated Gamma distributions (which has constant success probability in our setting), we can efficiently sample from Algorithm 1 in time  $O(n + d)$  with high probability.

---

### Algorithm 2: Private PCA using approximate inverse sensitivity

---

**Input:**  $\mathbf{x}$

- 1 Calculate  $\hat{v} = \mathcal{A}_1(\mathbf{x})$ ,  $\operatorname{GAP} = \operatorname{GAP}(\mathbf{x})$ ;
  - 2 Set  $\bar{v} = B\hat{v}$  for  $B \sim \operatorname{Uni}\{-1, +1\}$ ;
  - 3 Sample  $z$  from (M.2) using Algorithm 1 with  $p = 2$ ,  $R_i = \min(C_{\text{pca}}/(n \operatorname{GAP} - 2k), \sqrt{2})$ ;
  - 4 **return**  $v_{\text{out}} = \frac{\bar{v} + z}{\|\bar{v} + z\|_2}$ ;
- 

Following [18], we define the local sensitivity ( $k = 1$ ) while taking into consideration the vector sign

$$\operatorname{LS}(\mathbf{x}) = \sup_{\mathbf{x}': d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq 1} \min(\|\hat{V}(\mathbf{x}) - \hat{V}(\mathbf{x}')\|_2, \|\hat{V}(\mathbf{x}) + \hat{V}(\mathbf{x}')\|_2).$$

We build on the following key lemma that bounds the local sensitivity.

**Lemma 4.2** ([18], Theorem 5, Lemma 11). *If  $\operatorname{GAP}(\mathbf{x}) > 0$  then there is a universal constant  $C_{\text{pca}} < \infty$  such that  $\operatorname{LS}(\mathbf{x}) \leq \min(\frac{C_{\text{pca}}}{n \operatorname{GAP}(\mathbf{x})}, \sqrt{2})$ . Moreover,  $|\operatorname{GAP}(\mathbf{x}) - \operatorname{GAP}(\mathbf{x}')| \leq 2d_{\text{ham}}(\mathbf{x}, \mathbf{x}')/n$ .*

Using this bound and the guarantees of our approximate mechanism, we get the following proposition.

**Proposition 4.2.** *Assume  $n \geq 1/C_{\text{pca}}$ ,  $\beta > 0$  and  $\Omega(\frac{d}{\operatorname{GAP}(\mathbf{x}) \varepsilon}) \leq \frac{n}{\log n / \beta}$ . Algorithm 2 is  $\varepsilon$ -differentially private and with probability  $1 - \beta$ ,*

$$|F(v_{\text{out}}) - F(\hat{v})| \leq O\left(\frac{d \log n / \beta}{n \operatorname{GAP}(\mathbf{x}) \varepsilon} + \frac{1}{n^4}\right).$$

## 4.3 Linear regression

For our final example, we investigate the setting of linear regression where we have data points  $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ . Our goal here is to find  $\theta \in \Theta$  that minimizes

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} L_n(\theta; \mathbf{x}, \mathbf{y}) := \frac{1}{2n} \sum_{i=1}^n (\langle \theta, x_i \rangle - y_i)^2.$$

We let  $X \in \mathbb{R}^{n \times d}$  has  $x_i$  in the  $i$ 'th row,  $\Sigma_n = \frac{1}{n} X^T X$ , and  $\mathbf{y} \in \mathbb{R}^d$  denote the vector of  $y_i$ .



The gradient mechanism (4) have an efficient implementation in this setting. Assuming  $\Sigma_n \succ 0$ , we get that  $\nabla L_n(\theta) = \Sigma_n(\theta - \bar{\theta})$  where  $\bar{\theta} = \frac{1}{n}\Sigma_n^{-1}X^T\mathbf{y}$ . The gradient mechanism has density

$$\pi(\bar{\theta} + \Delta \mid \mathbf{x}) \propto e^{-\frac{n\varepsilon}{2L}\|\Sigma_n\Delta\|_2},$$

for  $\bar{\theta} + \Delta \in \Theta$  (and Lipschitz constant  $L$ ) which Algorithm 3 samples from (see Appendix F.1). The main difficulty in Algorithm 3 is calculating the non-private estimator  $\bar{\theta}$  in step 1 while the remaining steps (for private noise addition) only require sampling from simple distributions and matrix-vector products.

---

**Algorithm 3:** Gradient mechanism for linear regression

---

- Input:**  $(x_i, y_i)_{i=1}^n$
- 1 Calculate  $\Sigma_n = \frac{1}{n}X^T X$ ,  $\bar{\theta} = \frac{1}{n}\Sigma_n^{-1}X^T\mathbf{y}$ ;
  - 2 Sample  $R \sim \text{Gamma}(d, 1)$ ,  $U \sim \text{Uni}(\mathbb{S}^{d-1})$ ;
  - 3 Set  $\theta_{\text{out}} = \bar{\theta} + \frac{2L}{n\varepsilon}\Sigma_n^{-1} \cdot R \cdot U$ ;
  - 4 **if**  $\theta_{\text{out}} \notin \Theta$  **then** go to 3;
  - 5 **return**  $\theta_{\text{out}}$ ;
- 

The following proposition states the utility and privacy guarantees of Algorithm 3.

**Proposition 4.3.** *For the set of instances  $(x_i, y_i)_{i=1}^n$  with  $\Sigma_n(\mathbf{x}) \succ 0$  and Lipschitz constant  $L$ , Algorithm 3 is  $\varepsilon$ -differentially private. Moreover, if  $\inf_{\theta \in \text{bd } \Theta} \|\theta - \hat{\theta}_n\|_2 \geq \Omega\left(\frac{dL^2}{n^2\varepsilon^2} \text{tr}(\Sigma_n(\mathbf{x})^{-2})\right)$ ,*

$$\mathbb{E}\left[L_n(\theta_{\text{out}}; \mathbf{x}, \mathbf{y}) - L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y})\right] \leq O\left(\frac{dL^2 \text{tr}(\Sigma_n(\mathbf{x})^{-1})}{n^2\varepsilon^2}\right).$$

To appreciate the instance-specific upper bounds of Proposition 4.3, recall that existing private algorithms for empirical risk minimization of  $L$ -Lipschitz and  $\lambda$ -strongly convex functions achieve excess loss  $\mathbb{E}[L_n(\theta) - L_n(\hat{\theta}_n)] = O\left(\frac{d^2L^2}{n^2\varepsilon^2\lambda}\right)$  which is minimax optimal in some regimes [5]. In contrast, for natural instances where  $\Sigma_n(\mathbf{x})^{-1}$  has polynomially decaying eigenvalues  $\lambda_j = j^{-\alpha}$  for  $\alpha \in (0, 1]$ , Proposition 4.3 implies that Algorithm 3 achieves excess loss  $\tilde{O}\left(\frac{d^{2-\alpha}L^2}{n^2\varepsilon^2}\right)$  which can offer up to  $\tilde{O}(d)$  improvement. Finally, we note that Wang [34]—which focuses on approximate  $(\varepsilon, \delta)$ -DP—develops private algorithms for linear regression that exhibit good adaptivity to the difficulty of the underlying instance. There exists an extensive prior work on private linear regression and—as this is not the main focus of our work—we refer the reader to [34, 28] for a survey of results.

**Comparison with the smooth sensitivity framework** We conclude the paper with a short comparison of the smooth sensitivity framework and inverse sensitivity mechanisms. While smooth sensitivity mechanisms may not be instance-optimal in many settings, Asi and Duchi [4] show that the inverse sensitivity mechanism is (nearly) instance-optimal for most well-behaved functions and can offer quadratic improvement in sample complexity over smooth sensitivity mechanisms in certain settings. The inverse sensitivity mechanism also outperforms smooth Laplace uniformly for every instance for natural families of sample-monotone functions (see Section 4.3 in [4]). As our development in this paper shows, the approximate versions of the inverse sensitivity mechanism still enjoy similar advantages over smooth mechanisms. Proposition 2.2 shows that—for certain choices of approximations—the approximate inverse sensitivity mechanisms uniformly outperform the smooth Laplace mechanism for every instance. Moreover, the smooth sensitivity framework requires adding noise with heavy-tailed distributions and unbounded moments (such as Cauchy) to preserve  $\varepsilon$ -differential privacy, in contrast to the approximate inverse sensitivity mechanisms which (depending on the approximation and inverse sensitivity) has noise with exponentially decaying tails, resulting in better high-probability bounds and confidence intervals. The PCA example clearly demonstrates these advantages where the approximate inverse sensitivity mechanism enjoys a factor of  $\sqrt{d}$  improvement in sample complexity over smooth sensitivity mechanisms.

## Broader Impact

The substantial growth in data collection and analysis and the increasing awareness for privacy concerns has led to a growing body of work on privacy risks in both academic [16] and industrial settings [17, 3]. Differential privacy [16] has emerged as the standard method for preserving privacy and has enjoyed several applications including in statistical estimation [11], machine learning [5], and game theory [24].

Unfortunately, it is usually challenging to develop private algorithms that achieve satisfactory utility [11]. Therefore, while differential privacy has been successfully deployed in several industrial companies, most applications instantiate a large privacy parameter  $\epsilon$  to achieve acceptable utility, potentially compromising the privacy of users [1].

However, the standard approach in differential privacy to measure the performance of an algorithm is through its (worst case) minimax risk [11]. This—as our theory demonstrates—may be too pessimistic in general and may not capture the correct trade-off between privacy and utility for natural data that arises in real-life. An instance-specific understanding of this trade-off can therefore result in significant improvements in both utility and privacy.

We hope that this work—and instance-optimality in differential privacy in general [4]—can lead to a better understanding of the privacy-utility trade-off of private algorithms for the underlying data at hand. By exploiting the average-case nature of data in real life, we believe that the instance-optimal algorithms we develop can achieve satisfying utility with significantly stronger privacy protections for users.

## Funding Transparency Statement

Funding in direct support of this work: NSF CAREER CCF-1553086, ONR YIP N00014-19-2288, Sloan Foundation, NSF HDR 1934578 (Stanford Data Science Collaboratory), and Stanford DAWN Consortium.

## References

- [1] M. Abadi, U. Erlingsson, I. Goodfellow, H. B. McMahan, N. Papernot, I. Mironov, K. Talwar, and L. Zhang. On the protection of private information in machine learning systems: Two recent approaches. In *30th IEEE Computer Security Foundations Symposium (CSF)*, pages 1–6. IEEE, 2017.
- [2] T. W. Anderson. The integral of a symmetric unimodal function over a symmetric convex set and some probability inequalities. *Proceedings of the American Mathematical Society*, 6(2):170–176, 1955.
- [3] Apple Differential Privacy Team. Learning with privacy at scale, 2017. Available at <https://machinelearning.apple.com/2017/12/06/learning-with-privacy-at-scale.html>.
- [4] H. Asi and J. Duchi. Near instance-optimality in differential privacy. *arXiv:2005.10630 [cs.CR]*, 2020.
- [5] R. Bassily, A. Smith, and A. Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *55th Annual Symposium on Foundations of Computer Science*, pages 464–473, 2014.
- [6] M. Bun and T. Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography Conference (TCC)*, pages 635–658, 2016.
- [7] M. Bun and T. Steinke. Average-case averages: Private algorithms for smooth sensitivity and mean estimation. In *Advances in Neural Information Processing Systems 32*, pages 181–191, 2019.
- [8] M. Bun, G. Kamath, T. Steinke, and Z. S. Wu. Private hypothesis selection. In *Advances in Neural Information Processing Systems 32*, pages 156–167, 2019.
- [9] T. Cai and M. Low. A framework for estimating convex functions. *Statistica Sinica*, 25:423–456, 2015.
- [10] K. Chaudhuri, A. Sarwate, and K. Sinha. Near-optimal algorithms for differentially-private principal components. In *Advances in Neural Information Processing Systems 25*, 2012. URL <http://arxiv.org/abs/1207.2812>.
- [11] J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Minimax optimal procedures for locally private estimation (with discussion). *Journal of the American Statistical Association*, 113(521):182–215, 2018.

- [12] C. Dwork and J. Lei. Differential privacy and robust statistics. In *Proceedings of the Forty-First Annual ACM Symposium on the Theory of Computing*, pages 371–380, 2009.
- [13] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3 & 4):211–407, 2014.
- [14] C. Dwork and G. Rothblum. Concentrated differential privacy. *arXiv:1603.01887 [cs.DS]*, 2016.
- [15] C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor. Our data, ourselves: Privacy via distributed noise generation. In *Advances in Cryptology (EUROCRYPT 2006)*, 2006.
- [16] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the Third Theory of Cryptography Conference*, pages 265–284, 2006.
- [17] U. Erlingsson, V. Pihur, and A. Korolova. RAPPOR: Randomized aggregatable privacy-preserving ordinal response. In *Proceedings of the 21st ACM Conference on Computer and Communications Security (CCS)*, 2014.
- [18] A. Gonen and R. Gilad-Bachrach. Smooth sensitivity based approach for differentially private PCA. In *Algorithmic Learning Theory*, pages 438–450, 2018.
- [19] M. Hardt and K. Talwar. On the geometry of differential privacy. In *Proceedings of the Forty-Second Annual ACM Symposium on the Theory of Computing*, pages 705–714, 2010. URL <http://arxiv.org/abs/0907.3754>.
- [20] A. Johnson and V. Shmatikov. Privacy-preserving data exploration in genome-wide association studies. In *Proceedings of the 19th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 1079–1087, 2013.
- [21] M. Kapralov and K. Talwar. On differentially private low rank approximation. In *Proceedings of the Twenty-Fourth ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1395–1414, 2013.
- [22] S. P. Kasiviswanathan, K. Nissim, S. Raskhodnikova, and A. Smith. Analyzing graphs with node differential privacy. In A. Sahai, editor, *Theory of Cryptography*, volume 7785 of *Lecture Notes in Computer Science*, pages 457–476. Springer, 2013.
- [23] P. W. Laud, P. Damien, and T. S. Shively. Sampling some truncated distributions via rejection algorithms. *Communications in Statistics: Simulation and Computation*, 39(6):1111–1121, 2010.
- [24] F. McSherry and K. Talwar. Mechanism design via differential privacy. In *48th Annual Symposium on Foundations of Computer Science*, 2007.
- [25] K. Nissim, S. Raskhodnikova, and A. Smith. Smooth sensitivity and sampling in private data analysis. In *Proceedings of the Thirty-Ninth Annual ACM Symposium on the Theory of Computing*, 2007.
- [26] R. Okada, K. Fukuchi, and J. Sakuma. Differentially private analysis of outliers. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 458–473. Springer, 2015.
- [27] M. Reimherr and J. Awan. KNG: The K-norm gradient mechanism. In *Advances in Neural Information Processing Systems 32*, pages 10208–10219, 2019.
- [28] O. Sheffet. Differentially private ordinary least squares. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3105–3114, 2017.
- [29] O. Sheffet. Homework 1 for *Differential Privacy: Privacy-preserving Data Analysis*. University of Alberta course CMPUT651, 2018. URL <http://webdocs.cs.ualberta.ca/~osheffet/HW1W18.pdf>.
- [30] A. Smith and A. Thakurta. Differentially private feature selection via stability arguments, and the robustness of the Lasso. In *Proceedings of the Twenty Sixth Annual Conference on Computational Learning Theory*, pages 819–850, 2013. URL <http://proceedings.mlr.press/v30/Guha13.html>.
- [31] C. Stein. Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, pages 197–206, 1956.
- [32] J. Ullman and A. Sealfon. Efficiently estimating Erdos-Renyi graphs with node differential privacy. In *Advances in Neural Information Processing Systems 32*, pages 3765–3775, 2019.
- [33] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998.

- [34] Y. X. Wang. Revisiting differentially private linear regression: optimal and adaptive prediction & estimation in unbounded domain. In *Proceedings of the 34th Conference on Uncertainty in Artificial Intelligence*, pages 93–103, 2018.

# Appendix

## A Proofs of Section 2.1

### A.1 Proof of Theorem 1

Using the privacy guarantees of the exponential mechanism [24, Theorem 6], it is enough to prove that  $\overline{\text{len}}_f(\mathbf{x}; t)$  is 1-Lipschitz. Let  $\mathbf{x}, \mathbf{x}'$  be two neighboring datasets and let  $t \in \mathcal{T}$ . Assume w.l.o.g. that  $\ell = \overline{\text{len}}_f(\mathbf{x}; t) \leq \overline{\text{len}}_f(\mathbf{x}'; t)$ . We need to prove that  $\overline{\text{len}}_f(\mathbf{x}'; t) \leq \ell + 1$ . From the definition of  $\overline{\text{len}}$  in Equation (1), we have that  $\sum_{i=1}^{\ell} R_i(\mathbf{x}) \geq \|t - f(\mathbf{x}')\|_p$ , so the claim follows from conditions on  $R_\ell(\cdot)$  since

$$\|t - f(\mathbf{x}')\|_p \leq \|f(\mathbf{x}) - f(\mathbf{x}')\|_p + \|t - f(\mathbf{x})\|_p \leq R_1(\mathbf{x}') + \sum_{i=1}^{\ell} R_i(\mathbf{x}) \leq \sum_{i=1}^{\ell+1} R_i(\mathbf{x}').$$

### A.2 Proof of Proposition 2.1

To prove the claim about privacy, it is enough to show that  $R_\ell$  satisfy the conditions of Lemma 1. We clearly have  $\text{LS}(\mathbf{x}) \leq \overline{\text{LS}}(\mathbf{x}) \leq R_1(\mathbf{x})$  for all  $\mathbf{x} \in \mathcal{X}^n$ . Moreover, we have that  $R_\ell(\mathbf{x}) \leq R_{\ell+1}(\mathbf{x}')$  for all neighboring datasets  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}^n$  since

$$R_\ell(\mathbf{x}) = \sup_{\mathbf{x}_1: d_{\text{ham}}(\mathbf{x}, \mathbf{x}_1) \leq \ell} \overline{\text{LS}}(\mathbf{x}_1) \leq \sup_{\mathbf{x}_1: d_{\text{ham}}(\mathbf{x}', \mathbf{x}_1) \leq \ell+1} \overline{\text{LS}}(\mathbf{x}_1) = R_{\ell+1}(\mathbf{x}').$$

### A.3 Proof of Theorem 2

We begin with the exact inverse sensitivity mechanism. Let  $C^k = \{t : \text{len}_f^\rho(\mathbf{x}; t) = k\}$ . The definition of  $\text{len}^\rho$  implies that  $\text{len}_f^\rho(\mathbf{x}; t) = 0$  for  $t$  such that  $\|t - f(\mathbf{x})\|_2 \leq \rho$  and that  $\text{len}_f^\rho(\mathbf{x}; t) \geq K$  for any  $t$  such that  $\|t - f(\mathbf{x})\|_2 \geq \omega_f(\mathbf{x}; K) + \rho$ . Thus we have that

$$\begin{aligned} P(\|M_{\text{inv}}(\mathbf{x}) - f(\mathbf{x})\|_2 \geq \omega_f(\mathbf{x}; K) + \rho) &\leq \sum_{k=K}^n \mathbb{P}(M_{\text{inv}}(\mathbf{x}) \in C^k) \\ &\leq \frac{e^{-K\varepsilon/2} \sum_{k=K}^n \int_{s \in C^k} ds}{\int_{s \in \mathcal{T}} e^{-\text{len}_f^\rho(\mathbf{x}; t)\varepsilon/2} ds} \\ &\leq e^{-K\varepsilon/2} \frac{\text{Vol}\{t : \|t\|_2 \leq D\}}{\text{Vol}\{t : \|t\|_2 \leq \rho\}} \\ &\leq e^{-K\varepsilon/2} (D/\rho)^d. \end{aligned}$$

This gives the first part of the claim.

Now we prove the bounds for the approximate mechanism. First, we notice that the noise added by the approximate mechanism (M.2) satisfies

$$z(\mathbf{x}) := M(\mathbf{x}) - f(\mathbf{x}),$$

where  $\mathbb{P}(z(\mathbf{x}) = z) \propto e^{-k\varepsilon/2}$  for  $z \in B^k = \{z \in \mathbb{R}^d : \sum_{i=1}^{k-1} S_i \leq \|z\|_2 < \sum_{i=1}^k S_i\}$ . Noting that  $\|z(\mathbf{x})\|_2 \geq \sum_{i=1}^K R_i(\mathbf{x})$  implies that  $z(\mathbf{x}) \in B^k$  for  $k \geq K$ , we get that

$$\begin{aligned} \mathbb{P}\left(\|M_{\text{appr}}(\mathbf{x}) - f(\mathbf{x})\|_2 \geq \sum_{i=1}^K R_i(\mathbf{x})\right) &= \mathbb{P}\left(\|z(\mathbf{x})\|_2 \geq \sum_{i=1}^K R_i(\mathbf{x})\right) \\ &\leq \sum_{k=K}^n \mathbb{P}(z(\mathbf{x}) \in B^k) \\ &\leq e^{-K\varepsilon/2} \frac{\text{Vol}\{t : \|t\|_2 \leq nD\}}{e^{-1} \text{Vol}\{t : \|t\|_2 \leq \sum_{i=1}^{1/\varepsilon} R_i(\mathbf{x})\}} \\ &\leq e^{-K\varepsilon/2+1} \left(\frac{nD}{\sum_{i=1}^{1/\varepsilon} R_i(\mathbf{x})}\right)^d, \end{aligned}$$

where the last inequality follows using that the ratio of the volumes of two  $\ell_p$ -balls with radii  $r_1$  and  $r_2$  is  $(r_1/r_2)^d$ .

#### A.4 Proof of Proposition 2.2

The claim about privacy follows from identical arguments to the proof of Proposition 2.1. We now prove the claim about utility. We remove  $\mathbf{x}$  to simplify notation. First, we have that  $1 \leq \frac{R_i}{R_1} \leq e^{i\beta}$  from the definition of smooth sensitivity. Thus we have

$$\begin{aligned} \mathbb{E}[\|M(\mathbf{x}) - f(\mathbf{x})\|] &= \frac{\sum_{i=1}^n e^{-i\varepsilon/2} R_i \sum_{j=1}^i R_j}{\sum_{i=1}^n e^{-i\varepsilon/2} R_i} \\ &\leq \frac{\sum_{i=1}^n e^{-i\varepsilon/2} e^{i\beta} R_1^2 \sum_{j=1}^i e^{j\beta}}{\sum_{i=1}^n e^{-i\varepsilon/2} R_1} \\ &\leq R_1 \frac{\sum_{i=1}^n e^{-i\varepsilon/2} e^{i\beta} \frac{e^{i\beta}}{e^\beta - 1}}{\sum_{i=1}^n e^{-i\varepsilon/2}} \\ &= \frac{R_1}{e^{\varepsilon/8} - 1} \frac{\sum_{i=1}^n e^{-i\varepsilon/4}}{\sum_{i=1}^n e^{-i\varepsilon/2}} \\ &= \frac{R_1}{e^{\varepsilon/8} - 1} \frac{e^{\varepsilon/2} - 1}{e^{\varepsilon/4} - 1} = O\left(\frac{S^\beta(\mathbf{x})}{\varepsilon}\right), \end{aligned}$$

where the last equality follows since  $\varepsilon = O(1)$ .

## B Proofs of Section 2.2

### B.1 Proof of Lemma 2.1

We begin with some notation. Let  $(\mathbf{x}, \mathbf{y})$  and  $(\mathbf{x}', \mathbf{y}')$  be two neighboring instances, and denote their minimizers by  $\hat{\theta}_n$  and  $\hat{\theta}'_n$ , respectively. We let  $D = \text{diam}_2(\Theta)$ .

The following lemma bounds the distance between these minimizers.

**Lemma B.1.** *Under the assumptions of Proposition 2.1,*

$$\left\|\hat{\theta}_n - \hat{\theta}'_n\right\|_2 \leq \frac{2L}{\lambda n}.$$

To prove Lemma B.1, we first prove the following weaker version.

**Lemma B.2.** *Assume  $L_n(\theta; \mathbf{x}, \mathbf{y})$  is  $\lambda$ -strongly convex and  $L$ -Lipschitz. Then*

$$\left\|\hat{\theta}_n - \hat{\theta}'_n\right\|_2 \leq \frac{L}{\lambda n}.$$



**Proof** Since  $L_n$  is  $\lambda$ -strongly convex and  $L$ -Lipschitz, we have

$$\begin{aligned} \lambda \left\| \hat{\theta}_n - \hat{\theta}'_n \right\|_2^2 &\leq \langle \nabla L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y}) - \nabla L_n(\hat{\theta}'_n; \mathbf{x}, \mathbf{y}), \hat{\theta}_n - \hat{\theta}'_n \rangle \\ &\leq \left\| \nabla L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y}) - \nabla L_n(\hat{\theta}'_n; \mathbf{x}, \mathbf{y}) \right\|_2 \left\| \hat{\theta}_n - \hat{\theta}'_n \right\|_2. \end{aligned}$$

The claim now follows since  $\nabla L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y}) = 0$  and  $\left\| \nabla L_n(\hat{\theta}'_n; \mathbf{x}, \mathbf{y}) \right\|_2 \leq \frac{L}{n}$ .  $\square$

**Proof** [of Lemma B.1] First, we have that  $\nabla L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y}) = 0$  and  $\left\| \nabla L_n(\hat{\theta}'_n; \mathbf{x}, \mathbf{y}) \right\|_2 \leq \frac{L}{n}$  since  $\ell(\cdot; x_i)$  is  $L$ -Lipschitz. We split to cases whether  $\left\| \hat{\theta}_n - \hat{\theta}'_n \right\|_2 \leq \frac{\lambda}{2H}$ . First, if  $\left\| \hat{\theta}_n - \hat{\theta}'_n \right\|_2 \leq \frac{\lambda}{2H}$  then we know that the function  $L_n(\theta; \mathbf{x}, \mathbf{y})$  is  $\lambda/2$ -strongly convex on the set  $A = \{\theta : \left\| \hat{\theta}_n - \theta \right\|_2 \leq \frac{\lambda}{2H}\}$ . We have  $\hat{\theta}'_n \in A$ , and therefore Lemma B.2 implies that  $\left\| \hat{\theta}_n - \hat{\theta}'_n \right\|_2 \leq \frac{2L}{\lambda n}$ . Now assume that  $\left\| \hat{\theta}_n - \hat{\theta}'_n \right\|_2 > \frac{\lambda}{2H}$  and we get a contradiction. Indeed let  $\theta_t = (1-t)\hat{\theta}_n + t\hat{\theta}'_n$ . For any  $0 \leq t \leq 1$  such that  $\left\| \hat{\theta}_n - \theta_t \right\|_2 \leq \frac{\lambda}{2H}$ , we have that

$$\begin{aligned} \lambda \left\| \hat{\theta}_n - \theta_t \right\|_2^2 &\leq \langle \nabla L_n(\theta_t; \mathbf{x}, \mathbf{y}), \theta_t - \hat{\theta}_n \rangle \\ &\stackrel{(i)}{\leq} \langle \nabla L_n(\theta_1; \mathbf{x}, \mathbf{y}), \theta_1 - \hat{\theta}_n \rangle \\ &\leq \frac{L \left\| \hat{\theta}'_n - \hat{\theta}_n \right\|_2}{n}, \end{aligned}$$

where the third inequality follows from Cauchy-Schwartz inequality since  $\theta_1 = \hat{\theta}'_n$  and (i) follows from a monotonicity argument which we explain presently. This implies that  $\left\| \hat{\theta}'_n - \hat{\theta}_n \right\|_2 \geq \frac{n\lambda^2}{4LH^2}$  which is a contradiction.

Let us now explain why inequality (i) holds. First we denote  $u = \hat{\theta}'_n - \hat{\theta}_n$  and we notice that  $\theta_t = \hat{\theta}_n + tu$ . Define  $g(t) = L_n(\theta_t; \mathbf{x})$  which is convex in  $t$ . As  $g$  is convex with minimizer at  $t^* = 0$ , we have  $g'(0) = 0$  and  $0 \leq g'(t) \leq g'(s)$  for  $0 \leq t \leq s$ . Therefore we have that  $0 \leq \langle \nabla L_n(\theta_t; \mathbf{x}, \mathbf{y}), u \rangle \leq \langle \nabla L_n(\theta_s; \mathbf{x}, \mathbf{y}), u \rangle$ . Inequality (i) now follows since

$$\begin{aligned} \langle \nabla L_n(\theta_t; \mathbf{x}, \mathbf{y}), \theta_t - \hat{\theta}_n \rangle &= t \langle \nabla L_n(\theta_t; \mathbf{x}, \mathbf{y}), u \rangle \\ &\leq s \langle \nabla L_n(\theta_s; \mathbf{x}, \mathbf{y}), u \rangle \\ &= \langle \nabla L_n(\theta_s; \mathbf{x}, \mathbf{y}), \theta_s - \hat{\theta}_n \rangle. \end{aligned}$$

$\square$

Now we are ready to prove Proposition 2.1. First, Lemma B.1 implies that

$$\left\| \hat{\theta}_n - \hat{\theta}'_n \right\|_2 \leq \frac{2L}{\lambda n}.$$

Therefore as  $\nabla^2 L_n(\cdot; \mathbf{x}, \mathbf{y})$  is  $H$ -Lipschitz, we get

$$\left\| \nabla^2 L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y}) - \nabla^2 L_n(\hat{\theta}'_n; \mathbf{x}, \mathbf{y}) \right\|_2 \leq \frac{2LH}{\lambda n}.$$

As a consequence we have

$$\begin{aligned}
& \left| \left\| \nabla^2 L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y})(\theta - \hat{\theta}_n) \right\|_2 - \left\| \nabla^2 L_n(\hat{\theta}'_n; \mathbf{x}', \mathbf{y}')(\theta - \hat{\theta}'_n) \right\|_2 \right| \\
&= \left| \left\| \nabla^2 L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y})(\theta - \hat{\theta}_n) \right\|_2 - \left\| \nabla^2 L_n(\hat{\theta}'_n; \mathbf{x}', \mathbf{y}')(\theta - \hat{\theta}_n) + \nabla^2 L_n(\hat{\theta}'_n; \mathbf{x}', \mathbf{y}')(\hat{\theta}_n - \hat{\theta}'_n) \right\|_2 \right| \\
&\leq \left\| \nabla^2 L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y}) - \nabla^2 L_n(\hat{\theta}'_n; \mathbf{x}', \mathbf{y}') \right\|_2 \left\| \theta - \hat{\theta}_n \right\|_2 + \left\| \nabla^2 L_n(\hat{\theta}'_n; \mathbf{x}', \mathbf{y}')(\hat{\theta}_n - \hat{\theta}'_n) \right\|_2 \\
&\leq \frac{2DLH}{\lambda n} + \left\| \nabla L_n(\hat{\theta}_n; \mathbf{x}', \mathbf{y}') \right\|_2 + O(H \left\| \hat{\theta}_n - \hat{\theta}'_n \right\|_2^2) \\
&\leq \frac{2DLH}{\lambda n} + \frac{L}{n} + O\left(H \left(\frac{2LH}{\lambda n}\right)^2\right).
\end{aligned}$$

## B.2 Proof of Proposition 2.3

The claim about privacy is immediate from the exponential mechanism. Let us now argue about the claim for utility. To simplify notation, we let  $\Sigma = \nabla^2 L(\hat{\theta}_n; \mathbf{x}, \mathbf{y})$ . We later show (see Section F.1) that to sample from the distribution (5), one can sample  $R \sim \text{Gamma}(d, 1)$  and  $U \sim \text{Uni}(\mathbb{S}^{d-1})$  and then set  $\theta = \hat{\theta}_n + Z$  where  $Z = \frac{2\text{GS}_{\text{Hess}}}{\varepsilon} \Sigma^{-1} \cdot R \cdot U$ , and finally we accept  $\theta$  if  $\theta \in \Theta$ , otherwise we repeat the process, It is easy to show that for  $Z$  we have

$$\mathbb{E} \left[ \|Z\|_2^2 \right] = \frac{4\mathbb{E}[R^2]\text{GS}_{\text{Hess}}^2}{\varepsilon^2} \mathbb{E} \left[ \|\Sigma^{-1}U\|_2^2 \right] \leq \frac{Cd\text{GS}_{\text{Hess}}^2}{\varepsilon^2} \text{tr}(\Sigma^{-2}),$$

for a universal constant  $C$ . But we need to upper bound  $\mathbb{E} \left[ \|Z\|_2^2 \mid \hat{\theta}_n + Z \in \Theta \right]$  as this is the error of the mechanism. To finish the proof, we now prove that for every random variable  $W$ ,

$$\mathbb{E} \left[ \|W\|_2^2 \mid \hat{\theta}_n + Z \in \Theta \right] \leq 2\mathbb{E} \left[ \|W\|_2^2 \right]. \quad (9)$$

To this end, we let  $\rho^2 = \mathbb{E} \left[ \|Z\|_2^2 \right]$  and define three disjoint sets,  $S_1 = \{Z : \|Z\|_2 \leq 2\rho\}$ ,  $S_2 = \{Z : \hat{\theta}_n + Z \in \Theta, Z \notin S_1\}$ , and  $S_3 = \mathbb{R}^d \setminus (S_1 \cup S_2)$ . Clearly these sets are disjoint and the assumptions of the Proposition imply that  $S_1 \subseteq \Theta$  and therefore  $\Theta = S_1 \cup S_2$ . Using conditional expectation and denoting  $p_i = \mathbb{P}(Z \in S_i)$ , we have that

$$\mathbb{E} \left[ \|W\|_2^2 \right] = \sum_{i=1}^3 p_i \mathbb{E} \left[ \|W\|_2^2 \mid Z \in S_i \right].$$

Noting that  $p_1 \geq 1/2$  by Markov inequality, we now get that

$$\begin{aligned}
\mathbb{E} \left[ \|W\|_2^2 \mid \hat{\theta}_n + Z \in \Theta \right] &= \frac{p_1}{p_1 + p_2} \mathbb{E} \left[ \|W\|_2^2 \mid Z \in S_1 \right] + \frac{p_2}{p_1 + p_2} \mathbb{E} \left[ \|W\|_2^2 \mid Z \in S_2 \right] \\
&\leq 2p_1 \mathbb{E} \left[ \|W\|_2^2 \mid Z \in S_1 \right] + 2p_2 \mathbb{E} \left[ \|W\|_2^2 \mid Z \in S_2 \right] \\
&\leq 2\mathbb{E} \left[ \|W\|_2^2 \right].
\end{aligned}$$

The claim follows.

## C Proofs of Section 3 (lower bounds)

### C.1 Proofs of Theorem 3

We start with the lower bound for unbiased mechanisms. Fix  $\mathbf{x} \in \mathcal{X}^n$  and assume towards a contradiction that  $\mathbb{E}[\|M(\mathbf{x}) - f(\mathbf{x})\|_2] \leq \frac{1}{8}\beta_k$  where  $\beta_k = c \cdot e^{-2k\varepsilon/d} \omega_f(\mathbf{x}; k)$ . The definition of  $W_f(\mathbf{x}; k)$  implies that there exists a  $\beta_k/4$  packing, namely  $S$ , of  $W_f(\mathbf{x}; k)$  of size at least  $m_{\beta_k} \geq \left(\frac{4c\omega_f(\mathbf{x}; k)}{\beta_k}\right)^d \geq 4^d e^{2k\varepsilon}$ . The definition of  $W_f(\mathbf{x}; k)$  implies that there is an instance  $\mathbf{x}'$  such

that  $d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq k$  and  $f(\mathbf{x}') = t$  for every  $t \in S$ , hence we have a set  $A$  of size  $m_{\beta_k}$  of datasets  $\mathbf{x}'$  such that  $f(\mathbf{x}') \in S$ . For every  $\mathbf{x}' \in A$ , we define

$$B_{\mathbf{x}'} = \{y : \|y - f(\mathbf{x}')\|_2 \leq \beta_k/4\}.$$

We now have that:

$$\begin{aligned} \mathbb{P}(M(\mathbf{x}) \in B_{\mathbf{x}'}) &\stackrel{(i)}{\geq} \mathbb{P}(M(\mathbf{x}') \in B_{\mathbf{x}'})e^{-k\varepsilon} \stackrel{(ii)}{\geq} \mathbb{P}(M(\mathbf{x}') \in B_{\mathbf{x}})e^{-k\varepsilon} \\ &\stackrel{(iii)}{\geq} \mathbb{P}(M(\mathbf{x}) \in B_{\mathbf{x}})e^{-2k\varepsilon} \stackrel{(iv)}{\geq} \frac{e^{-2k\varepsilon}}{2}, \end{aligned}$$

where (i) and (iii) follow from the definition of differential privacy, (ii) follows since  $M$  is  $\|\cdot\|_2$ -unbiased, and (iv) follows from Markov inequality. As the sets  $B_{\mathbf{x}'}$  are disjoint for  $\mathbf{x}' \in A$ , we have a contradiction

$$1 \geq \sum_{\mathbf{x}' \in A} \mathbb{P}(M(\mathbf{x}) \in B_{\mathbf{x}'}) \geq \frac{1}{2} m_{\beta_k} e^{-2k\varepsilon} \geq \frac{4^d}{2}.$$

To prove the first part of the claim (i.e., the minimax lower bound), we use similar ideas while starting from the assumption that for every  $\mathbf{x}$  we have  $\mathbb{E}[\|M(\mathbf{x}) - f(\mathbf{x})\|_2] \leq \frac{1}{8}\beta_k$  where  $\beta_k = c \cdot \sup_{\mathbf{x}'} e^{-k\varepsilon/d} \omega_f(\mathbf{x}'; k)$ . We again define a packing  $A$  (now with  $m_{\beta_k} \geq 4^d e^{k\varepsilon}$ ) and we get using Markov inequality and the definition of differential privacy that  $\mathbb{P}(M(\mathbf{x}) \in B_{\mathbf{x}'}) \geq \mathbb{P}(M(\mathbf{x}') \in B_{\mathbf{x}'})e^{-k\varepsilon} \geq \frac{e^{-k\varepsilon}}{2}$ . This gives a contradiction similarly to our argument above.

## C.2 Local-minimax lower bound and proof of Theorem 4

First, we start by explaining why  $r = \Omega(d/\varepsilon)$  is necessary to exclude trivial mechanisms in the local minimax definition (6). Assume that we choose  $r \ll d/\varepsilon$ . Then for an instance  $\mathbf{x} \in \mathcal{X}^n$ , consider the trivial constant mechanism that sets  $M_{\text{triv}}(\mathbf{x}') = f(\mathbf{x})$  for every  $\mathbf{x}' \in \mathcal{X}^n$ . Clearly this mechanism is  $\varepsilon$ -differentially private and its local-minimax risk for  $\mathbf{x}$  is

$$\sup_{\mathbf{x}': d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq r} \mathbb{E}[\|M_{\text{triv}}(\mathbf{x}') - f(\mathbf{x}')\|_2] = \sup_{\mathbf{x}': d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq r} \|f(\mathbf{x}) - f(\mathbf{x}')\|_2 = \omega_f(\mathbf{x}; r).$$

Our lower bounds of Lemma C.1 on the local minimax risk with radius  $r$  show this is the optimal risk—up to constant factors that do not depend on  $r$ —for  $\mathbf{x}$  whenever  $r \leq d/\varepsilon$ . Therefore we need to pick a larger value of  $r$  such that the optimal mechanism is not the trivial constant mechanism. Picking  $r = C \cdot d/\varepsilon$ , our lower bounds are roughly  $\sup_{1 \leq i \leq C} e^{-i} \omega_f(\mathbf{x}; id/\varepsilon)$  (which is usually maximized at  $i = 1$  resulting in  $\omega_f(\mathbf{x}; d/\varepsilon)$ ) and so the trivial mechanism does not achieve this for  $C$  large enough.

Moreover, when  $r \ll d/\varepsilon$ , no single mechanism  $M$  can be instance-optimal according to this definition: if there exists  $M^*$  such that  $\mathbb{E}[\|M^*(\mathbf{x}) - f(\mathbf{x})\|_2] \leq O(1)\mathcal{R}(\mathbf{x}; r) \leq O(1)\omega_f(\mathbf{x}; r)$  for every  $\mathbf{x} \in \mathcal{X}^n$ , then we get that

$$\sup_{\mathbf{x} \in \mathcal{X}^n} \mathbb{E}[\|M^*(\mathbf{x}) - f(\mathbf{x})\|_2] \leq O(1) \sup_{\mathbf{x} \in \mathcal{X}^n} \omega_f(\mathbf{x}; r) \ll \sup_{\mathbf{x} \in \mathcal{X}^n} \omega_f(\mathbf{x}; d/\varepsilon),$$

which contradict the minimax lower bounds of Theorem 3.

Theorem 4 follows from the following lemma by setting  $r = d/\varepsilon$ .

**Lemma C.1.** *Let the assumptions of Theorem 4 hold. Then for any  $r \geq 1$ ,  $\mathcal{R}(\mathbf{x}; r) \geq \frac{c}{8} \sup_{k \leq r} e^{-k\varepsilon/d} \omega_f(\mathbf{x}; k)$ .*

**Proof** The proof follows similar arguments to those we had in the proof of Theorem 3. Assume toward a contradiction that  $\mathcal{R}(\mathbf{x}; r) \leq \frac{1}{8}\beta_k$  where  $\beta_k = c \cdot e^{-k\varepsilon/d} \omega_f(\mathbf{x}; k)$  for  $k \leq r$ . This implies that there exists a mechanism  $M$  such that for every  $\mathbf{x}'$  such that  $d_{\text{ham}}(\mathbf{x}, \mathbf{x}') \leq r$ , we have  $\mathbb{E}[\|M(\mathbf{x}') - f(\mathbf{x}')\|_2] \leq \frac{1}{8}\beta_k$ . Repeating the arguments of the proof of Theorem 3 with  $W_f(\mathbf{x}; k)$  for  $k \leq r$  proves the claim.  $\square$

## D Proofs and further details of Section 4.1 (mean estimation)

In this section we provide proofs for the claims in Section 4.1 and give our algorithm. We begin by giving the full version of Lemma 4.1, which we prove in Appendix D.2.

**Lemma** (Full version of Lemma 4.1). *Let  $f(\mathbf{x}) = [\text{trim}_m(\mathbf{x})]_{[a,b]}$ . Then, for any  $t \in [a, b]$ , if  $t \geq \text{trim}_m(\mathbf{x})$*

$$\text{len}_f(\mathbf{x}; t) = \min\{k : k \leq m, |t - f(\mathbf{x})| \leq \frac{1}{n-2m} \sum_{i=1}^k (x_{(n-m+i)} - x_{(m+i)})\} \cup \{m+1\}.$$

Moreover, if  $t \leq \text{trim}_m(\mathbf{x})$

$$\text{len}_f(\mathbf{x}; t) = \min\{k : k \leq m, |t - f(\mathbf{x})| \leq \frac{1}{n-2m} \sum_{i=1}^k (x_{(n-m+1-i)} - x_{(m+1-i)})\} \cup \{m+1\}.$$

### D.1 Sampling from the inverse sensitivity mechanism

We describe an algorithm for sampling from the inverse sensitivity mechanism with  $\rho > 0$  for the mean estimation problem of Section 4.1. Our goal is to sample from

$$\pi_{M_{\text{inv}}(\mathbf{x})}(t) = \frac{e^{-\text{len}^\rho(\mathbf{x}; t)\varepsilon/2}}{\int_{\mathcal{T}} e^{-\text{len}^\rho(\mathbf{x}; s)\varepsilon/2} dS}.$$

Algorithm 4 shows how to sample from this distribution using Lemma 4.1. To see this, note that Lemma 4.1 implies that  $S_k$  in Algorithm 4 is exactly the set  $\{t : \text{len}^\rho(\mathbf{x}; t) = k\}$ . And so all values  $t \in S_k$  have the same probability. Moreover, the probability of sampling a value from the set  $S_k$  is  $\text{Vol}(S_k)e^{-k\varepsilon/2}$  using the definition of the mechanism.

---

#### Algorithm 4: Inverse sensitivity for mean estimation

---

**Input:**  $\mathbf{x} \in \mathbb{R}^n$ ,  $m, \rho, a, b$

- 1 Calculate  $\hat{x}_t = [\text{trim}_m(\mathbf{x})]_{[a,b]}$ ;
  - 2 Calculate  $u_k = \min\left(\rho + \frac{1}{n-2m} \sum_{i=1}^k (x_{(n-m+i)} - x_{(m+i)}), b - \hat{x}_t\right)$  for  $0 \leq k \leq m$ ;
  - 3 Calculate  $\ell_k = \min\left(\rho + \frac{1}{n-2m} \sum_{i=1}^k (x_{(n-m+1-i)} - x_{(m+1-i)}), \hat{x}_t - a\right)$  for  $0 \leq k \leq m$ ;
  - 4 Set  $u_{m+1} = b - \hat{x}_t$  and  $\ell_{m+1} = \hat{x}_t - a$ ;
  - 5 Set  $S_k = [-\ell_{k+1}, -\ell_k] \cup [u_k, u_{k+1}]$ ;
  - 6 Sample  $k \propto \text{Vol}(S_k)e^{-k\varepsilon/2}$ ;
  - 7 Sample  $z \sim \text{Uni}(S_k)$ ;
  - 8 **return**  $\hat{x} = \hat{x}_t + z$
- 

### D.2 Proof of Lemma 4.1

We only prove the case  $t \geq \text{trim}_m(\mathbf{x})$  as the other one is similar. In this case, to make the value of the trimmed mean  $t$  by changing  $k$  values, we must change the value of the  $k$  smallest samples  $x_{(1)}, \dots, x_{(k)}$  and set their value to  $\infty$ . Denote the resulting sample by  $\mathbf{x}'$ . The trimmed mean  $\mathbf{x}'$  is

$$\text{trim}_m(\mathbf{x}') = \frac{x'_{(m+1)} + x'_{(m+2)} + \dots + x'_{(n-m)}}{n-2m}.$$

We split to two cases. First, if  $k \geq m+1$ , then we get that  $x'_{(n-m)} = \infty$  and therefore  $\text{trim}_m(\mathbf{x}') = \infty$ . This means that for any  $t > \text{trim}_m(\mathbf{x})$  we can set suitable new values to  $x_{(1)}, \dots, x_{(k)}$  (instead of  $\infty$ ) such that  $\text{trim}_m(\mathbf{x}') = t$ . Now assume  $k \leq m$ . In this case, we get that

$$\text{trim}_m(\mathbf{x}') = \frac{x_{(m+k+1)} + x_{(m+k+2)} + \dots + x_{(n+k-m)}}{n-2m}.$$

The claim follows as we have

$$\text{trim}_m(\mathbf{x}') - \text{trim}_m(\mathbf{x}) = \frac{\sum_{i=1}^k (x_{(n-m+i)} - x_{(m+i)})}{n-2m}.$$

### D.3 Proof of Proposition 4.1

The privacy guarantees of Algorithm 4 follow from Proposition 3.2 in [4].

To prove the claim about utility, we use the following result from [7] which upper bounds the error of the trimmed mean estimator.

**Lemma D.1** (Bun and Steinke [7], Proposition 10). *Let  $x_i \stackrel{\text{iid}}{\sim} P$  where  $P$  has mean  $\mu$  and variance  $\sigma^2$ . Then*

$$\mathbb{E}[(\text{trim}_m(\mathbf{x}) - \mu)^2] \leq \frac{n(1 + \sqrt{8m})}{(n - 2m)^2} \sigma^2 = O\left(\frac{m}{n}\right) \sigma^2.$$

Moreover, if  $P$  is symmetric about its mean then

$$\mathbb{E}[(\text{trim}_m(\mathbf{x}) - \mu)^2] \leq \left(1 + O\left(\frac{m}{n}\right)\right) \frac{\sigma^2}{n}.$$

We have that

$$\begin{aligned} \mathbb{E}[(\hat{x} - \mu)^2] &\leq \mathbb{E}[(\text{trim}_m(\mathbf{x})_{[a,b]} - \mu)^2] + \mathbb{E}[z^2] \\ &\leq \mathbb{E}[(\text{trim}_m(\mathbf{x}) - \mu)^2] + \mathbb{E}[z^2], \end{aligned}$$

where the second inequality follows since  $\mu \in [a, b]$ , and so a projection to  $[a, b]$  cannot increase error. Thus, given the bound of Lemma D.1, now we only need to upper bound  $\mathbb{E}[z^2]$ .

We begin with the following lemma for a fixed  $\mathbf{x}$ .

**Lemma D.2.** *Let  $\mathbf{x} \in \mathbb{R}^n$  and  $K \leq m$ . Let  $L(\mathbf{x}) = \max_{1 \leq k \leq K} (u_{k+1} - u_k, \ell_{k+1} - \ell_k)$ . Then*

$$\mathbb{E}[z^2] \leq O\left(\frac{L(\mathbf{x})^2}{\varepsilon^2}\right) + e^{-K\varepsilon/2} \frac{m(b-a)^2}{\rho}.$$

**Proof** Let  $v_k = u_k + \ell_k$ . Then we have that  $v_{k+1} - v_k \leq 2L$  for  $k \leq K$  and  $v_0 \geq \rho$ , hence we get

$$\begin{aligned} \mathbb{E}[z^2] &\leq \frac{\sum_{k=0}^{m+1} e^{-k\varepsilon/2} v_k^2 (v_k - v_{k-1})}{\sum_{k=0}^{m+1} e^{-k\varepsilon/2} (v_k - v_{k-1})} \\ &\leq \frac{\sum_{k=0}^K e^{-k\varepsilon/2} v_k^2 (v_k - v_{k-1})}{\sum_{k=0}^{m+1} e^{-k\varepsilon/2} (v_k - v_{k-1})} + e^{-K\varepsilon/2} \frac{m(b-a)^2}{\rho} \\ &\stackrel{(i)}{\leq} O\left(\frac{L^2}{\varepsilon^2}\right) + e^{-K\varepsilon/2} \frac{m(b-a)^2}{\rho}. \end{aligned}$$

Inequality (i) follows from similar arguments to the proof of Proposition 4.3 in [4]: let  $T \leq K$  be the smallest such that  $v_T \geq \frac{L}{\varepsilon}$ . If no such  $T$  exists, then inequality (i) clearly holds. Note that  $v_T \leq \frac{L}{\varepsilon} + L \leq O\left(\frac{L}{\varepsilon}\right)$ . Thus we get

$$\frac{\sum_{k=1}^K e^{-k\varepsilon/2} v_k^2 (v_k - v_{k-1})}{\sum_{k=1}^{m+1} e^{-k\varepsilon/2} (v_k - v_{k-1})} \leq v_T + \frac{\sum_{k=T}^K 2e^{-k\varepsilon/2} k^2 L^3}{e^{-T\varepsilon/2} L/\varepsilon} = O\left(\frac{L^2}{\varepsilon^2}\right).$$

□

Now we are ready to finish the proof of Proposition 4.1. We notice that for any  $\mathbf{x}$

$$L(\mathbf{x}) \leq \frac{x_{(n)} - x_{(1)}}{n - m} \leq \frac{2 \max_{1 \leq i \leq n} |x_i|}{n - m}.$$

Therefore using standard bound on the expectation of maximum of subgaussian variables (Lemma 45 in full version [7])

$$\mathbb{E}[L(\mathbf{x})^2] \leq \frac{4\mathbb{E}[\max_{1 \leq i \leq n} x_i^2]}{(n - m)^2} \leq \frac{8\sigma^2 \log n}{(n - m)^2}.$$

Overall we have that using  $K = m$  in Lemma D.2 implies that for a constant  $C$ ,

$$\mathbb{E}[(\hat{x} - \mu)^2] \leq \mathbb{E}[(\text{trim}_m(\mathbf{x}) - \mu)^2] + \frac{C\sigma^2 \log n}{(n-m)^2 \varepsilon^2} + e^{-m\varepsilon/2} \frac{m(b-a)^2}{\rho}.$$

As  $\rho = \frac{\sigma^2}{n^2}$  and  $m \leq n$ , setting  $m = \frac{12 \log(n(b-a)/\sigma^2)}{\varepsilon}$  and using Lemma D.1 proves the claim.

## E Proofs of Section 4.2 (PCA)

Here we prove Proposition 4.2. First, we prove the claim about privacy then we proceed to show the utility analysis.

### E.1 Proof of Proposition 4.2 (privacy)

We only need to prove that  $w(\mathbf{x}) = \bar{v} + z$  is  $\varepsilon$ -DP as the claim for Algorithm 2 then follows since post-processing preserves privacy [13, Proposition 2.1]. To this end, first we note that Lemma 4.2 implies that the choice of  $R_i(\mathbf{x})$  in Proposition 4.2 satisfies the conditions of Theorem 1. Now assume we have two neighboring instances  $\mathbf{x}, \mathbf{x}'$  with leading eigenvectors  $u_1, u_2 = -u_1$  and  $u'_1, u'_2 = -u'_1$  respectively and assume without loss of generality  $\|u_1 - u'_1\|_2 \leq \text{LS}(\mathbf{x})$ . We get that  $w(\mathbf{x}) = w_1(\mathbf{x}) = u_1 + z$  with probability  $1/2$  and  $w(\mathbf{x}) = w_2(\mathbf{x}) = u_2 + z$  otherwise. Similarly we have  $w_1(\mathbf{x}')$  and  $w_2(\mathbf{x}')$ . Theorem 1 now implies that the densities of  $w_1(\mathbf{x})$  and  $w_1(\mathbf{x}')$  are  $\varepsilon$ -DP (i.e.,  $\frac{\pi_{w_1(\mathbf{x})}(t)}{\pi_{w_1(\mathbf{x}')} (t)} \leq e^\varepsilon$ ) and similarly for the densities of  $w_2(\mathbf{x})$  and  $w_2(\mathbf{x}')$ , therefore by quasi convexity we get that  $w(\mathbf{x})$  is  $\varepsilon$ -DP.

### E.2 Proof of Proposition 4.2 (utility)

To facilitate notation, we drop  $\mathbf{x}$  from our analysis. First, we bound the norm of the noise  $z$  that the algorithm adds. We claim that there exists a universal constant  $C_1 > 0$  such that the noise  $z$  in step 4 of Algorithm 2 has with probability  $1 - \beta$

$$\|z\|_2 \leq C_1 \left( \frac{d \log n / \beta}{n \text{GAP}(\mathbf{x}) \varepsilon} + \frac{1}{n^4} \right).$$

Deferring the proof of this, we can now complete the proof of the claim. We assume that  $n$  is large enough so that  $\|z\|_2 \leq 1/2$ . Notice that in our setting ( $k = 1$ ) we have that  $F(v) = v^T \Sigma v = \|\Sigma^{1/2} v\|_2^2$ . Therefore denoting  $\lambda = \frac{1}{\|\bar{v}+z\|_2}$  we get that

$$\begin{aligned} |F(v_{\text{out}}) - F(\bar{v})| &\leq |F(v_{\text{out}}) - F(\bar{v} + z)| + |F(\bar{v} + z) - F(\bar{v})| \\ &= |\lambda^2 - 1| F(\bar{v} + z) + \left| \|\Sigma^{1/2}(\bar{v} + z)\|_2^2 - \|\Sigma^{1/2}\bar{v}\|_2^2 \right| \\ &\leq |\lambda^2 - 1| \|\bar{v} + z\|_2 + \left\| \Sigma^{1/2} z \right\|_2^2 \left( \|\Sigma^{1/2}(\bar{v} + z)\|_2^2 + \|\Sigma^{1/2}\bar{v}\|_2^2 \right) \\ &\leq 2|\lambda^2 - 1| + 3\|z\|_2^2, \end{aligned}$$

where we use the fact that  $\|z\|_2 \leq 1/2$ ,  $\|\bar{v}\|_2 = 1$ , and  $\|x_i\|_2 \leq 1$  so that  $\|\Sigma^{1/2} u\|_2 \leq \|u\|_2$  for every  $u$ . Now we only need to upper bound  $|\lambda^2 - 1|$ . As  $\lambda \leq 2$ , we have

$$|\lambda^2 - 1| \leq 3|\lambda - 1| = 3 \frac{|1 - \|\bar{v} + z\|_2|}{\|\bar{v} + z\|_2} \leq 6\|z\|_2.$$

Therefore overall we have that

$$|F(v_{\text{out}}) - F(\hat{v})| \leq 15\|z\|_2,$$

which proves the claim.

Now we return to prove the claim about the norm of  $z$ . We use Theorem 2 with  $K = \frac{c_1 d \log n / \beta}{\varepsilon}$  to get

$$\mathbb{P} \left( \|z\|_2 \geq \sum_{i=1}^K R_i(\mathbf{x}) \right) \leq \frac{e^{-K\varepsilon/2}}{e} \left( \frac{n\sqrt{2}}{\sum_{i=1}^K R_i(\mathbf{x})} \right)^d.$$



Assuming  $K \leq n\text{GAP}(\mathbf{x})/4$  as we can take  $n$  large enough in the assumption of the proposition, we get that  $R_i(\mathbf{x}) \leq \frac{2C_{\text{pca}}}{n\text{GAP}(\mathbf{x})}$  for  $i \leq K$  and therefore  $\sum_{i=1}^K R_i(\mathbf{x}) = O\left(\frac{C_{\text{pca}}d \log n/\beta}{n\text{GAP}(\mathbf{x})\varepsilon}\right)$ . Since  $R_i(\mathbf{x}) \geq \frac{C_{\text{pca}}}{n\text{GAP}(\mathbf{x})}$ , setting  $c_1$  large enough we get that

$$\mathbb{P}\left(\|z\|_2 \geq \sum_{i=1}^K R_i(\mathbf{x})\right) \leq \frac{e^{-K\varepsilon/2}}{e} \left(\frac{n^2\text{GAP}(\mathbf{x})\sqrt{2}}{C_{\text{pca}}}\right)^d \leq \beta.$$

## F Proofs of Section 4.3 (linear regression)

### F.1 Sampling from the gradient mechanism

In this section, we show how Algorithm 3 is basically sampling from the distribution of the gradient mechanism. To this end, we show how to sample a vector  $t \in \mathbb{R}^d$  with density  $\pi(t) = \exp(-\|At\|)$  for a matrix  $A \succ 0$ . The change of variables  $u = At$  and then using rotational symmetry gives that

$$\begin{aligned} \int \pi(t) dt &= \frac{1}{\det(A)} \int \exp(-\|u\|) du = \frac{1}{\det(A)} \int_0^\infty \exp(-r) \text{Vol}_{d-1}(r\mathbb{S}^{d-1}) dr \\ &= \frac{1}{\det(A)} \frac{d\pi^{d/2}}{\Gamma(\frac{d}{2} + 1)} \int_0^\infty r^{d-1} e^{-r} dr = \frac{d\pi^{d/2}\Gamma(d)}{\det(A)\Gamma(\frac{d}{2} + 1)}. \end{aligned}$$

In particular, to sample  $T$  with the density  $\pi(t) = \exp(-\|At\|)$ , we draw  $R \sim \text{Gamma}(d, 1)$ , then  $U \sim \text{Uni}(\mathbb{S}^{d-1})$ , and set  $T = RA^{-1}U$ .

Recall that the gradient mechanism has  $\pi(t) = \exp(-\|At\|)$  only for  $t \in S$  for some set  $S \subset \mathbb{R}^d$  and  $\pi(t) = 0$  otherwise. Therefore we apply rejection sampling until we get  $t \in S$ . This shows that Algorithm 3 is sampling from the gradient mechanism.

### F.2 Proof of Proposition 4.3

Similarly to the proof of Proposition 2.3, and letting  $Z = \frac{2L}{n\varepsilon}\Sigma_n^{-1} \cdot R \cdot U$  for  $R \sim \text{Gamma}(d, 1)$ , and  $U \sim \text{Uni}(\mathbb{S}^{d-1})$ , we note that Algorithm 3 sets  $\theta_{\text{out}} = \bar{\theta} + Z$  and accepts it if  $\theta_{\text{out}} \in \Theta$ . We also have

$$L_n(\theta_{\text{out}}; \mathbf{x}, \mathbf{y}) - L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y}) = (\theta_{\text{out}} - \hat{\theta}_n)^T \Sigma_n (\theta_{\text{out}} - \hat{\theta}_n).$$

Since  $\hat{\theta}_n \in \text{int}(\Theta)$ , we get that  $\bar{\theta} = \hat{\theta}_n$  and thus the excess loss of the algorithm is

$$\mathbb{E} \left[ L_n(\theta_{\text{out}}; \mathbf{x}, \mathbf{y}) - L_n(\hat{\theta}_n; \mathbf{x}, \mathbf{y}) \right] \leq \mathbb{E} \left[ Z^T \Sigma_n Z \mid \hat{\theta}_n + Z \in \Theta \right].$$

Using inequality (9) in the proof of Proposition 2.3, the claim follows since

$$\mathbb{E} \left[ Z^T \Sigma_n Z \mid \bar{\theta} + Z \in \Theta \right] \leq 2\mathbb{E} \left[ Z^T \Sigma_n Z \right] = \frac{8\mathbb{E}[R^2]L^2}{n^2\varepsilon^2} \mathbb{E} \left[ U \Sigma_n^{-1} U \right] = O \left( \frac{dL^2}{n^2\varepsilon^2} \text{tr}(\Sigma_n^{-1}) \right).$$