

---

# Private Stochastic Convex Optimization: Optimal Rates in $\ell_1$ Geometry

---

Hilal Asi<sup>1</sup> Vitaly Feldman<sup>2</sup> Tomer Koren<sup>3</sup> Kunal Talwar<sup>2</sup>

## Abstract

Stochastic convex optimization over an  $\ell_1$ -bounded domain is ubiquitous in machine learning applications such as LASSO but remains poorly understood when learning with differential privacy. We show that, up to logarithmic factors the optimal excess population loss of any  $(\varepsilon, \delta)$ -differentially private optimizer is  $\sqrt{\log(d)/n} + \sqrt{d}/\varepsilon n$ . The upper bound is based on a new algorithm that combines the iterative localization approach of Feldman et al. (2020a) with a new analysis of private regularized mirror descent. It applies to  $\ell_p$  bounded domains for  $p \in [1, 2]$  and queries at most  $n^{3/2}$  gradients improving over the best previously known algorithm for the  $\ell_2$  case which needs  $n^2$  gradients. Further, we show that when the loss functions satisfy additional smoothness assumptions, the excess loss is upper bounded (up to logarithmic factors) by  $\sqrt{\log(d)/n} + (\log(d)/\varepsilon n)^{2/3}$ . This bound is achieved by a new variance-reduced version of the Frank-Wolfe algorithm that requires just a single pass over the data. We also show that the lower bound in this case is the minimum of the two rates mentioned above.

## 1. Introduction

Convex optimization is one of the most well-studied problems in private data analysis. Existing works have largely studied optimization problems over  $\ell_2$ -bounded domains. However several machine learning applications, such as LASSO and minimization over the probability simplex, involve optimization over  $\ell_1$ -bounded domains. In this work we study the problem of differentially private stochastic convex optimization (DP-SCO) over  $\ell_1$ -bounded domains.

---

<sup>1</sup>Department of Electrical Engineering, Stanford University  
<sup>2</sup>Apple <sup>3</sup>Blavatnik School of Computer Science, Tel Aviv University, and Google Research Tel Aviv. Correspondence to: Hilal Asi <asi@stanford.edu>.

In this problem (DP-SCO), given  $n$  i.i.d. samples  $z_1, \dots, z_n$  from a distribution  $P$ , we wish to release a private solution  $x \in \mathcal{X} \subseteq \mathbb{R}^d$  that minimizes the population loss  $F(x) = \mathbb{E}_{z \sim P}[f(x; z)]$  for a convex function  $f$  over  $x$ . The algorithm's performance is measured using the excess population loss of the solution  $x$ , that is  $F(x) - \min_{y \in \mathcal{X}} F(y)$ . The optimal algorithms and rates for this problem—even without privacy—have a crucial dependence on the geometry of the constraint set  $\mathcal{X}$  and in this work we focus on sets with bounded  $\ell_1$ -diameter. Without privacy constraints, there exist standard and efficient algorithms, such as mirror descent and exponentiated gradient descent, that achieve the optimal excess loss  $O(\sqrt{\log(d)/n})$  (Shalev-Shwartz & Ben-David, 2014). The landscape of the problem, however, with privacy constraints is not fully understood yet.

Most prior work on private convex optimization has focused on minimization of the empirical loss  $\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n f(x; z_i)$  over  $\ell_2$ -bounded domains (Chaudhuri et al., 2011; Bassily et al., 2014; 2019). Bassily et al. (2014) show that the optimal excess empirical loss in this setting is  $\Theta(\sqrt{d}/\varepsilon n)$  up to log factors. More recently, Bassily et al. (2019) give an asymptotically tight bound of  $1/\sqrt{n} + \sqrt{d}/(\varepsilon n)$  on the excess population loss in this setting using noisy gradient descent. Under mild smoothness assumptions, Feldman et al. (2020a) develop algorithms that achieve the optimal excess population loss using  $n$  gradient computations.

In contrast, existing results for private optimization in  $\ell_1$ -geometry do not achieve the optimal rates for the excess population loss (Kifer et al., 2012; Jain & Thakurta, 2014; Talwar et al., 2015). For the empirical loss, Talwar et al. (2015) develop private algorithms with  $\tilde{O}(1/(n\varepsilon)^{2/3})$  excess empirical loss for smooth functions and provide tight lower bounds when the dimension  $d$  is sufficiently high. These bounds can be converted into bounds on the excess population loss using standard techniques of uniform convergence of empirical loss to population loss, however these techniques can lead to suboptimal bounds as there are settings where uniform convergence is lower bounded by  $\Omega(\sqrt{d}/n)$  (Feldman, 2016). Moreover, the algorithm of Talwar et al. (2015) has runtime  $O(n^{5/3})$  in the moderate privacy regime ( $\varepsilon = \Theta(1)$ ) which is prohibitive in practice.

On the other hand, Jain & Thakurta (2014) develop algorithms for the population loss, however, their work is limited to generalized linear models and achieves a sub-optimal rate  $\tilde{O}(1/n^{1/3})$ .

In this work we develop private algorithms that achieve the optimal excess population loss in  $\ell_1$ -geometry, demonstrating that significant improvements are possible when the functions are smooth, in contrast to  $\ell_2$ -geometry where smoothness does not lead to better bounds. Specifically, for non-smooth functions, we develop an iterative localization algorithm, based on noisy mirror descent which achieves the optimal rate  $\sqrt{\log(d)/n} + \sqrt{d}/\varepsilon n$ . With additional smoothness assumptions, we show that rates with logarithmic dependence on the dimension are possible using a private variance-reduced Frank-Wolfe algorithm which obtains the rate  $\sqrt{\log(d)/n} + (\log(d)/\varepsilon n)^{2/3}$  and runs in linear (in  $n$ ) time. This shows that privacy is essentially free in this setting even when  $d \gg n$  and  $\varepsilon$  is as small as  $n^{-1/4}$ . Moreover, we show that similar rates are possible for general  $\ell_p$ -geometries for non-smooth functions when  $1 \leq p \leq 2$ . Finally, our algorithms query at most  $O(n^{3/2})$  gradients which improves over the best known algorithms for the non-smooth case in  $\ell_2$ -geometry which require  $n^2$  gradients (Feldman et al., 2020a).

The following two theorems summarize our upper bounds.

**Theorem 1** (non-smooth functions). *Let  $\mathcal{X} \subset \mathbb{R}^d$  be a convex body with  $\ell_1$  diameter less than 1. Let  $f(\cdot; z)$  be convex, Lipschitz with respect to  $\|\cdot\|_1$  for any  $z \in \mathcal{Z}$ . There is an  $(\varepsilon, \delta)$ -DP algorithm that takes a dataset  $S \in \mathcal{Z}^n$ , queries at most  $O(\log n \cdot \min(n^{3/2}\sqrt{\log d}, n^2\varepsilon/\sqrt{d}))$  and outputs a solution  $\hat{x}$  that has*

$$\mathbb{E}[F(\hat{x})] \leq \min_{x \in \mathcal{X}} F(x) + \tilde{O} \left( \sqrt{\frac{\log d}{n}} + \frac{\sqrt{d} \log^{3/2} d}{n\varepsilon} \right),$$

where the expectation is over the random choice of  $S$  and the randomness of the algorithm.

**Theorem 2** (smooth functions). *Let  $\mathcal{X} = \{x \in \mathbb{R}^d : \|x\|_1 \leq 1\}$  be the  $\ell_1$ -ball. Let  $f(\cdot; z)$  be convex, Lipschitz and smooth with respect to  $\|\cdot\|_1$  for any  $z \in \mathcal{Z}$ . There is an  $(\varepsilon, \delta)$ -DP linear time algorithm that takes a dataset  $S \in \mathcal{Z}^n$  and outputs a solution  $\hat{x}$  that has*

$$\mathbb{E}[F(x_K)] \leq \min_{x \in \mathcal{X}} F(x) + \tilde{O} \left( \sqrt{\frac{\log d}{n}} + \left( \frac{\log d}{n\varepsilon} \right)^{2/3} \right),$$

where the expectation is over the random choice of  $S$  and the randomness of the algorithm.

We also show how to improve these rates for strongly convex functions in Appendix C.

Before proceeding to review our algorithmic techniques, we briefly explain why the approaches used to obtain optimal rates in  $\ell_2$ -geometry (Bassily et al., 2019; Feldman et al., 2020a) do not work in our setting. One of the most natural approaches to proving bounds for private stochastic optimization is to use the generalization properties of differential privacy to derive population loss bounds for a private ERM algorithm. This approach fails to give asymptotically optimal bounds for the  $\ell_2$  case (Bassily et al., 2014), and similarly gives suboptimal bounds for the  $\ell_1$  case. Broadly, there are two approaches that have been used to get optimal bounds in the  $\ell_2$  case. An approach due to Bassily et al. (2019) uses stability of SGD on sufficiently smooth losses (Hardt et al., 2016) to get population loss bounds. These stability results rely on contractivity of gradient descent steps. However, as we show in an example that appears in Appendix A, the versions of mirror descent that are relevant to our setting do not have this property. Feldman et al. (2020a) derive generalization properties of their one pass algorithms from online-to-batch conversion. However, their analysis still relies on contractivity to prove the privacy guarantees of their algorithm. For their iterative localization approach Feldman et al. (2020a) use stability of the optimal solution to ERM in a different way to determine the scale of the noise added in each phase of the algorithm. In  $\ell_1$  geometry the norm of the noise added via this approach would overwhelm the signal (we discuss this in detail below).

We overview the key techniques we use to overcome these challenges below.

**Mirror descent based Iterative Localization.** In the non-smooth setting, we build on the iterative localization framework of Feldman et al. (2020a). In this framework in each phase a non-private optimization algorithm is used to solve a regularized version of the optimization problem. Regularization ensures that the output solution has small sensitivity and thus addition of Gaussian noise guarantees privacy. By appropriately choosing the noise and regularization scales, each phase reduces the distance to an approximate minimizer by a multiplicative factor. Thus after a logarithmic number of phases, the current iterate has the desired guarantees. Unfortunately, addition of Gaussian noise (and other output perturbation techniques) results in sub-optimal bounds in  $\ell_1$ -geometry since the  $\ell_1$ -error due to noise grows linearly with  $d$ . In contrast, the  $\ell_2$ -error grows as  $\sqrt{d}$ .

Instead of using output perturbation, we propose to use a private optimization algorithm in each phase. Using stability properties of strongly convex functions, we show that if the output of the private algorithm has sufficiently small empirical excess loss, then it has to be close to an approximate minimizer. Specifically, we reduce the distance to a minimizer by a multiplicative factor (relative to the initial

conditions at that phase). We show that a private version of mirror descent for strongly convex empirical risk minimization achieves sufficiently small excess empirical loss giving us an algorithm that achieves the optimal rate for non-smooth loss functions. More generally, this technique reduces the problem of DP-SCO to the problem of DP-ERM with strongly convex objectives. We provide details and analysis of this approach in Section 3.

**Dyadic variance-reduced Frank-Wolfe.** Our second algorithm is based on recent progress in stochastic optimization. Yurtsever et al. (2019) developed (non-private) variance-reduced Frank-Wolfe algorithm that achieves the optimal  $\tilde{O}(1/\sqrt{n})$  excess population loss improving on the standard implementations of Frank-Wolfe that achieve population loss  $\tilde{O}(1/n^{1/3})$ . The improvement relies on a novel variance reduction techniques that uses previous samples to improve the gradient estimates at future iterates (Fang et al., 2018). This frequent reuse of samples is the main challenge in developing a private version of the algorithm.

Inspired by the binary tree technique in the privacy literature (Dwork et al., 2010; 2015), we develop a new binary-tree-based variance reduction technique for the Frank-Wolfe algorithm. At a high level, the algorithm constructs a binary tree and allocates a set of samples to each vertex. The gradient at each vertex is then estimated using the samples of that vertex and the gradients along the path to the root. We assign more samples (larger batch sizes) to vertices that are closer to the root, to account for the fact that they are reused in more steps of the algorithm. This ensures that the privacy budget of samples in any vertex is not exceeded.

Using this privacy-aware design of variance-reduction, we rely on two tools to develop and analyze our algorithm. First, similarly to the private Frank-Wolfe for ERM (Talwar et al., 2015), we use the exponential mechanism to privatize the updates. A Frank-Wolfe update chooses one of the vertices of the constraint set ( $2d$  possibilities including signs for  $\ell_1$ -balls) and therefore the application of the exponential mechanism leads to a logarithmic dependence on the dimension  $d$ . This tool together with the careful accounting of privacy losses across the nodes, suffices to get the optimal bounds for the pure  $\varepsilon$ -DP case ( $\delta = 0$ ). To get the optimal rates for  $(\varepsilon, \delta)$ -DP, we rely on recent amplification by shuffling result for private local randomizers (Feldman et al., 2020b). To amplify privacy, we view our algorithm as a sequence of local randomizers, each operating on a different subset of the tree. Section 4 contains details of this algorithm.

In independent and concurrent work, Bassily et al. (2021) study differentially private algorithms for stochastic optimization in  $\ell_p$ -geometry. Similarly to our work, they build on mirror descent and variance-reduced Frank-Wolfe al-

gorithms to design private procedures for DP-SCO albeit without the iterative localization scheme and the binary-tree-based sample allocation technique we propose. As a result, their algorithms achieve sub-optimal rates in some of the parameter regimes: in  $\ell_1$ -geometry, they achieve excess loss of roughly  $\log(d)/\varepsilon\sqrt{n}$  in contrast to the  $\sqrt{\log(d)}/\sqrt{n} + \log(d)/(\varepsilon n)^{2/3}$  rate of our algorithms. For  $1 < p < 2$ , their algorithms have excess loss of (up to log factors)  $\min(d^{1/4}/\sqrt{n}, \sqrt{d}/(\varepsilon n^{3/4}))$ , whereas our algorithms achieve the rate of  $\sqrt{d}/\varepsilon n$ . On the other hand, Bassily et al. (2021) develop a generalized Gaussian mechanism for adding noise in  $\ell_p$ -geometry. Their mechanism improves over the standard Gaussian mechanism and can improve the rates of our algorithms for  $\ell_p$ -geometry (Theorem 5) by a  $\sqrt{\log d}$  factor. Moreover, they prove a lower bound for  $\ell_p$ -geometries with  $1 < p < 2$  that establishes the optimality of our upper bounds for  $1 < p < 2$ .

## 2. Preliminaries

### 2.1. Stochastic Convex Optimization

We let  $\mathcal{S} = (z_1, \dots, z_n)$  denote datasets where  $z_i \in \mathcal{Z}$  are drawn i.i.d. from a distribution  $P$  over the domain  $\mathcal{Z}$ . Let  $\mathcal{X} \subseteq \mathbb{R}^d$  be a convex set that denotes the set of parameters for the optimization problem. Given a loss function  $f(x; z) : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}$  that is convex in  $x$  (for every  $z$ ), we define the population loss  $F(x) = \mathbb{E}_{z \sim P}[f(x; z)]$ . The excess population loss of a parameter  $x \in \mathcal{X}$  is then  $F(x) - \min_{y \in \mathcal{X}} F(y)$ . We also consider the empirical loss  $\hat{F}(x; S) = \frac{1}{n} \sum_{i=1}^n f(x; z_i)$  and the excess empirical loss of  $x \in \mathcal{X}$  is  $\hat{F}(x; S) - \min_{y \in \mathcal{X}} \hat{F}(y; S)$ . For a set  $\mathcal{X}$ , we will denote its  $\ell_p$  diameter by  $\text{diam}_p(\mathcal{X}) = \sup_{x, y \in \mathcal{X}} \|x - y\|_p$ .

As we are interested in general geometries, we define the standard properties (e.g., Lipschitz, smooth and strongly convex) with respect to a general norm which are frequently used in the optimization literature (Duchi, 2018).

**Definition 2.1** (Lipschitz continuity). *A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is  $L$ -Lipschitz with respect to a norm  $\|\cdot\|$  over  $\mathcal{X}$  if for every  $x, y \in \mathcal{X}$  we have  $|f(x) - f(y)| \leq L \|x - y\|$ .*

A standard result is that  $L$ -Lipschitz continuity is equivalent to bounded (sub)-gradients, namely that  $\|g\|_* \leq L$  for all  $x \in \mathcal{X}$  and sub-gradient  $g \in \partial f(x)$  where  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$ .

**Definition 2.2** (smoothness). *A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is  $\beta$ -smooth with respect to a norm  $\|\cdot\|$  over  $\mathcal{X}$  if for every  $x, y \in \mathcal{X}$  we have  $\|\nabla f(x) - \nabla f(y)\|_* \leq \beta \|x - y\|$ .*

**Definition 2.3** (strong convexity). *A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is  $\lambda$ -strongly convex with respect to a norm  $\|\cdot\|$  over  $\mathcal{X}$  if for any  $x, y \in \mathcal{X}$  we have  $f(x) + \langle \nabla f(x), y - x \rangle + \frac{\lambda}{2} \|y - x\|^2 \leq f(y)$ .*

Since we develop private versions of mirror descent, we define the Bregman divergence associated with a differentiable convex function  $h : \mathcal{X} \rightarrow \mathbb{R}$  to be  $D_h(x, y) = h(x) - h(y) - \langle \nabla h(y), x - y \rangle$ . We require a definition of strong convexity relative to a function which has been used in several works in the optimization literature (Duchi et al., 2010; Lu et al., 2018).

**Definition 2.4** (relative strong convexity). *A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is  $\lambda$ -strongly convex relative to  $h : \mathcal{X} \rightarrow \mathbb{R}$  if for any  $x, y \in \mathcal{X}$ ,  $f(x) + \langle \nabla f(x), y - x \rangle + \lambda D_h(y, x) \leq f(y)$ .*

Note that if  $h(x)$  is convex, then  $h(x)$  is 1-strongly convex relative to  $h(x)$  according to this definition. Moreover, the function  $f(x) = g(x) + h(x)$  is also 1-strongly convex relative to  $h(x)$  for any convex function  $g(x)$ .

## 2.2. Differential Privacy

We recall the definition of  $(\varepsilon, \delta)$ -differential privacy.

**Definition 2.5** (Dwork et al., 2006b;a). *A randomized algorithm  $\mathcal{A}$  is  $(\varepsilon, \delta)$ -differentially private  $((\varepsilon, \delta)$ -DP) if, for all datasets  $\mathcal{S}, \mathcal{S}' \in \mathcal{Z}^n$  that differ in a single data element and for all events  $\mathcal{O}$  in the output space of  $\mathcal{A}$ , we have*

$$\Pr[\mathcal{A}(\mathcal{S}) \in \mathcal{O}] \leq e^\varepsilon \Pr[\mathcal{A}(\mathcal{S}') \in \mathcal{O}] + \delta.$$

When  $\delta = 0$ , we use the shorter notation  $\varepsilon$ -DP. We also use the following privacy composition results.

**Lemma 2.1** (Basic composition Dwork & Roth, 2014). *If  $\mathcal{A}_1, \dots, \mathcal{A}_k$  are randomized algorithms that each is  $\varepsilon$ -DP, then their composition  $(\mathcal{A}_1(\mathcal{S}), \dots, \mathcal{A}_k(\mathcal{S}))$  is  $k\varepsilon$ -DP.*

**Lemma 2.2** (Advanced composition Dwork & Roth, 2014). *If  $\mathcal{A}_1, \dots, \mathcal{A}_k$  are randomized algorithms that each is  $(\varepsilon, \delta)$ -DP, then their composition  $(\mathcal{A}_1(\mathcal{S}), \dots, \mathcal{A}_k(\mathcal{S}))$  is  $(\sqrt{2k \log(1/\delta')}\varepsilon + k\varepsilon(e^\varepsilon - 1), \delta' + k\delta)$ -DP.*

## 3. Algorithms for Non-Smooth Functions

In this section, we develop an algorithm that builds on the iterative localization techniques of Feldman et al. (2020a) to achieve optimal excess population loss for non-smooth functions over the  $\ell_1$ -ball. Instead of using output perturbation to solve the regularized optimization problems, our algorithm uses general private algorithms for solving strongly convex ERM problems. This essentially reduces the problem of privately minimizing the population loss to that of privately minimizing a strongly convex empirical risk. In Section 3.1 we develop private versions of mirror descent that achieve optimal bounds for strongly convex ERM problems, and in Section 3.2 we use these algorithms in an iterative localization framework to obtain optimal bounds for the population loss.

### 3.1. Private Algorithms for Strongly Convex ERM

In this section, we consider empirical risk minimization for strongly convex functions and achieve optimal excess empirical loss using noisy mirror descent (Algorithm 1).

---

#### Algorithm 1 Noisy Mirror Descent

---

**Require:** Dataset  $\mathcal{S} = (z_1, \dots, z_n) \in \mathcal{Z}^n$ , convex set  $\mathcal{X}$ , convex function  $h : \mathcal{X} \rightarrow \mathbb{R}$ , step sizes  $\{\eta_k\}_{k=1}^T$ , batch size  $b$ , initial point  $x_0$ , number of iterations  $T$ ;

- 1: **for**  $k = 1$  to  $T$  **do**
  - 2:   Sample  $S_1, \dots, S_b \sim \text{Unif}(\mathcal{S})$
  - 3:   Set  $\hat{g}_k = \frac{1}{b} \sum_{i=1}^b \nabla f(x_k; S_i) + \zeta_i$  where  $\zeta_i \sim \mathcal{N}(0, \sigma^2 I_d)$  with  $\sigma = 100L\sqrt{d \log(1/\delta)}/b\varepsilon$
  - 4:   Find  $x_{k+1} := \operatorname{argmin}_{x \in \mathcal{X}} \{ \langle \hat{g}_k, x - x_k \rangle + \frac{1}{\eta_k} D_h(x, x_k) \}$
  - 5: **end for**
  - 6: **return**  $\bar{x}_T = \frac{1}{T} \sum_{k=1}^T x_k$  (convex)
  - 7: **return**  $\hat{x}_T = \frac{2}{T(T+1)} \sum_{k=1}^T k x_k$  (strongly convex)
- 

**Theorem 3.** *Let  $h : \mathcal{X} \rightarrow \mathbb{R}$  be 1-strongly convex with respect to  $\|\cdot\|_1$ ,  $x^* = \operatorname{argmin}_{x \in \mathcal{X}} \hat{F}(x; \mathcal{S})$ , and assume  $D_h(x^*, x_0) \leq D^2$ . Let  $f(x; z)$  be convex and  $L$ -Lipschitz with respect to  $\|\cdot\|_1$  for all  $z \in \mathcal{Z}$ . Setting  $1 \leq b, T = \frac{n^2}{b^2}$  and  $\eta_k = \frac{D}{\sqrt{T}} \frac{1}{\sqrt{L^2 + 2\sigma^2 \log d}}$ , Algorithm 1 is  $(\varepsilon, \delta)$ -DP and*

$$\mathbb{E}[\hat{F}(\bar{x}_T; \mathcal{S}) - \hat{F}(x^*; \mathcal{S})] \leq LD \cdot \mathcal{O} \left( \frac{b}{n} + \frac{\sqrt{d \log d \log \frac{1}{\delta}}}{n\varepsilon} \right).$$

Moreover, if  $f(x; z)$  is  $\lambda$ -strongly convex relative to  $h(x)$ , then setting  $\eta_k = \frac{2}{\lambda(k+1)}$

$$\mathbb{E}[\hat{F}(\hat{x}_T; \mathcal{S}) - \hat{F}(x^*; \mathcal{S})] \leq \mathcal{O} \left( \frac{L^2 b^2}{\lambda n^2} + \frac{L^2 d \log d \log \frac{1}{\delta}}{\lambda n^2 \varepsilon^2} \right).$$

To prove Theorem 3, we need the following standard results for the convergence of stochastic mirror descent for convex and strongly convex functions.

**Lemma 3.1** (Duchi, 2018, Corollary 4.2.11). *Assume  $h(x)$  is 1-strongly convex with respect to  $\|\cdot\|_1$ . Let  $f(x)$  be a convex function and  $x^* = \operatorname{argmin}_{x \in \mathcal{X}} f(x)$ . Consider the stochastic mirror descent update  $x_{k+1} = \operatorname{argmin}_{x \in \mathcal{X}} \{ \langle g_k, x - x_k \rangle + \frac{1}{\eta_k} D_h(x, x_k) \}$  where  $\mathbb{E}[g_k] \in \partial f(x_k)$  with  $\mathbb{E}[\|g_k\|_\infty^2] \leq L^2$ . If  $\eta_k = \eta$  for all  $k$  then the average iterate  $\bar{x}_T = \frac{1}{T} \sum_{i=1}^T x_i$  has  $\mathbb{E}[f(\bar{x}_T) - f(x^*)] \leq \frac{D_h(x^*, x_1)}{T\eta} + \frac{\eta L^2}{2}$ .*

We also need the following result which states the rates of stochastic mirror descent for strongly convex functions. Similar results appear in the optimization literature (Lacoste-Julien et al., 2012), though as the statement we require is less common, we provide a proof in Appendix D.2.

**Lemma 3.2.** *Under the same notation of Lemma 3.1, if  $f(x)$  is  $\lambda$ -strongly convex relative to  $h(x)$ , then setting  $\eta_k = \frac{2}{\lambda(k+1)}$  the weighted average  $\hat{x}_T = \frac{2}{T(T+1)} \sum_{k=1}^T kx_k$  has  $\mathbb{E}[f(\hat{x}_T) - f(x^*)] \leq \frac{L^2}{\lambda(T+1)}$ .*

We are now ready to prove Theorem 3.

*Proof.* The privacy guarantees follow from standard properties of the Gaussian mechanism and Moments accountant (Abadi et al., 2016); we provide full details in Appendix D. Now we prove the utility of the algorithm. To this end, we have that  $\mathbb{E}[\|\hat{g}_k\|_\infty^2] \leq 2L^2 + 2\mathbb{E}[\|\zeta_k\|_\infty^2] \leq 2L^2 + 4\sigma^2 \log d$ . Lemma 3.1 now implies that

$$\begin{aligned} \mathbb{E}[\hat{F}(\bar{x}_T; S) - \hat{F}(x^*; S)] &\leq \frac{D^2}{T\eta} + \eta L^2 + 2\eta\sigma^2 \log d \\ &\leq 2D\sqrt{(L^2 + 2\sigma^2 \log d)/T} \\ &\leq LD \cdot O\left(\frac{b}{n} + \frac{\sqrt{d \log d \log \frac{1}{\delta}}}{n\varepsilon}\right), \end{aligned}$$

where the second inequality follows from the choice of  $\eta$ . For the second part, Lemma 3.2 implies that

$$\mathbb{E}[\hat{F}(\hat{x}_T; S) - \hat{F}(x^*; S)] \leq \frac{L^2}{\lambda} O\left(\frac{b^2}{n^2} + \frac{d \log d \log \frac{1}{\delta}}{n^2 \varepsilon^2}\right). \quad \square$$

### 3.2. Private Algorithms for SCO

Building on the noisy mirror descent algorithm of Section 3.1, in this section we develop a localization based algorithm for the population loss that achieves the optimal bounds in  $\ell_1$  geometry. The algorithm iteratively solves a regularized version of the (empirical) objective function using noisy mirror decent (Algorithm 1). The output of each iterate is accurate enough to allow to shrink the diameter of the domain at the next iterate (increasing regularization), hence making the optimization problem easier.

We present the full details in Algorithm 2 which enjoys the following guarantees.

**Theorem 4.** *Assume  $\text{diam}_1(\mathcal{X}) \leq D$  and  $f(x; z)$  is convex and  $L$ -Lipschitz with respect to  $\|\cdot\|_1$  for all  $z \in \mathcal{Z}$ . If we set*

$$\eta = \frac{D}{L} \min \left\{ \sqrt{\log(d)/n}, \varepsilon / \sqrt{d \log d \log \frac{1}{\delta}} \right\},$$

*then Algorithm 2 is  $(\varepsilon, \delta)$ -DP, uses  $O(\log n \cdot \min(n^{3/2} \sqrt{\log d}, n^2 \varepsilon / \sqrt{d}))$  gradients and its output has*

$$\mathbb{E}[F(x_k) - F(x^*)] = LD \cdot O\left(\frac{\sqrt{\log d}}{\sqrt{n}} + \frac{\sqrt{d \log^3 d \log \frac{1}{\delta}}}{n\varepsilon}\right).$$

### Algorithm 2 Localized Noisy Mirror Descent

**Require:** Dataset  $\mathcal{S} = (z_1, \dots, z_n) \in \mathcal{Z}^n$ , constraint set  $\mathcal{X}$ , step size  $\eta$ , initial point  $x_0$ ;

- 1: Set  $k = \lceil \log n \rceil$ ,  $p = 1 + 1/\log d$
- 2: **for**  $i = 1$  to  $k$  **do**
- 3:   Set  $n_i = 2^{-i}n$ ,  $\eta_i = 2^{-4i}\eta$
- 4:   Apply Algorithm 1 with  $(\varepsilon, \delta)$ -DP, batch size  $b_i = \max(\sqrt{n_i/\log d}, \sqrt{d/\varepsilon})$ ,  $T = n_i^2/b_i^2$  and  $h_i(x) = \frac{1}{p-1} \|x - x_{i-1}\|_p^2$  for solving the ERM over  $\mathcal{X}_i = \{x \in \mathcal{X} : \|x - x_{i-1}\|_p \leq 2L\eta_i n_i(p-1)\}$ :  

$$F_i(x) = \frac{1}{n_i} \sum_{j=1}^{n_i} f(x; z_j) + \frac{\|x - x_{i-1}\|_p^2}{\eta_i n_i (p-1)}$$
- 5:   Let  $x_i$  be the output of the private algorithm
- 6: **end for**
- 7: **return** the final iterate  $x_k$

We begin with the following lemma which bounds the distance of the private minimizer to the true minimizer at each iteration.

**Lemma 3.3.** *Let  $\hat{x}_i = \arg\min_{x \in \mathcal{X}} F_i(x)$ . Then ,*

$$\mathbb{E}[\|x_i - \hat{x}_i\|_p^2] \leq O\left(\frac{L^2 \eta_i^2 n_i}{\log d} + L^2 \eta_i^2 d \log d \log(1/\delta)/\varepsilon^2\right).$$

The next lemma follows from Shalev-Shwartz et al. (2009).

**Lemma 3.4.** *Let  $\hat{x}_i = \arg\min_{x \in \mathcal{X}_i} F_i(x)$  and  $y \in \mathcal{X}$ . If  $f(x; z)$  is  $L$ -Lipschitz with respect to  $\|\cdot\|_1$ , then  $\mathbb{E}[F(\hat{x}_i)] - F(y) \leq \frac{\mathbb{E}[\|y - x_{i-1}\|_p^2]}{\eta_i n_i (p-1)} + O(L^2 \eta_i)$ .*

We are now ready to prove Theorem 4.

*Proof.* First, note that the algorithm is  $(\varepsilon, \delta)$ -DP as each sample is used in one iterate, hence the privacy claim follows from the guarantees of Algorithm 1 and post-processing. Now we prove the claim about the number of queried gradients. Algorithm 1 requires  $n_i^2/b_i$  gradients hence since  $b_i = \max(\sqrt{n_i/\log d}, \sqrt{d/\varepsilon})$  we get that the number of gradients at each stage is at most  $\min(n^{3/2} \sqrt{\log d}, n^2 \varepsilon / \sqrt{d})$ , implying the claim as we have  $\log n$  iterates. Next, we prove utility which is similar to the proof of Theorem 4.4 in (Feldman et al., 2020a). Letting  $\hat{x}_0 = x^*$ , we have:

$$\begin{aligned} \mathbb{E}[F(x_k)] - F(x^*) &= \sum_{i=1}^k \mathbb{E}[F(\hat{x}_i) - F(\hat{x}_{i-1})] + \mathbb{E}[F(x_k) - F(\hat{x}_k)]. \end{aligned}$$

First, note that Lemma 3.3 implies

$$\begin{aligned}
 \mathbb{E}[F(x_k) - F(\hat{x}_k)] &\leq L\mathbb{E}[\|x_k - \hat{x}_k\|_1] \\
 &\leq L\sqrt{\mathbb{E}[2\|x_k - \hat{x}_k\|_p^2]} \\
 &\leq CL^2\eta_k(\sqrt{n_i/\log d} + \sqrt{d\log d\log(1/\delta)}/\varepsilon) \\
 &\leq C2^{-2k}L^2\eta(\sqrt{n/\log d} + \sqrt{d\log d\log(1/\delta)}/\varepsilon) \\
 &\leq CLD/n^2,
 \end{aligned}$$

where the last inequality follows since  $\eta \leq \frac{D}{L} \min(\sqrt{\log(d)/n}, \varepsilon/\sqrt{d\log d\log(1/\delta)})$ . Lemmas 3.4 and 3.3 imply

$$\begin{aligned}
 &\sum_{i=1}^k \mathbb{E}[F(\hat{x}_i) - F(\hat{x}_{i-1})] \\
 &\leq \sum_{i=1}^k \frac{\mathbb{E}[\|\hat{x}_{i-1} - x_{i-1}\|_p^2]}{\eta_i n_i (p-1)} + CL^2\eta_i \\
 &\leq \frac{D^2}{\eta n (p-1)} + \sum_{i=2}^k C(2L^2\eta_i + \frac{L^2\eta_i d\log d\log(1/\delta)}{n_i \varepsilon^2 (p-1)}) \\
 &\leq \frac{D^2}{\eta n (p-1)} + 4CL^2\eta + C \sum_{i=2}^k 2^{-i} \frac{L^2\eta d\log d\log(1/\delta)}{n \varepsilon^2 (p-1)} \\
 &\leq \frac{D^2}{\eta n (p-1)} + 2C \frac{L^2\eta d\log d\log(1/\delta)}{n \varepsilon^2 (p-1)} + 4CL^2\eta.
 \end{aligned}$$

The claim now follows by setting the value of  $\eta$ .  $\square$

Finally, we can extend Algorithm 2 to work for general  $\ell_p$  geometries for  $1 < p \leq 2$ , resulting in the following theorem. We defer full details to Appendix B.

**Theorem 5.** *Let  $1 < p \leq 2$ . Assume  $\text{diam}_p(\mathcal{X}) \leq D$  and  $f(x; z)$  is convex and  $L$ -Lipschitz with respect to  $\|\cdot\|_p$  for all  $z \in \mathcal{Z}$ . Then there is an  $(\varepsilon, \delta)$ -DP algorithm that uses  $O(\log n \cdot \min(n^{3/2}\sqrt{\log d}, n^2\varepsilon/\sqrt{d}))$  and outputs  $\hat{x}$  such that*

$$\mathbb{E}[F(\hat{x}) - F(x^*)] = LD \cdot O\left(\frac{1}{\sqrt{(p-1)n}} + \frac{\sqrt{d\log d\log \frac{1}{\delta}}}{(p-1)n\varepsilon}\right).$$

If  $p = 2$  then the output  $\hat{x}$  has

$$\mathbb{E}[F(\hat{x}) - F(x^*)] = LD \cdot O\left(\frac{1}{\sqrt{n}} + \frac{\sqrt{d\log \frac{1}{\delta}}}{n\varepsilon}\right).$$

#### 4. Efficient Algorithms for Smooth Functions

Having established tight bounds for the non-smooth case, in this section we turn to the smooth setting and develop

linear-time private Frank-Wolfe algorithms with variance-reduction that achieve the optimal rates. Specifically, our algorithms achieve the rate  $\tilde{O}(1/\sqrt{n\varepsilon})$  for pure  $\varepsilon$ -DP and  $\tilde{O}(1/\sqrt{n} + 1/(n\varepsilon)^{2/3})$  for  $(\varepsilon, \delta)$ -DP. These results imply that the optimal (non-private) statistical rate  $\tilde{O}(1/\sqrt{n})$  is achievable with strong privacy guarantees—whenever  $\varepsilon \geq \tilde{\Omega}(1/n^{1/4})$  for  $(\varepsilon, \delta)$ -DP—even for high dimensional functions with  $d \gg n$ .

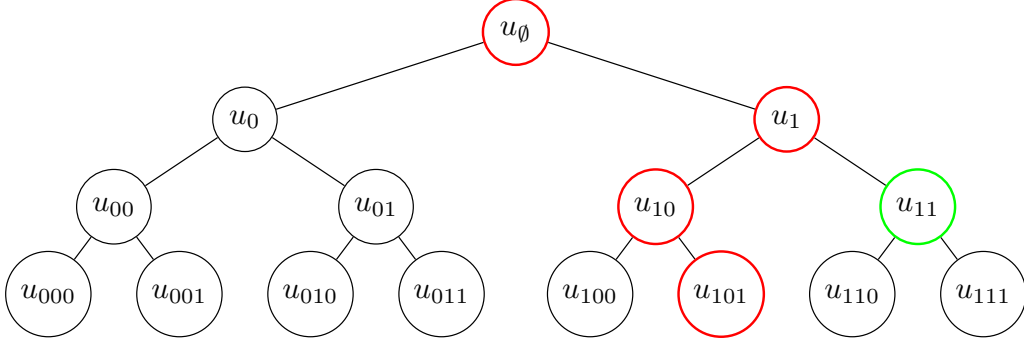
The starting point of our algorithms is the recent non-private Frank-Wolfe algorithm of Yurtsever et al. (2019) which uses variance-reduction techniques to achieve the (non-private) optimal rates. Due to the high reuse of samples, a direct approach to privatizing their algorithm would result in sub-optimal bounds. To overcome this, we design a new binary-tree scheme for variance reduction that allows for more noise-efficient private algorithms.

We describe our private Frank-Wolfe procedure in Algorithm 3. We present the algorithm in a more general setting where  $\mathcal{X}$  can be an arbitrary convex body with  $m$  vertices. The algorithm has  $T$  phases (outer iterations) indexed by  $1 \leq t \leq T$  and each phase  $t$  has a binary tree of depth  $t$ . We will denote vertices by  $u_s$  where  $s \in \{0, 1\}^{\leq t}$  is the path to the vertex; i.e.,  $u_\emptyset$  will denote the root of the tree,  $u_{01}$  will denote the right child of  $u_\emptyset$ . Each vertex  $u_s$  is associated with a parameter  $x_{t,s}$ , a gradient estimate  $v_{t,s}$ , and a set of samples  $S_{t,s}$  of size  $2^{-j}b$  where  $j$  is the depth of the vertex. Roughly, the idea is to improve the gradient estimate at a vertex (reduce the variance) using the gradient estimates at vertices along the path to the root. Crucially, the large sample size at vertices with smaller depth allows these vertices to apply gradient corrections for their relatively large sub-trees with mild privacy cost.

More precisely, the algorithm traverses through the graph vertices according to the Depth-First-Search (DFS) approach. At each vertex, the algorithm improves the gradient estimate at the current vertex using the estimate at the parent vertex. When the algorithm visits a leaf vertex, it also updates the current iterate using the Frank-Wolfe step with the gradient estimate at the leaf.

For notational convenience, we let  $\text{DFS}(t)$  denote the DFS order of the vertices in a binary tree of depth  $t$  (root not included), i.e., for  $t = 2$  we have  $\text{DFS}(t) = \{u_0, u_{00}, u_{01}, u_1, u_{10}, u_{11}\}$ . Moreover, for  $s \in \{0, 1\}^t$  we let  $\ell(s)$  denote the integer whose binary representation is  $s$ . In the description of the algorithm, we denote iterates by  $x_{t,s}$  where  $t$  is the phase and  $s \in \{0, 1\}^t$  is the path from the root. In our proofs, we sometimes use the equivalent notation  $x_k$  where  $k = 2^{t-1} + \ell(s)$ .

We analyze Algorithm 3 for pure and approximate DP.



**Figure 1.** Binary tree at phase  $t = 3$  of the algorithm. At the leaf  $u_{101}$ , the algorithm has the gradient estimate  $v_{t,101}$  which is calculated along the path to the root where every right son applies a correction step to the estimate. Using the gradient estimate  $v_{t,101}$ , the algorithm applies a Frank-Wolfe step to calculate the next iterate and put its value in the next DFS vertex, namely  $u_{t,11}$ .

---

**Algorithm 3** Private Variance Reduced Frank-Wolfe

**Require:** Dataset  $\mathcal{S} = (z_1, \dots, z_n) \in \mathcal{Z}^n$ , constraint set  $\mathcal{X} = \text{conv}\{c_1, \dots, c_m\}$ , number of phases  $T$ , batch size  $b$ , initial point  $x_0$ ;

- 1: **for**  $t = 1$  to  $T$  **do**
  - 2:   Set  $x_{t,\emptyset} = x_{t-1,L_{t-1}}$
  - 3:   Draw  $b$  samples to the set  $S_{t,\emptyset}$
  - 4:    $v_{t,\emptyset} \leftarrow \nabla f(x_{t,\emptyset}; S_{t,\emptyset})$
  - 5:   **for**  $u_s \in \text{DFS}[2^t]$  **do**
  - 6:     Let  $s = s'c$  where  $c \in \{0, 1\}$  and  $j = |s|$
  - 7:     **if**  $c = 0$  **then**
  - 8:        $v_{t,s} \leftarrow v_{t,s'}$ ;  $x_{t,s} \leftarrow x_{t,s'}$
  - 9:     **else**
  - 10:       Draw  $2^{-j}b$  samples to the set  $S_{t,s}$
  - 11:        $v_{t,s} \leftarrow v_{t,s'} + \nabla f(x_{t,s}; S_{t,s}) - \nabla f(x_{t,s'}; S_{t,s})$
  - 12:     **end if**
  - 13:     **if**  $j = t$  **then**
  - 14:       Let  $s_+$  be the next vertex in the DFS iteration
  - 15:        $w_{t,s} \leftarrow \text{argmin}_{c_i: 1 \leq i \leq m} \langle c_i, v_{t,s} \rangle + \zeta_i$  where  $\zeta_i \sim \text{Laplace}(\lambda_{t,s})$
  - 16:        $x_{t,s_+} \leftarrow (1 - \eta_{t,s})x_{t,s} + \eta_{t,s}w_{t,s}$  where  $\eta_{t,s} = \frac{2}{2^{t-1} + \ell(s) + 1}$
  - 17:     **end if**
  - 18:   **end for**
  - 19: **end for**
  - 20: **return** the final iterate  $x_K$
- 

#### 4.1. Pure Differential Privacy

The following theorem summarizes our guarantees for pure privacy. We defer missing proofs to Appendix E.

**Theorem 6.** Assume that  $\text{diam}_1(\mathcal{X}) \leq D$ ,  $m \leq O(d)$  and that  $f(x; z)$  is convex,  $L$ -Lipschitz and  $\beta$ -smooth with respect to  $\|\cdot\|_1$ . Assume also that  $\frac{L \log m \log^2 n}{n \varepsilon D} \leq \beta \leq \frac{nL \log m}{\varepsilon D \log^2 n}$ . Setting  $b = n / \log^2 n$ ,  $\lambda_{t,s} = \frac{2LD2^t}{b\varepsilon}$  and  $T = \frac{1}{2} \log \left( \frac{b\varepsilon\beta D}{L \log m} \right)$ , Algorithm 3 is  $\varepsilon$ -DP, queries  $n$  gradi-

ents, and has

$$\begin{aligned} \mathbb{E}[F(x_K) - F(x^*)] \\ \leq O \left( D(L + \beta D) \frac{\sqrt{\log d \log n}}{\sqrt{n}} + \frac{\sqrt{\beta L D^3 \log d \log n}}{\sqrt{n\varepsilon}} \right). \end{aligned}$$

Moreover, if  $\beta \leq \frac{L \log m \log^2 n}{n \varepsilon D}$  then setting  $T = 1$  and  $b = n$ , Algorithm 3 is  $\varepsilon$ -DP, queries  $n$  gradients, and has

$$\mathbb{E}[F(x_K) - F(x^*)] \leq DL \cdot O \left( \frac{\sqrt{\log d}}{\sqrt{n}} + \frac{\log d}{n\varepsilon} \right) + O(\beta D^2).$$

To prove the theorem, we begin with the following lemma that gives pure privacy guarantees.

**Lemma 4.1.** Assume  $2^T \leq b$ . Setting  $\lambda_{t,s} = \frac{2LD2^t}{b\varepsilon}$ , Algorithm 3 is  $\varepsilon$ -DP with  $\varepsilon \leq 1$ . Moreover,  $\mathbb{E}[\langle v_{t,s}, w_{t,s} \rangle] \leq \mathbb{E}[\min_{w \in \mathcal{X}} \langle v_{t,s}, w \rangle] + O(\frac{LD2^t}{b\varepsilon} \log m)$ .

*Proof.* (sketch) Let  $\mathcal{S} = (z_1, \dots, z_{n-1}, z_n)$  and  $\mathcal{S}' = (z_1, \dots, z_{n-1}, z'_n)$  be two neighboring datasets. We prove privacy given that the  $n$ 'th sample belongs to the set  $S_{t,s}$ , which will imply our general privacy guarantee as this holds for every choice of  $t$  and  $s$ . The main idea is that each sample in the set  $S_{t,s}$  is (directly) involved in the calculation of a Frank-Wolfe update at most  $N_{t,s} = 2^{t-|s|}$  times. Hence, setting the noise level  $\lambda_{t,s}$  large enough to guarantee that each iterate is  $\varepsilon/N_{t,s}$ -DP, basic composition implies the final output is  $\varepsilon$ -DP.  $\square$

The next lemma upper bounds the variance of the gradients.

**Lemma 4.2.** At the vertex  $(t, s)$ , we have

$$\mathbb{E} \|v_{t,s} - \nabla F(x_{t,s})\|_\infty \leq (L + \beta D) \cdot O \left( \sqrt{\log(d)/b} \right).$$

Using the previous two lemmas, we can prove Theorem 6.

*Proof.* The setting of the parameters and the condition on  $\beta$  ensures that  $2^T \leq b$  hence Lemma 4.1 implies the claim about privacy. Now we proceed to prove utility. In this proof, we use the equivalent representation  $k = 2^{t-1} + \ell(s)$  for a leaf vertex  $(t, s)$  where  $\ell(s)$  is the number whose binary representation is  $s$ . By smoothness we get,

$$\begin{aligned} & F(x_{k+1}) \\ & \leq F(x_k) + \langle \nabla F(x_k), x_{k+1} - x_k \rangle + \beta \|x_{k+1} - x_k\|_1^2 / 2 \\ & \leq F(x_k) + \eta_k \langle \nabla F(x_k), w_k - x_k \rangle + \beta \eta_k^2 D^2 / 2 \\ & = F(x_k) + \eta_k \langle \nabla F(x_k), x^* - x_k \rangle + \eta_k \langle v_k, w_k - x^* \rangle \\ & \quad + \eta_k \langle \nabla F(x_k) - v_k, w_k - x^* \rangle + \beta \eta_k^2 D^2 / 2 \\ & \leq F(x_k) + \eta_k (F(x^*) - F(x_k)) + \eta_k D \|\nabla F(x_k) - v_k\|_\infty \\ & \quad + \eta_k (\langle v_k, w_k \rangle - \min_{w \in \mathcal{X}} \langle v_k, w \rangle) + \beta \eta_k^2 D^2 / 2. \end{aligned}$$

Subtracting  $F(x^*)$  from each side, using Lemmas 4.1 and 4.2 and taking expectations, we have

$$\begin{aligned} & \mathbb{E}[F(x_{k+1}) - F(x^*)] \\ & \leq (1 - \eta_k) \mathbb{E}[F(x_k) - F(x^*)] + \eta_k D(L + \beta D) \sqrt{\frac{\log d}{b}} \\ & \quad + \frac{\eta_k^2}{2} \beta D^2 + \eta_k D L \frac{2^t \log m}{b\varepsilon}. \end{aligned}$$

Letting  $\alpha_k = \eta_k D(L + \beta D) \sqrt{\frac{\log d}{b}} + \frac{\eta_k^2}{2} \beta D^2 + \eta_k D L \frac{2^t \log m}{b\varepsilon}$ , we have

$$\begin{aligned} \mathbb{E}[F(x_K) - F(x^*)] & \leq \sum_{k=1}^K \alpha_k \prod_{i>k} (1 - \eta_i) \\ & = \sum_{k=1}^K \alpha_k \frac{(k-1)k}{K(K+1)} \leq \sum_{k=1}^K \alpha_k \frac{k^2}{K^2}. \end{aligned}$$

Since  $t \leq T$  and  $K = 2^T$ , simple algebra now yields

$$\begin{aligned} & \mathbb{E}[F(x_K) - F(x^*)] \\ & \leq O\left(D(L + \beta D) \frac{\sqrt{\log d}}{\sqrt{b}} + \frac{\beta D^2}{2^T} + DL \frac{2^T \log m}{b\varepsilon}\right). \end{aligned}$$

The number of samples in the algorithm is upper bounded by  $T^2 \cdot b$  hence the first part of the claim follows by setting  $b = n / \log^2 n$  and  $T = \frac{1}{2} \log\left(\frac{b\varepsilon\beta D}{L \log m}\right)$ . The condition on  $\beta$  ensures that the term inside the log is greater than 1. The second part follows similarly using  $T = 1$  and  $b = n$ .  $\square$

## 4.2. Approximate Differential Privacy

The previous section achieves the optimal non-private rate  $1/\sqrt{n}$  only for  $\varepsilon = \Theta(1)$ . In this section we show that for approximate differential privacy, it is possible to achieve the

optimal rates when  $\varepsilon \geq \Omega(n^{-1/4})$ . The first approach to improve the privacy analysis is to use advanced composition for approximate DP. Unfortunately, it is not enough by itself and we use amplification by shuffling results to achieve the optimal bounds. The following theorem summarizes the guarantees of Algorithm 3 for approximate privacy.

**Theorem 7.** *Let  $\delta \leq 1/n$  and assume that  $\text{diam}_1(\mathcal{X}) \leq D$ ,  $m \leq O(d)$  and that  $f(x; z)$  is convex,  $L$ -Lipschitz and  $\beta$ -smooth with respect to  $\|\cdot\|_1$ . Assume  $\frac{L \log(n/\delta) \log m \log^2 n}{n\varepsilon D} \leq \beta \leq \frac{\sqrt{n} L \log(n/\delta) \log m}{\varepsilon D \log n}$  and  $\varepsilon \leq \frac{(L \log(n/\delta) \log m)^{1/4} \sqrt{\log n}}{(n\beta D)^{1/4}}$ .*

*Let  $\lambda_{t,s} = \frac{LD2^{T/2} \log(n/\delta)}{b\varepsilon}$ ,  $b = n / \log^2 n$ , and  $T = \frac{2}{3} \log\left(\frac{b\varepsilon\beta D}{L \log(n/\delta) \log m}\right)$ , then Algorithm 3 is  $(\varepsilon, \delta)$ -DP, queries  $n$  gradients, and has*

$$\begin{aligned} \mathbb{E}[F(x_K) - F(x^*)] & \leq O\left(D(L + \beta D) \frac{\sqrt{\log d \log n}}{\sqrt{n}}\right) \\ & \quad + O\left(\frac{\sqrt{\beta} LD^2 \log(1/\delta) \log m \log^2 n}{n\varepsilon}\right)^{2/3}. \end{aligned}$$

The following lemma proves privacy in this setting.

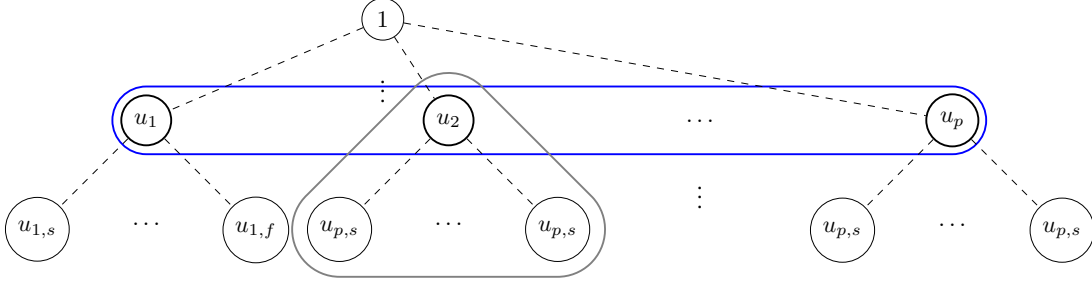
**Lemma 4.3.** *Let  $2^T \leq b$ ,  $\delta \leq 1/n$  and  $\varepsilon \leq \sqrt{2^{-T} \log(1/\delta)}$ . Setting  $\lambda_{t,s} = \frac{LD2^{T/2} \log(n/\delta)}{b\varepsilon}$ , Algorithm 3 is  $(O(\varepsilon), 21\delta)$ -DP. Moreover,  $\mathbb{E}[\langle v_{t,s}, w_{t,s} \rangle] \leq \mathbb{E}[\min_{w \in \mathcal{X}} \langle v_{t,s}, w \rangle] + O(LD2^{T/2} \log(n/\delta) \log(m) / b\varepsilon)$ .*

*Proof (sketch).* Let  $\mathcal{S} = (z_1, \dots, z_{n-1}, z_n)$  and  $\mathcal{S}' = (z_1, \dots, z_{n-1}, z'_n)$  denote two neighboring datasets with iterates  $x = (x_1, \dots, x_K)$  and  $x' = (x'_1, \dots, x'_K)$ . Here, we prove privacy after conditioning on the event that the  $n$ 'th sample is sampled at phase  $t$  and depth  $j$ . We only need to prove privacy for the iterates at phase  $t$  as the iterates before phase  $t$  do not depend on the  $n$ 'th sample and the iterates after phase  $t$  are  $(\varepsilon, \delta)$ -DP by post-processing.

Let us focus on the iterates at phase  $t$ . Let  $u_1, \dots, u_p$  denote the vertices at level  $j$  that has samples  $S_1, \dots, S_p$  each of size  $|S_i| = 2^{-j}b$ . Let  $\mathcal{A}_i$  denote the algorithm that outputs the iterates corresponding to the descendant of the vertex  $u_i$ . The proof has two steps: first, we use advanced composition to show that each  $\mathcal{A}_i$  is  $(\varepsilon_0, \delta_0)$ -DP where roughly  $\varepsilon_0 = 2^{j/2}\varepsilon$ . Then, as we have  $p = 2^j$  vertices at depth  $j$  with random samples (that is, shuffled between vertices), we use the amplification by shuffling result (Feldman et al., 2020b) (see Lemma E.2) to argue that the final algorithm is  $(\varepsilon, \delta)$ -DP (see Fig. 2 for a demonstration of the shuffling).  $\square$

Theorem 7 now follows using similar arguments to the proof of Theorem 6 (see Appendix E.4).





**Figure 2.** We view Algorithm 3 as a sequence of algorithms  $\mathcal{A}_1, \dots, \mathcal{A}_p$ , each  $\mathcal{A}_i$  operating on the subtree of the vertex  $u_i$  using the outputs of the previous algorithms. The gray set denotes the subtree over which  $\mathcal{A}_2$  operates; its outputs are the iterates corresponding to the leafs of this subtree. If each  $\mathcal{A}_i$  is  $(\varepsilon_0, \delta_0)$ -DP, then shuffling the samples at nodes of depth  $j$  (in blue) amplifies the privacy to roughly  $(\varepsilon_0 \sqrt{\log(1/\delta)}/2^j, \delta + n\delta_0)$ -DP.

## 5. Lower Bounds

We conclude the paper with tight lower bounds. Our lower bounds are for the excess empirical loss but these can be translated to lower bounds for excess population loss using a simple bootstrapping approach (Bassily et al., 2019).

### 5.1. Lower Bounds for Non-Smooth Functions

In this section, we prove tight lower bounds for non-smooth functions using bounds for estimating the sign of the mean. In this problem, given a dataset  $\mathcal{S} = (z_1, \dots, z_n)$  with mean  $\bar{z}$ , we aim to design private algorithms that estimate  $\text{sign}(\bar{z})$ . The following lemma provides a lower bound for this problem. We defer the proof to Appendix F.1.

**Lemma 5.1.** *Let  $\mathcal{S} = (z_1, \dots, z_n)$  where  $z_i \in \mathcal{Z} = \{-D/d, D/d\}^d$  and let  $\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$ . Then any  $(\varepsilon, \delta)$ -DP algorithm  $\mathcal{A} : \mathcal{Z} \rightarrow \{-1, +1\}^d$  has*

$$\max_{\mathcal{S} \in \mathcal{Z}^n} \mathbb{E} \left[ \sum_{j=1}^d |\bar{z}_j| \mathbb{1}\{\mathcal{A}(\mathcal{S})_j \neq \text{sign}(\bar{z}_j)\} \right] \geq \Omega \left( \frac{D\sqrt{d}}{n\varepsilon \log d} \right).$$

The previous lemma implies our desired lower bound.

**Theorem 8.** *Let  $f(x; z_i) = L \|x - z_i\|_1$  where  $z_i \in \mathcal{Z} = \{-D/d, D/d\}^d$ ,  $\hat{F}(x; \mathcal{S}) = \frac{1}{n} \sum_{i=1}^n \|x - z_i\|_1$ , and  $\mathcal{X} = \{x : \|x\|_1 \leq D\}$ . Then any  $(\varepsilon, \delta)$ -DP algorithm  $\mathcal{A}$  has*

$$\max_{\mathcal{S} \in \mathcal{Z}^n} \mathbb{E} \left[ \hat{F}(\mathcal{A}(\mathcal{S}); \mathcal{S}) - \min_{x \in \mathcal{X}} \hat{F}(x; \mathcal{S}) \right] \geq \Omega \left( \frac{LD\sqrt{d}}{n\varepsilon \log d} \right).$$

*Proof.* First, note that  $f(x; z_i)$  is  $L$ -Lipschitz with respect to  $\|\cdot\|_1$ . Moreover, it is immediate to see that the minimizer of  $\hat{F}(\cdot; \mathcal{S})$  is  $x^* = \text{sign}(\bar{z})D/d$  where  $\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$  is the mean. Letting  $\hat{x} = \mathcal{A}(\mathcal{S})$ , simple algebra yields

$$\hat{F}(\hat{x}; \mathcal{S}) - \hat{F}(x^*; \mathcal{S}) \geq L \sum_{j=1}^d |\bar{z}_j| \mathbb{1}\{\text{sign}(\hat{x}_j) \neq \text{sign}(\bar{z}_j)\}.$$

The claim now follows from Lemma 5.1 as  $\text{sign}(\mathcal{A}(\mathcal{S}))$  is differentially private by post-processing.  $\square$

### 5.2. Lower Bounds for Smooth Functions

In this section we prove tight lower bounds for smooth function. Specifically, we focus on  $\beta$ -smooth functions with  $\beta \approx L/D$ ; such an assumption holds for many applications including LASSO (linear regression). Our results in this section build on the lower bounds of Talwar et al. (2015) which show tight bounds for private Lasso for sufficiently large dimension. We have the following lower bound for smooth functions which we prove in Appendix F.2.

**Theorem 9.** *Let  $\mathcal{X} = \{x \in \mathbb{R}^d : \|x\|_1 \leq D\}$ . There is family of convex functions  $f : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}$  that is  $L$ -Lipschitz and  $\beta$ -smooth with  $\beta \leq L/D$  such that any  $(\varepsilon, \delta)$ -DP algorithm  $\mathcal{A}$  with  $\delta = n^{-\omega(1)}$  has*

$$\begin{aligned} & \sup_{\mathcal{S} \in \mathcal{Z}^n} \mathbb{E} \left[ \hat{F}(\mathcal{A}(\mathcal{S}); \mathcal{S}) - \min_{x \in \mathcal{X}} \hat{F}(x; \mathcal{S}) \right] \\ & \geq LD \cdot \tilde{\Omega} \left( \min \left( \frac{1}{(n\varepsilon)^{2/3}}, \frac{\sqrt{d}}{n\varepsilon} \right) \right). \end{aligned}$$

The lower bound of Theorem 9 implies the optimality of our upper bounds; if  $d \geq \tilde{O}((n\varepsilon)^{2/3})$  then the lower bound is essentially  $1/(n\varepsilon)^{2/3}$  which is achieved by the private Frank-Wolfe algorithm of Section 4, otherwise  $d \leq \tilde{O}((n\varepsilon)^{2/3})$  and the lower bound is  $\sqrt{d}/n\varepsilon$  which is the same bound that private mirror descent (Section 3) obtains.

### Acknowledgements

TK has been supported in part by the Israeli Science Foundation (ISF) grant 2549/19, by the Len Blavatnik and the Blavatnik Family foundation, and by the Yandex Initiative in Machine Learning.

### References

Abadi, M., Chu, A., Goodfellow, I., McMahan, B., Mironov, I., Talwar, K., and Zhang, L. Deep learning with differential privacy. In *23rd ACM Conference on Computer*

- and Communications Security (ACM CCS), pp. 308–318, 2016.
- Bassily, R., Smith, A., and Thakurta, A. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *55th Annual Symposium on Foundations of Computer Science*, pp. 464–473, 2014.
- Bassily, R., Feldman, V., Talwar, K., and Thakurta, A. Private stochastic convex optimization with optimal rates. In *Advances in Neural Information Processing Systems*, volume 32, pp. 11282–11291, 2019.
- Bassily, R., Guzman, C., and Nandi, A. Non-euclidean differentially private stochastic convex optimization. *arXiv:2103.01278 [cs.LG]*, 2021.
- Chaudhuri, K., Monteleoni, C., and Sarwate, A. D. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12:1069–1109, 2011.
- Duchi, J. C. Introductory lectures on stochastic convex optimization. In *The Mathematics of Data*, IAS/Park City Mathematics Series. American Mathematical Society, 2018.
- Duchi, J. C. Information theory and statistics. Lecture Notes for Statistics 311/EE 377, Stanford University, 2019. URL <http://web.stanford.edu/class/stats311/lecture-notes.pdf>. Accessed May 2019.
- Duchi, J. C., Shalev-Shwartz, S., Singer, Y., and Tewari, A. Composite objective mirror descent. In *Proceedings of the Twenty Third Annual Conference on Computational Learning Theory*, 2010.
- Dwork, C. and Roth, A. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3 & 4):211–407, 2014.
- Dwork, C., Kenthapadi, K., McSherry, F., Mironov, I., and Naor, M. Our data, ourselves: Privacy via distributed noise generation. In *Advances in Cryptology (EUROCRYPT 2006)*, 2006a.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. In *Proceedings of the Third Theory of Cryptography Conference*, pp. 265–284, 2006b.
- Dwork, C., Naor, M., Pitassi, T., and Rothblum, G. N. Differential privacy under continual observation. In *Proceedings of the Forty-Second Annual ACM Symposium on the Theory of Computing*, pp. 715–724, 2010.
- Dwork, C., Naor, M., Reingold, O., and Rothblum, G. N. Pure differential privacy for rectangle queries via private partitions. In *International Conference on the Theory and Application of Cryptology and Information Security*, pp. 735–751, 2015.
- Fang, C., Li, C. J., Lin, Z., and Zhang, T. SPIDER: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In *Advances in Neural Information Processing Systems*, volume 31, pp. 689–699, 2018.
- Feldman, V. Generalization of ERM in stochastic convex optimization: The dimension strikes back. In *Advances in Neural Information Processing Systems*, volume 29, pp. 3576–3584, 2016.
- Feldman, V., Koren, T., and Talwar, K. Private stochastic convex optimization: optimal rates in linear time. In *Proceedings of the 52nd Annual ACM on the Theory of Computing*, pp. 439–449, 2020a.
- Feldman, V., McMillan, A., and Talwar, K. Hiding among the clones: A simple and nearly optimal analysis of privacy amplification by shuffling. *arXiv:2012.12803 [cs.LG]*, 2020b.
- Hardt, M., Recht, B., and Singer, Y. Train faster, generalize better: Stability of stochastic gradient descent. In *ICML*, pp. 1225–1234, 2016. URL <http://jmlr.org/proceedings/papers/v48/hardt16.html>.
- Jain, P. and Thakurta, A. (Near) dimension independent risk bounds for differentially private learning. In *Proceedings of the 31st International Conference on Machine Learning*, pp. 476–484, 2014.
- Kifer, D., Smith, A., and Thakurta, A. Private convex empirical risk minimization and high-dimensional regression. In *Proceedings of the Twenty Fifth Annual Conference on Computational Learning Theory*, pp. 25–1, 2012.
- Lacoste-Julien, S., Schmidt, M., and Bach, F. A simpler approach to obtaining an  $o(1/t)$  convergence rate for the projected stochastic subgradient method. *arXiv preprint arXiv:1212.2002*, 2012.
- Lu, H., Freund, R. M., and Nesterov, Y. Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.
- Shalev-Shwartz, S. and Ben-David, S. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.
- Shalev-Shwartz, S., Shamir, O., Srebro, N., and Sridharan, K. Stochastic convex optimization. In *Proceedings of the Twenty Second Annual Conference on Computational Learning Theory*, 2009.

Steinke, T. and Ullman, J. Between pure and approximate differential privacy. *Journal of Privacy and Confidentiality*, 7(2):3–22, 2017.

Talwar, K., Thakurta, A., and Zhang, L. Nearly optimal private Lasso. In *Advances in Neural Information Processing Systems*, volume 28, pp. 3025–3033, 2015.

Yurtsever, A., Sra, S., and Cevher, V. Conditional gradient methods via stochastic path-integrated differential estimator. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97, pp. 7282–7291, 2019.