

# Towards Practical Differentially Private Convex Optimization

Roger Iyengar  
Carnegie Mellon University

Joseph P. Near  
University of California, Berkeley

Dawn Song  
University of California, Berkeley

Om Thakkar  
Boston University

Abhradeep Thakurta  
University of California, Santa Cruz

Lun Wang  
Peking University

**Abstract**—Building useful predictive models often involves learning from sensitive data. Training models with differential privacy can guarantee the privacy of such sensitive data. For convex optimization tasks, several differentially private algorithms are known, but none has yet been deployed in practice.

In this work, we make two major contributions towards practical differentially private convex optimization. First, we present **Approximate Minima Perturbation**, a novel algorithm that can leverage any off-the-shelf optimizer. We show that it can be employed without any hyperparameter tuning, thus making it an attractive technique for practical deployment. Second, we perform an extensive empirical evaluation of the state-of-the-art algorithms for differentially private convex optimization, on a range of publicly available benchmark datasets, and real-world datasets obtained through an industrial collaboration. We release open-source implementations of all the differentially private convex optimization algorithms considered, and benchmarks on as many as nine public datasets, four of which are high-dimensional.

## I. INTRODUCTION

Building useful predictive models often involves learning from sensitive data. To preserve the privacy of such sensitive data, many systems first carry out the learning task over the data, and then release just the final “learned” model. However, as many recent works [1], [2], [3] indicate, a model can leak information about the sensitive data it was trained on, even though the data might have never been made public. To prevent such information leakage, **differential privacy (DP)** [4], [5] has been recently used as a gold standard for performing learning tasks over sensitive data. It has also been adopted by large-scale corporations like Google [6], Apple [7], etc. Intuitively, DP prevents an adversary from confidently making any conclusions about whether a sample was used in training a model, even while having access to the model and any external side information.

The authors are ordered alphabetically.

With the recent advancements in machine learning and big data, private convex optimization has proven to be useful for large-scale learning over sensitive user data that has been collected by organizations. **The objective of this work is to provide insight into practical differentially private convex optimization, with a specific focus on the classical technique of objective perturbation** [8], [9]. Our main technical contribution is to design a new algorithm for private convex optimization that is amenable to real-world scenarios, and provide privacy and utility guarantees for it. In addition, we conduct a broad empirical evaluation of approaches for private convex optimization, including our new approach. Our evaluation is more **extensive** than prior works [8], [10], [11], [12]; it includes nine public datasets, four of which are *high-dimensional*. Apart from these, we also consider four real-world use cases, obtained in collaboration with Uber Technologies, Inc. We provide advice and resources for practitioners, including open-source implementations [13] of the algorithms evaluated, and benchmarks on all the public datasets considered.

### A. *Objective Perturbation and its Practical Feasibility*

**We focus our attention on the technique of objective perturbation**, because prior works [8], [9] as well as our own preliminary empirical results have hinted at its superior performance. The standard technique of objective perturbation [8] consists of a two-stage process: “perturbing” the objective function by adding a random linear term, and releasing the minima of the perturbed objective. It has been shown [8], [9] that releasing such a minima is sufficient for achieving DP guarantees.

However, objective perturbation provides privacy guarantees *only if* the output of the mechanism is the *exact minima* of the perturbed objective. Practical algorithms for convex optimization often involve the use of first-order iterative methods, such as gradient descent or

stochastic gradient descent (SGD) due to their scalability. However, such methods typically offer convergence rates that depend on the number of iterations carried out by the method, so they are *not* guaranteed to reach the exact minima in finite time. As a result, it is not clear if objective perturbation in its current form can be applied in a practical setting, where one is usually constrained by resources such as computing power, and reaching the minima might not be feasible.

We address the fundamental question of whether an alternative to standard objective perturbation exists which provides privacy and utility guarantees when the released model is *not necessarily* the minima of the perturbed objective. In this work, we answer this question in the positive: *it is possible to get privacy and utility guarantees even if the system releases a noisy “approximate” minima of the perturbed objective.* A major implication of this result is that one can use first-order iterative methods in combination with such an approach. This can be highly beneficial in terms of the time taken to obtain a private solution, as first-order methods often succeed to find a “good” solution very quickly.

### B. Our Approach: Approximate Minima Perturbation

We propose Approximate Minima Perturbation (AMP), a strengthened alternative to objective perturbation that provides privacy and utility guarantees when the released model is a noisy “approximate” minima for the perturbed objective. We measure the convergence of a model in terms of the Euclidean norm of the gradient of the perturbed objective at the model. The scale of the noise added to the approximate minima contains a parameter representing the maximum tolerable gradient norm. This results in a trade-off between the gradient norm bound (consequently, the amount of noise to be added to the approximate minima), and the difficulty, in practice, of being able to obtain an approximate minima within the norm bound. This can be useful in settings where a limited computing power is available, which can in turn act as a guide for setting an appropriate norm bound. We note that if the norm bound is set to zero, then this approach reduces to the setting of standard objective perturbation. Approximate Minima Perturbation also brings with it certain distinct advantages, which we will describe next.

**Approximate Minima Perturbation works for all convex objective functions:** While previous works [8], [9] provide privacy guarantees for objective perturbation *only* when the objective is a loss function of a

generalized linear model<sup>1</sup>, the guarantees provided for AMP hold for *any* objective function.<sup>2</sup> In both cases, the objective functions are assumed to possess standard properties like *Lipschitz continuity*, and *smoothness*.

**Approximate Minima Perturbation is the first feasible approach that can leverage any off-the-shelf optimizer:** AMP can accommodate any off-the-shelf optimizer as a black-box for carrying out its optimization step. This enables a simple implementation that inherits the scalability properties of the optimizer used, which can be particularly important in situations where high-performance optimizers are available. AMP is the *only* known feasible algorithm for DP convex optimization that allows the use of any off-the-shelf optimizer.

**Approximate Minima Perturbation has a competitive hyperparameter-free variant:** To ensure privacy for an algorithm, its hyperparameters must be chosen either independently of the data, or by using a differentially private hyperparameter tuning algorithm. Previous work [8], [14], [15] has shown this to be a challenging task. AMP has only four hyperparameters: one related to the Lipschitz continuity of the objective, and the other three related to splitting the privacy budget within the algorithm. In Section V-A, we present a data-independent method for setting all of them. Our empirical evaluation demonstrates that the resulting hyperparameter-free variant of AMP yields comparable accuracy to the standard variant with its hyperparameters tuned in a data-dependent manner (which may be non-private).

### C. Empirical Evaluation, & Resources for Practitioners

This work also reports on an extensive and broad empirical study of the state-of-the-art differentially private convex optimization techniques. In addition to AMP and its hyperparameter-free variant, we evaluate four existing algorithms: private gradient descent with minibatching [16], [15], both the variants (convex, and *strongly convex*) of the private Permutation-based SGD algorithm [12], and the private Frank-Wolfe algorithm [17].

Our evaluation is the largest to date, including a total of 13 datasets. We include datasets with a variety of different properties, including four high-dimensional datasets and four real-world use cases represented by datasets obtained in collaboration with Uber.

The results of our empirical evaluation demonstrate three key findings. First, we confirm the expectation that the cost of privacy decreases as dataset size increases.

<sup>1</sup>The loss function of a generalized linear model is parameterized by an inner product of the feature vector of the data, and the model.

<sup>2</sup>Our analysis, in particular, extends to objective perturbation as well.

For *all* the real-world use cases in our evaluation, we obtain differentially private models that achieve an accuracy within 4% of the non-private baseline even for very conservative settings of the privacy parameters. For reasonable values of the privacy parameters, the accuracy of the best private model is within 2% of the baseline for two of these datasets, essentially identical to the baseline for one of them, and even slightly higher than the baseline for one of the datasets! This provides empirical evidence to further the claims of previous works [18], [19] that DP can also act as a type of regularization, reducing the generalization error.

Second, our results demonstrate a general ordering of algorithms in terms of empirical accuracy. Our results show that AMP generally outperforms all the other algorithms across all the considered datasets. Under specific conditions like high-dimensionality of the dataset and sparsity of the optimal predictive model for it, we see that private Frank-Wolfe provides the best performance.

Third, our results show that a hyperparameter-free variant of AMP achieves nearly the same accuracy as the standard variant with its hyperparameters tuned in a data-dependent manner. Approximate Minima Perturbation is therefore simple to deploy in practice as it can leverage any off-the-shelf optimizer, and it has a competitive variant that does not require any hyperparameter tuning.

We provide an open-source implementation [13] of the algorithms evaluated, including our Approximate Minima Perturbation, and a complete set of benchmarks used in producing our empirical results. In addition to enabling the reproduction of our results, this set of benchmarks will provide a standard point of reference for evaluating private algorithms proposed in the future. Our open-source release represents the first benchmark for differentially private convex optimization.

#### D. Main Contributions

The main contributions of this work are as follows:

- We propose Approximate Minima Perturbation, a strengthened alternative to objective perturbation that provides privacy guarantees even for an approximate minima of the perturbed objective, and therefore allows the use of *any* off-the-shelf optimizer. No previous approach provides this capability. Compared to previous approaches, AMP also provides improved utility in practice, and works with any convex loss function under standard assumptions.
- We conduct the largest empirical study to date of state-of-the-art DP convex optimization approaches, including as many as nine public datasets, four of which are high-dimensional. Our results demonstrate that AMP generally provides the best accuracy.
- We evaluate DP convex optimization on four real-world use cases, obtained in collaboration with Uber. Our results suggest that for the large-scale datasets used in practice, privacy-preserving models can obtain essentially the same accuracy as non-private models for reasonable values of the privacy parameters. In one case, we show that a DP model achieves a *higher* accuracy than the non-private baseline.
- We present a competitive hyperparameter-free variant of AMP, allowing the approach to be deployed without the need for tuning on publicly available datasets, or by a DP hyperparameter tuning algorithm.
- We release open-source implementations [13] of all the algorithms we evaluate, and the first benchmarks for differentially private convex optimization algorithms on as many as nine public datasets.

## II. RELATED WORK

Convex optimization in the non-private setting has a long history, with several excellent resources providing a good overview [20], [21]. A lot of recent advances have been made in the field of convex Empirical Risk Minimization (ERM) as well. A comprehensive list of works on *stochastic* convex ERM has been provided in [22], whereas [23] provides dimension-dependent lower bounds for the sample complexity required for stochastic convex ERM and uniform convergence.

A large body of existing work examines the problem of differentially private convex ERM. The techniques of output perturbation and objective perturbation were first proposed in [8]. Near dimension-independent risk bounds for both the techniques were provided in [11]; however, the bounds are achieved for the standard settings of the techniques, which provide privacy guarantees only for the minima of their respective objective functions. A private SGD algorithm was first given in [10], and optimal risk bounds were provided for a later version of private SGD in [16]. A variant of output perturbation was proposed in [12] that requires the use of permutation-based SGD, and reduces sensitivity using properties of that algorithm. Several works [9], [24] deal with DP convex ERM in the setting of high-dimensional sparse regression, but the algorithms in these works also require obtaining the minima. The Frank-Wolfe algorithm [25] has also seen a resurgence lately [26], [27], [28], [29]. We study the performance of a DP version of Frank-Wolfe [17] in our empirical analysis.

There are also works in DP convex optimization apart from the ERM model. Many recent works [30], [31], [32] examine the setting of online learning, whereas high-dimensional kernel learning is considered in [33]; these settings are quite different from ours, and the results are incomparable. There have also been works [34], [35] on DP regression analysis, a subset of DP convex optimization. However, the privacy guarantees in these hold only if the algorithms are able to find some minima. There have also been advances in DP non-convex optimization, including deep learning [36], [15]. A broad survey of works in DP machine learning has been provided in [37].

Previous empirical evaluations have provided limited insight into the practical performance of the various algorithms for DP convex optimization. Output perturbation and objective perturbation are evaluated on two datasets in [8] and [11], and private SGD is evaluated in [10]. Wu et al. [12] perform the broadest comparison, including their own approach, and two variants of private SGD [10], [16] on six datasets, but they do not include objective perturbation. No prior evaluation considers the state-of-the-art algorithms from all three major lines of work in the area (output perturbation, objective perturbation, and private SGD). Moreover, none of the prior evaluations considers high-dimensional data—a maximum of 75 dimensions is considered in [12].

Our empirical evaluation is the most complete to date. We evaluate state-of-the-art algorithms from all 3 lines of work on 9 public datasets and 4 real-world use cases. We consider low-dimensional and high-dimensional (as many as 47,236 dimensions) datasets. In addition, we release open-source implementations for all algorithms, and benchmarking scripts to reproduce our results [13].

### III. PRELIMINARIES

In this section, we formally define the notation, important definitions, and results used in this work.

Given an  $n$ -element dataset  $D = \{d_1, d_2, \dots, d_n\}$ , s.t.  $d_i \in \mathcal{D}$  for  $i \in [n]$ , the objective is to get a model  $\theta$  from the following unconstrained optimization problem:

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \mathcal{L}(\theta; D),$$

where  $\mathcal{L}(\theta; D) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; d_i)$  is the empirical risk,  $p > 0$ , and  $\ell(\theta; d_i)$  is defined as a loss function for  $d_i$  that is convex in the first parameter  $\theta \in \mathbb{R}^p$ . This formulation falls under the framework of ERM, which is useful in various settings, including the widely applicable problem of classification in machine learning via linear regression, logistic regression, or support vector machines. The

notation  $\|x\|$  is used to represent the  $L_2$ -norm of a vector  $x$ . Next, we define certain basic properties of functions that will be helpful in further sections.

**Definition III.1.** A function  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  :

- is a convex function if for all  $\theta_1, \theta_2 \in \mathbb{R}^p$ ,  $f(\theta_1) - f(\theta_2) \geq \langle \nabla f(\theta_2), \theta_1 - \theta_2 \rangle$ .
- is a  $\xi$ -strongly convex function if for all  $\theta_1, \theta_2 \in \mathbb{R}^p$ ,  $f(\theta_1) \geq f(\theta_2) + \langle \nabla f(\theta_2), \theta_1 - \theta_2 \rangle + \frac{\xi}{2} \|\theta_1 - \theta_2\|^2$ , or equivalently,  $\langle \nabla f(\theta_1) - \nabla f(\theta_2), \theta_1 - \theta_2 \rangle \geq \xi \|\theta_1 - \theta_2\|^2$ .
- has  $L_q$ -Lipschitz constant  $L$  if for all  $\theta_1, \theta_2 \in \mathbb{R}^p$ ,  $|f(\theta_1) - f(\theta_2)| \leq L \cdot \|\theta_1 - \theta_2\|_q$ .
- is  $\beta$ -smooth if for all  $\theta_1, \theta_2 \in \mathbb{R}^p$ ,  $\|\nabla f(\theta_1) - \nabla f(\theta_2)\| \leq \beta \cdot \|\theta_1 - \theta_2\|$ .

To establish the notion of DP, we first define neighboring datasets. We will refer to a pair of datasets  $D, D' \in \mathcal{D}^n$  as neighbors, if  $D'$  can be obtained from  $D$  by modifying one sample  $d_i \in D$  for some  $i \in [n]$ .

**Definition III.2** ( $(\epsilon, \delta)$ -Differential Privacy [4], [5]). A (randomized) algorithm  $M$  with input domain  $\mathcal{D}^n$  and output range  $\mathcal{R}$  is  $(\epsilon, \delta)$ -differentially private if for all pairs of neighboring inputs  $D, D' \in \mathcal{D}^n$ , and every measurable  $S \subseteq \mathcal{R}$ , we have with probability over the coin flips of  $M$  that:

$$\Pr(M(D) \in S) \leq e^\epsilon \cdot \Pr(M(D') \in S) + \delta.$$

One of the most common techniques for achieving differential privacy is the Gaussian mechanism, for which we first need the  $L_2$ -sensitivity of a function.

**Definition III.3** ( $L_2$ -sensitivity). A function  $f : \mathcal{D}^n \rightarrow \mathbb{R}^p$  has  $L_2$ -sensitivity  $\Delta$  if

$$\max_{\substack{D, D' \in \mathcal{D}^n \text{ s.t.} \\ (D, D') \text{ neighbors}}} \|f(D) - f(D')\| = \Delta.$$

**Lemma III.1** (Gaussian mechanism [38]). If a function  $f : \mathcal{D}^n \rightarrow \mathbb{R}^p$  has  $L_2$ -sensitivity  $\Delta$ , then the mechanism  $M$ , which on input  $D \in \mathcal{D}^n$  outputs  $f(D) + b$ , where  $b \sim \mathcal{N}(0, \sigma^2 I_{p \times p})$  and  $\sigma = \frac{\Delta}{\epsilon} \left(1 + \sqrt{2 \log \frac{1}{\delta}}\right)$ , satisfies  $(\epsilon, \delta)$ -differential privacy. Here,  $\mathcal{N}(0, \sigma^2 I_{p \times p})$  denotes the  $p$ -dimensional zero-mean Gaussian distribution with each dimension having variance  $\sigma^2$ .

Lastly, we define Generalized Linear Models (GLMs).

**Definition III.4** (Generalized Linear Model). For a model space  $\theta \in \mathbb{R}^p$ , where  $p > 0$ , the sample space  $\mathcal{D}$  in a GLM is defined as the cartesian product of a  $p$ -dimensional feature space  $\mathcal{X} \subseteq \mathbb{R}^p$  and a label space  $\mathcal{Y}$ , i.e.,  $\mathcal{D} = \mathcal{X} \times \mathcal{Y}$ . Thus, each data sample  $d_i \in \mathcal{D}$

can be decomposed into a feature vector  $x_i \in \mathcal{X}$ , and a label  $y_i \in \mathcal{Y}$ . Moreover, the loss function  $\ell(\theta; d_i)$  for a GLM is a function of  $x_i^T \theta$  and  $y_i$ .

#### IV. APPROXIMATE MINIMA PERTURBATION

In this section, we will describe Approximate Minima Perturbation, a strengthened alternative to objective perturbation that provides DP guarantees in the case even when the output of the algorithm is not the actual minima of the perturbed objective function. The perturbed objective takes the form  $\mathcal{L}(\theta; D) + \frac{\Lambda}{2} \|\theta\|^2 + \langle b, \theta \rangle$ , where  $b$  is a random variable drawn from an appropriate distribution, and  $\Lambda$  is an appropriately chosen regularization constant. We make two crucial improvements over the original objective perturbation algorithm [8], [9]:

- The privacy guarantee of objective perturbation holds only at the *exact* minima of the underlying optimization problem, which is *never guaranteed in practice given finite time*. We show that AMP provides a privacy guarantee even for an *approximate* solution.
- Earlier privacy analyses for objective perturbation [8], [9] hold only when the loss function  $\ell(\theta; d)$  is a loss for a GLM (see Definition III.4), as they implicitly make a rank-one assumption on the Hessian of the loss  $\nabla^2 \ell(\theta; d)$ . Via a careful perturbation analysis of the Hessian, we extend the analysis to any convex loss function under standard assumptions. It is important to note that AMP reduces to objective perturbation if the “approximate” minima condition is tightened to getting the actual minima of the perturbed objective.

**Algorithmic description:** Given a dataset  $D = \{d_1, d_2, \dots, d_n\}$ , where each  $d_i \in \mathcal{D}$ , we consider (objective) functions of the form  $\mathcal{L}(\theta; D) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; d_i)$ , where  $\theta \in \mathbb{R}^p$  is a model, loss  $\ell(\theta; d_i)$  has  $L_2$ -Lipschitz constant  $L$  for all  $d_i$ , is convex in  $\theta$ , has a continuous Hessian, and is  $\beta$ -smooth in both the parameters.

At a high level, Approximate Minima Perturbation provides a convergence-based solution for objective perturbation. In other words, once the algorithm finds a model  $\theta_{approx}$  for which the norm of the gradient of the perturbed objective  $\nabla \mathcal{L}_{priv}(\theta_{approx}; D)$  is within a pre-determined threshold  $\gamma$ , it outputs a noisy version of  $\theta_{approx}$ , denoted by  $\theta_{out}$ . Since the perturbed objective is strongly convex, it is sufficient to add Gaussian noise, with standard deviation  $\sigma_2$  having a linear dependence on the norm bound  $\gamma$ , to  $\theta_{approx}$  to ensure DP.

Details of AMP are provided in Algorithm 1. Note that although we get a relaxed constraint on the regularization parameter  $\Lambda$  (in Algorithm 1) if the loss function  $\ell$  is a

loss for a GLM, the privacy guarantees hold for general convex loss functions as well. The parameters  $(\epsilon_1, \delta_1)$  within the algorithm represent the amount of the privacy budget dedicated to perturbing the objective, with the rest of the budget  $(\epsilon_2, \delta_2)$  being used for adding noise to the approximate minima  $\theta_{approx}$ . On the other hand, the parameter  $\epsilon_3$  intuitively represents the part of the privacy budget  $\epsilon_1$  allocated to scaling the noise added to the objective function. The remaining budget  $(\epsilon_1 - \epsilon_3)$  is used to set the amount of regularization used.

---

#### Algorithm 1: Approximate Minima Perturbation

---

**Input:** Dataset:  $D = \{d_1, \dots, d_n\}$ ; loss function:  $\ell(\theta; d_i)$  that has  $L_2$ -Lipschitz constant  $L$ , is convex in  $\theta$ , has a continuous Hessian, and is  $\beta$ -smooth for all  $\theta \in \mathbb{R}^p$  and all  $d_i$ ; Hessian rank bound parameter:  $r$  which is the minimum of  $p$  and twice the upper bound on the rank of  $\ell$ 's Hessian; privacy parameters:  $(\epsilon, \delta)$ ; gradient norm bound:  $\gamma$ .

- 1 Set  $\epsilon_1, \epsilon_2, \epsilon_3, \delta_1, \delta_2 > 0$  such that  $\epsilon = \epsilon_1 + \epsilon_2$ ,  $\delta = \delta_1 + \delta_2$ , and  $0 < \epsilon_1 - \epsilon_3 < 1$
  - 2 Set  $\Lambda \geq \frac{r\beta}{\epsilon_1 - \epsilon_3}$
  - 3  $b_1 \sim \mathcal{N}(0, \sigma_1^2 I_{p \times p})$ , where  $\sigma_1 = \frac{(\frac{2L}{n})(1 + \sqrt{2 \log \frac{1}{\delta_1}})}{\epsilon_3}$
  - 4 Let  $\mathcal{L}_{priv}(\theta; D) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; d_i) + \frac{\Lambda}{2n} \|\theta\|^2 + b_1^T \theta$
  - 5  $\theta_{approx} \leftarrow \theta$  such that  $\|\nabla \mathcal{L}_{priv}(\theta; D)\| \leq \gamma$
  - 6  $b_2 \sim \mathcal{N}(0, \sigma_2^2 I_{p \times p})$ , where  $\sigma_2 = \frac{(\frac{n\gamma}{\Lambda})(1 + \sqrt{2 \log \frac{1}{\delta_2}})}{\epsilon_2}$
  - 7 Output  $\theta_{out} = \theta_{approx} + b_2$
- 

**Privacy and utility guarantees:** Here, we provide the privacy and utility guarantees for Algorithm 1. While we provide a complete privacy analysis (Theorem 1), we only state the utility guarantee (Theorem 2) as it is a slight modification from previous work [9]. We provide a proof for it in Appendix A for completeness.

**Theorem 1** (Privacy guarantee). *Algorithm 1 is  $(\epsilon, \delta)$ -differentially private.*

*Proof Idea.* For obtaining an  $(\epsilon, \delta)$ -DP guarantee for Algorithm 1, we first split the output of the algorithm into two parts: one being the exact minima of the perturbed objective, whereas the other containing the exact minima, the approximate minima obtained in Step 5 of the algorithm, as well as the Gaussian noise added to it. For the first part, we bound the ratio of the density of the exact minima taking any particular value, under any two neighboring datasets, by  $e^{\epsilon_1}$  with probability at

least  $1 - \delta_1$ . We first simplify such a ratio, as done in [8] via the function inverse theorem, by transforming it to two ratios: one involving only the density of a function of the minima value and the input dataset, and the other involving the determinant of this function's Jacobian. For the former ratio, we start by bounding the sensitivity of the function using the  $L_2$ -Lipschitz constant  $L$  of the loss function. Then, we use the guarantees of the Gaussian mechanism to obtain a high-probability bound (shown in Lemma IV.1). We bound the latter ratio (in Lemma IV.2) via a novel approach that uses the  $\beta$ -smoothness property of the loss. Next, we use the gradient norm bound  $\gamma$ , and the strong convexity of the perturbed objective to obtain an  $(\epsilon_2, \delta_2)$ -DP guarantee for the second part of the split output. Lastly, we use the general composition property of DP to get the statement of the theorem.  $\square$

*Proof.* Define  $\theta_{\min} = \arg \min_{\theta \in \mathbb{R}^p} \mathcal{L}_{\text{priv}}(\theta; D)$ . Fix a pair of neighboring datasets  $D^*, D' \in \mathcal{D}^n$ , and some  $\alpha \in \mathbb{R}^p$ . First, we will show that:

$$\frac{\text{pdf}_{D^*}(\theta_{\min} = \alpha)}{\text{pdf}_{D'}(\theta_{\min} = \alpha)} \leq e^{\epsilon_1} \text{ w.p. } \geq 1 - \delta_1. \quad (1)$$

We define  $b(\theta; D) = -\nabla \mathcal{L}(\theta; D) - \frac{\Lambda \theta}{n}$  for  $D \in \mathcal{D}^n$  and  $\theta \in \mathbb{R}^p$ . Changing variables according to the function inverse theorem (Theorem 17.2 in [39]), we get

$$\frac{\text{pdf}_{D^*}(\theta_{\min} = \alpha)}{\text{pdf}_{D'}(\theta_{\min} = \alpha)} = \frac{\text{pdf}(b(\alpha; D^*); \epsilon_1, \delta_1, L)}{\text{pdf}(b(\alpha; D'); \epsilon_1, \delta_1, L)} \cdot \frac{|\det(\nabla b(\alpha; D'))|}{|\det(\nabla b(\alpha; D^*))|} \quad (2)$$

We will bound the ratios of the densities and the determinants separately. First, we will show that for  $\epsilon_3 < \epsilon_1$ ,  $\frac{\text{pdf}(b(\alpha; D^*); \epsilon_1, \delta_1, L)}{\text{pdf}(b(\alpha; D'); \epsilon_1, \delta_1, L)} \leq e^{\epsilon_3}$  w.p. at least  $1 - \delta_1$ , and then we will show that  $\frac{|\det(\nabla b(\alpha; D'))|}{|\det(\nabla b(\alpha; D^*))|} \leq e^{\epsilon_1 - \epsilon_3}$  if  $\epsilon_1 - \epsilon_3 < 1$ .

**Lemma IV.1.** *We define  $b(\theta; D) = -\nabla \mathcal{L}(\theta; D) - \frac{\Lambda \theta}{n}$  for  $D \in \mathcal{D}^n$ , and  $\theta \in \mathbb{R}^p$ . Then, for any pair of neighboring datasets  $D^*, D' \in \mathcal{D}^n$ , and  $\epsilon_3 < \epsilon_1$ , we have*

$$\frac{\text{pdf}(b(\alpha; D^*); \epsilon_1, \delta_1, L)}{\text{pdf}(b(\alpha; D'); \epsilon_1, \delta_1, L)} \leq e^{\epsilon_3} \text{ w.p. at least } 1 - \delta_1.$$

Here,  $\mathcal{L}_{\text{priv}}(\alpha; \cdot)$  is defined as in Algorithm 1.

*Proof.* Assume w.l.o.g. that  $d_i \in D^*$  has been replaced by  $d'_i$  in  $D'$ . We first bound the  $L_2$ -sensitivity of  $b(\alpha; \cdot)$ :

$$\|b(\alpha; D^*) - b(\alpha; D')\| \leq \frac{\|\nabla \ell(\alpha; d'_i) - \nabla \ell(\alpha; d_i)\|}{n} \leq \frac{2L}{n},$$

where the last inequality follows as  $\|\nabla \ell(\alpha; \cdot)\| \leq L$ .

Setting  $\sigma_1 \geq \frac{(\frac{2L}{n})(1 + \sqrt{2 \log \frac{1}{\delta_1}})}{\epsilon_3}$  for  $\epsilon_3 < \epsilon_1$ , we get the statement of the lemma from the guarantees of the Gaussian mechanism [38].  $\square$

**Lemma IV.2.** *Let  $b(\theta; D)$  be defined as in Lemma IV.1. Then for any pair of neighboring datasets  $D^*, D' \in \mathcal{D}^n$ , if  $\epsilon_1 - \epsilon_3 < 1$ , we have*

$$\frac{|\det(\nabla b(\alpha; D'))|}{|\det(\nabla b(\alpha; D^*))|} \leq e^{\epsilon_1 - \epsilon_3}.$$

*Proof.* Assume w.l.o.g. that  $d_i \in D^*$  is replaced by  $d'_i$  in  $D'$ . Let  $A = n \nabla^2 \mathcal{L}(\alpha; D^*)$ , and  $E = \nabla^2 \ell(\alpha; d'_i) - \nabla^2 \ell(\alpha; d_i)$ . As the  $(p \times p)$  matrix  $E$  is the difference of the Hessians of the loss of two individual samples, we can define a bound  $r$  on the rank of  $E$  as follows:

$$r = \min \{p, 2 \cdot (\text{upper bound on rank of } \nabla^2 \ell(\alpha; \cdot))\}$$

Let  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_p$  be the eigenvalues of  $A$ , and  $\lambda'_1 \leq \lambda'_2 \leq \dots \leq \lambda'_p$  be the eigenvalues of  $A + E$ . Thus,

$$\begin{aligned} \frac{|\det(\nabla b(\alpha; D'))|}{|\det(\nabla b(\alpha; D^*))|} &= \frac{\det\left(\frac{A+E+\Lambda \mathbb{I}_p}{n}\right)}{\det\left(\frac{A+\Lambda \mathbb{I}_p}{n}\right)} = \frac{\prod_{i=1}^p (\lambda'_i + \Lambda)}{\prod_{i=1}^p (\lambda_i + \Lambda)} \\ &= \prod_{i=1}^p \left(1 + \frac{\lambda'_i - \lambda_i}{\lambda_i + \Lambda}\right) \\ &\leq \prod_{i=1}^p \left(1 + \frac{|\lambda'_i - \lambda_i|}{\Lambda}\right) \\ &= 1 + \sum_{i=1}^p \frac{|\lambda'_i - \lambda_i|}{\Lambda} \\ &\quad + \sum_{\substack{i,j \in [p], \\ i \neq j}} \frac{\prod_{k \in \{i,j\}} |\lambda'_k - \lambda_k|}{\Lambda^2} + \dots \\ &\leq 1 + \frac{r\beta}{\Lambda} + \frac{(r\beta)^2}{\Lambda^2} + \dots \leq \frac{\Lambda}{\Lambda - r\beta} \end{aligned}$$

The first inequality follows since  $A$  is a positive semi-definite matrix (as  $\ell$  is convex) and thus,  $\lambda_j \geq 0$  for all  $j \in [p]$ . The second inequality follows as i) the rank of  $E$  is at most  $r$ , ii) both  $A$  and  $A + E$  are positive semi-definite (so  $\lambda_j, \lambda'_j \geq 0$  for all  $j \in [p]$ ), and iii) we have an upper bound  $\beta$  on the eigenvalues of  $A$  and  $A + E$  due to  $\ell(\theta; d_j)$  being convex in  $\theta$ , having a continuous Hessian, and being  $\beta$ -smooth. The last inequality follows if  $\Lambda > r\beta$ . Also, we want  $\frac{\Lambda}{\Lambda - r\beta} \leq \exp(\epsilon_1 - \epsilon_3)$ , which implies  $\Lambda \geq \frac{r\beta}{1 - \exp(\epsilon_3 - \epsilon_1)} \geq \frac{r\beta}{\epsilon_1 - \epsilon_3}$ . Both conditions are satisfied by setting  $\Lambda = \frac{r\beta}{\epsilon_1 - \epsilon_3}$  as  $\epsilon_1 - \epsilon_3 < 1$ .  $\square$

From Equation 2, and Lemmas IV.1 and IV.2, we get that  $\frac{\text{pdf}_{D^*}(\theta_{\min} = \alpha)}{\text{pdf}_{D'}(\theta_{\min} = \alpha)} \leq e^{\epsilon_1}$  w.p.  $\geq 1 - \delta_1$ . In other words,  $\theta_{\min}$  is  $(\epsilon_1, \delta_1)$ -differentially private.

Now, since we can write  $\theta_{out}$  as  $\theta_{out} = \theta_{min} + (\theta_{approx} - \theta_{min} + b_2)$ , we will prove that releasing  $(\theta_{approx} - \theta_{min} + b_2)$  is  $(\epsilon_2, \delta_2)$ -differentially private.

**Lemma IV.3.** *For  $D \in \mathcal{D}^n$ , let  $\gamma \geq 0$  be chosen independent of  $D$ , and let  $\theta_{min} = \arg \min_{\theta \in \mathbb{R}^p} \mathcal{L}_{priv}(\theta; D)$ . If  $\theta_{approx} \in \mathbb{R}^p$  s.t.  $\|\nabla \mathcal{L}_{priv}(\theta_{approx}; D)\| \leq \gamma$ , then releasing  $(\theta_{approx} - \theta_{min} + b_2)$ , where  $b_2 \sim \mathcal{N}(0, \sigma_2^2 I_{p \times p})$  for  $\sigma_2 = \frac{(\frac{n\gamma}{\Lambda})(1 + \sqrt{2 \log \frac{1}{\delta_2}})}{\epsilon_2}$ , is  $(\epsilon_2, \delta_2)$ -DP.*

*Proof.* We start by bounding the  $L_2$ -norm of  $(\theta_{approx} - \theta_{min})$ :

$$\begin{aligned} \|\theta_{approx} - \theta_{min}\| &\leq \frac{n \|\nabla \mathcal{L}_{priv}(\theta_{approx}; D) - \nabla \mathcal{L}_{priv}(\theta_{min}; D)\|}{\Lambda} \\ &\leq \frac{n\gamma}{\Lambda} \end{aligned} \quad (3)$$

The first inequality follows as  $\mathcal{L}_{priv}$  is  $\frac{\Lambda}{n}$ -strongly convex, and the second inequality follows as  $\nabla \mathcal{L}_{priv}(\theta_{min}; D) = 0$  and  $\nabla \mathcal{L}_{priv}(\theta_{approx}; D) \leq \gamma$ .

Now, setting  $\sigma_2 = \frac{(\frac{n\gamma}{\Lambda})(1 + \sqrt{2 \log \frac{1}{\delta_2}})}{\epsilon_2}$ , we get the statement of the lemma by the properties of the Gaussian mechanism [38].  $\square$

As  $\epsilon_1 + \epsilon_2 = \epsilon$ , and  $\delta_1 + \delta_2 = \delta$ , we get the privacy guarantee of Algorithm 1 by Equation 1, Lemma IV.3, and the general composition property of DP [40].  $\square$

Next, we provide the utility guarantee (in terms of excess empirical risk) for Algorithm 1 in Theorem 2. For completeness, we provide a proof for it in Appendix A.

**Theorem 2** (Utility guarantee (adapted from [9])). *Let  $\hat{\theta}$  be the true unconstrained minimizer of  $\mathcal{L}(\theta; D) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; d_i)$ , and  $r = \min\{p, 2 \cdot (\text{upper bound on rank of } \ell\text{'s Hessian})\}$ . In Algorithm 1, if  $\epsilon_i = \frac{\epsilon}{2}$  for  $i \in \{1, 2\}$ ,  $\epsilon_3 = \max\{\frac{\epsilon_1}{2}, \epsilon_1 - 0.99\}$ ,  $\delta_j = \frac{\delta}{2}$  for  $j \in \{1, 2\}$ , and we set the regularization parameter  $\Lambda = \Theta\left(\frac{1}{\|\hat{\theta}\|} \left(\frac{L\sqrt{rp \log 1/\delta}}{\epsilon} + \sqrt{\frac{n^2 L \gamma \sqrt{p \log 1/\delta}}{\epsilon}}\right)\right)$  such that it satisfies the constraint in Step 2, then the following is true:*

$$\begin{aligned} &\mathbb{E}(\mathcal{L}(\theta_{out}; D) - \mathcal{L}(\hat{\theta}; D)) \\ &= O\left(\frac{\|\hat{\theta}\| L \sqrt{rp \log \frac{1}{\delta}}}{n\epsilon} + \|\hat{\theta}\| \sqrt{\frac{L \gamma \sqrt{p \log \frac{1}{\delta}}}{\epsilon}}\right). \end{aligned}$$

**Remark:** For loss functions of Generalized Linear Models, we have  $r = 2$ . Here, for small values of  $\gamma$  (for example,  $\gamma = O(\frac{1}{n^2})$ ), the excess empirical risk of Approximate Minima Perturbation is asymptotically the same as that of objective perturbation [9], and has a better dependence on  $n$  than that of Private Permutation-based SGD [12]. Specifically, the dependence is  $\propto \frac{1}{n}$  for Approximate Minima Perturbation, and  $\propto \frac{1}{\sqrt{n}}$  for Private PSGD.

### Towards Hyperparameter-free Approximate Minima Perturbation:

AMP can be considered to have the following hyperparameters: the Lipschitz constant  $L$ , and the privacy parameters  $\epsilon_2, \delta_2$ , and  $\epsilon_3$  which split the privacy budget within the algorithm. A data-independent approach for setting these parameters can eliminate the need for hyperparameter tuning with this approach, making it convenient to deploy in practice.

For practical applications, given a sensitive dataset and a convex loss function, the  $L$  hyperparameter can be thought of as a trade-off between the sensitivity of the loss and the amount of external interference required to achieve that sensitivity, for instance, *sample clipping* (defined in Section V-A) on the data. In the next section, we provide a hyperparameter-free variant of AMP that has performance comparable to the standard variant in which all the hyperparameters are tuned.

## V. EXPERIMENTAL RESULTS

Our evaluation seeks to answer two major research questions:

- 1) **What is the cost (to accuracy) of privacy?** How close can a DP model come to the non-private baseline? For real-world use cases, is the cost of privacy low enough to make DP learning practical?
- 2) **Which algorithm provides the best accuracy in practice?** Is there a total order on the available algorithms? Does this ordering differ for datasets with different properties?

Additionally, we also attempt to answer the following question which can result in a significant advantage for the deployment of a DP model in practice:

- 3) **Can Approximate Minima Perturbation be deployed without hyperparameter tuning?** Can its hyperparameters ( $L, \epsilon_2, \delta_2$ , and  $\epsilon_3$ ) be set in a data-independent manner?

**Summary of results:** *Question #1:* Our results demonstrate that for datasets of sufficient size, *the cost of privacy is negligible*. Experiments 1 (on low-dimensional datasets), 2 (on high-dimensional datasets), and 3 (on

real-world datasets obtained in collaboration with Uber) evaluate the cost of privacy using logistic loss. Our results show that for large datasets, a DP model exists that approaches the accuracy of the non-private baseline at reasonable privacy budgets. Experiment 3 shows that for the larger datasets common in practical settings, a privacy-preserving model can produce even *better* accuracy than the non-private one, which suggests that privacy-preserving learning is indeed practical.

We also present the performance of private algorithms using Huber SVM loss (on all the datasets mentioned above) in Appendix B. The general trends from the experiments using Huber SVM loss are identical to those obtained using logistic loss.

*Question #2:* Our experiments demonstrate that AMP generally provides the best accuracy among all the evaluated algorithms. Moreover, experiment 2 shows that under specific conditions, private Frank-Wolfe can provide the best accuracy. In all the regimes, the results generally show an improvement over other approaches.

*Question #3:* Our experiments also demonstrate that a simple data-independent method can be used to set  $L, \epsilon_2, \delta_2$ , and  $\epsilon_3$  for AMP, and that this method provides good accuracy across datasets. For most values of  $\epsilon$ , our data-independent approach provides nearly the same accuracy as the version tuned using a grid search (which may be non-private).

#### A. Experiment Setup

**Algorithms evaluated:** Our evaluation includes one algorithm drawn from each of the major approaches to private convex optimization: objective perturbation, output perturbation, private gradient descent, and the private Frank-Wolfe algorithm. For each approach, we select the best-known algorithm and configuration.

For objective perturbation, we implement AMP (Algorithm 1) as it is the only practically feasible objective perturbation approach<sup>3</sup>. For all the experiments, we tune the value of the hyperparameters  $L, \epsilon_2, \delta_2$ , and  $\epsilon_3$  using the grid search described later.

We also evaluate a hyperparameter-free variant of AMP that sets the hyperparameters  $L, \epsilon_2, \delta_2$  and  $\epsilon_3$  independent of the data. We describe the strategy in detail towards the end of this subsection.

For output perturbation, we implement Private Perturbation-based SGD (PSGD) [12], as it is the only practically feasible variant of output perturbation<sup>3</sup>. We evaluate both the variants, with minibatching, proposed

<sup>3</sup>For all variants pertaining to the standard regime in [8], obtaining some exact minima is necessary for achieving a privacy guarantee.

TABLE I  
DATASETS USED IN OUR EVALUATION

Dataset	# Samples	# Dim.	# Classes
<b>Low-Dimensional Datasets (Public)</b>			
Synthetic-L	10,000	20	2
Adult	45,220	104	2
KDDCup99	70,000	114	2
Covertime	581,012	54	7
MNIST	65,000	784	10
<b>High-Dimensional Datasets (Public)</b>			
Synthetic-H	5,000	5,000	2
Gisette	6,000	5,000	2
Real-sim	72,309	20,958	2
RCV1	50,000	47,236	2
<b>Real-World Datasets (Uber)</b>			
Dataset #1	4m	23	2
Dataset #2	18m	294	2
Dataset #3	18m	20	2
Dataset #4	19m	70	2

in [12]: convex (Algorithm 3), and strongly convex (Algorithm 4). For the convex variant, we evaluate all three proposed learning rate schemes (constant learning rate, decreasing learning rate, and square-root learning rate). We include results only for constant learning rate, as our experiments show that this scheme produces the most accurate models.

For private gradient descent, we implement a variant of the private SGD algorithm originally proposed in [16]. Our variant (Algorithm 2) leverages the Moments Accountant [15], incorporates minibatching, and sets the noise parameter based on the desired number of iterations (as compared to a fixed  $n^2$  iterations in [16]).

For private Frank-Wolfe, we implement the version (Algorithm 5) originally proposed in [17]. This algorithm performs constrained optimization by design. Following the advice of [26], we use a decreasing learning rate for better accuracy guarantees. Unlike the other algorithms, private Frank-Wolfe has nearly dimension-independent error bounds, so it should be expected to perform comparatively better on high-dimensional datasets.

Pseudocodes for all the evaluated algorithms are provided in Appendix C for reference. For each algorithm, we evaluate the variant that provides  $(\epsilon, \delta)$ -differential privacy. Most algorithms can also provide  $\epsilon$ -differential privacy at an additional cost to accuracy.

**Datasets:** Table I lists the public datasets used in our experimental evaluation. Each of these datasets is available for download, and our open-source release contains scripts for downloading and pre-processing these datasets. It also contains scripts for generating both the synthetic datasets. As RCV1 has multi-label classification over 103 labels (with most of the labels being used



for a very small proportion of the dataset), for this dataset we consider the task of predicting whether a sample is categorized under the most frequently used label or not.

The selected datasets include both *low-dimensional* and *high-dimensional* datasets. We define low-dimensional datasets to be ones where  $n \gg p$  (where  $n$  is the number of samples and  $p$  is the number of dimensions). High-dimensional datasets are defined as those for which  $n$  and  $p$  are on roughly the same scale, i.e.  $n \leq p$  (or nearly so). We consider the Synthetic-H, Gisette, Real-sim, and RCV1 datasets to be high-dimensional.

To obtain training and testing sets, we randomly shuffle the dataset, take the first 80% as the training set, and the remaining 20% as the testing set.

**Sample clipping:** Each of the algorithms we evaluate has the requirement that the loss have a Lipschitz constant. We can enforce this requirement for the loss functions we consider by bounding the norm for each sample. We can accomplish this by pre-processing the dataset, but it must be done carefully to preserve DP.

For all the algorithms except private Frank-Wolfe, to make the loss have an  $L_2$ -Lipschitz constant  $L$ , we bound the influence of each sample  $(x_i, y_i)$  by clipping the feature vector  $x_i$  to  $\left(x_i \cdot \min\left(1, \frac{L}{\|x_i\|}\right)\right)$ . This transformation is independent of other samples, and thus preserves DP; it has also been previously used, e.g. in [15]. As the private Frank-Wolfe algorithm requires the loss to have a relaxed  $L_1$ -Lipschitz constant  $L$ , it suffices (using Theorem 1 from [41]) to bound the  $L_\infty$ -norm of each sample  $(x_i, y_i)$  by  $L$ . We achieve this by clipping each dimension  $x_{i,j}$ , where  $j \in [d]$ , to  $\min(x_{i,j}, L)$ .

**Hyperparameters:** Each of the evaluated algorithms has at least one hyperparameter. The values for these hyperparameters should be tuned to provide the best accuracy, but the tuning should be done privately in order to guarantee end-to-end differential privacy. Although a number of differentially private hyperparameter tuning algorithms have been proposed [8], [14], [15] to address this problem, they add more variance in the performance of each algorithm, thus making it more difficult to compare the performance across different algorithms.

In order to provide a fair comparison between algorithms, we use a grid search to determine the *best* value for each hyperparameter. Our grid search considers the hyperparameter values listed in Table II. In addition to the standard algorithm hyperparameters  $(\Lambda, \eta, T, k)$ , we tune the clipping parameter  $L$  used in pre-processing the datasets, and the constraint on the model space

TABLE II  
HYPERPARAMETER & PRIVACY PARAMETER VALUES

Hyperparameter	Values Considered
$\Lambda$ (regularization factor)	$10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 0$
$\eta$ (learning rate)	0.001, 0.01, 0.1, 1
$T$ (number of iterations)	5, 10, 100, 1000, 5000
$k$ (minibatch size)	50, 100, 300
$L$ (clipping threshold)	0.1, 1, 10, 100
$C$ (model constraint)	1, 10
$f$ (output budget fraction)	0.001, 0.01, 0.1, 0.5
$f_1$ (privacy budget fraction)	0.9, 0.92, 0.95, 0.98, 0.99
Privacy Parameter	Values Considered
$\epsilon$	$10^{-2}, 10^{-\frac{3}{2}}, 10^{-1}, 10^{-\frac{1}{2}}, 10^0, 10^{\frac{1}{2}}, 10^1$
$\delta$	$\frac{1}{n^2}$

used by private Frank-Wolfe, Private SGD when using regularized loss, and Private strongly convex PSGD. The parameter  $C$  controls the size of the  $L_1/L_2$ -ball from which models are selected by private Frank-Wolfe/the other algorithms respectively. For AMP, we set  $\epsilon_2 = f \cdot \epsilon$ ,  $\delta_2 = f \cdot \delta$ , and tune for  $f$ . Here,  $f$  denotes the fraction of the budget  $(\epsilon, \delta)$  that is allocated to  $(\epsilon_2, \delta_2)$ . Also, since the valid range of the hyperparameter  $\epsilon_3$  depends on the value of  $\epsilon_1$ , we set  $\epsilon_3 = f_1 \cdot \epsilon_1$ , and tune for  $f_1$ . We also ensure that the constraint on  $\epsilon_3$  in Line 3 of Algorithm 1 is satisfied. Note that tuning hyperparameters may be non-private, but it enables a direct comparison of the algorithms themselves.

We consider a range of values for the privacy parameter  $\epsilon$ . Following Wu et al. [12], we set the privacy parameter  $\delta = \frac{1}{n^2}$ , where  $n$  is the size of the training data. The complete set of values considered is listed in Table II. For multiclass classification datasets such as MNIST and Covertype, we implement the one-vs-all strategy by training a binary classifier for each class, and split  $\epsilon$  and  $\delta$  equally among the binary classifiers so that we can achieve an overall  $(\epsilon, \delta)$ -DP guarantee by using general composition [40].

**Algorithm Implementations:** The implementations used in our evaluation correspond to the pseudocode listings in Appendix C, are written in Python, and are available in our open source release [13]. For Approximate Minima Perturbation, we define the loss and gradient according to Algorithm 1, and leverage SciPy’s `minimize` procedure to find the approximate minima.

For all datasets, our implementation is able to achieve  $\gamma = \frac{1}{n^2}$ , where  $n$  is the size of the training data. For low-dimensional datasets, our implementation of AMP uses SciPy’s `BFGS` solver, for which we can specify the desired norm bound  $\gamma$ . The `BFGS` algorithm stores

the full Hessian of the objective, which does not fit in memory for the sparse high-dimensional datasets in our study. For these, we define an alternative low-memory implementation using SciPy’s `L-BFGS-B` solver, which does not store the full Hessian.

**Experiment procedure:** Our experiment setup is designed to find the best possible accuracy achievable for a given setting of the privacy parameters. To ensure a fair comparison, we begin every run of each algorithm with the initial model  $0^p$ . Because each of the evaluated algorithms introduces randomness due to noise, we train 10 independent models for each combination of the hyperparameter setting. We report the mean accuracy and standard deviation for the combination of the hyperparameter setting with the highest mean accuracy over the 10 independent runs.<sup>4</sup>

**Differences with the setting in [12]:** Although both the studies have 3 datasets in common (Coverttype, KDDCup99, and MNIST), our setting is slightly different from [12] for all 3 of them. For Coverttype, our study uses all 7 classes, while [12] uses a binary version. For KDDCup99, we use a 10% sample of the full dataset (as in [42]), while [12] uses the full dataset. For MNIST, we use all 784 dimensions, while [12] uses random projection to reduce the dimensionality to 50.

The results we obtain for both the variants of the Private PSGD algorithm [12] are based on faithful implementations of those algorithms. We tune the hyperparameters for both, using the grid search described earlier.

**Non-private baseline:** Note that one of the main objectives of this study is to determine the cost of privacy in practice for convex optimization. Hence, to provide a point of comparison for our results, we also train a non-private baseline model for each experiment. We use Scikit-learn’s `LogisticRegression` class to train this model on the same training data as the private algorithms, and test its accuracy on the same testing data as the private algorithms. We do not perform sample clipping when training this model.

**Strategy for Hyperparameter-free Approximate Minima Perturbation:** Now, we describe a data-independent approach for setting Approximate Minima Perturbation’s only hyperparameters,  $L$ ,  $\epsilon_2$ ,  $\delta_2$ , and  $\epsilon_3$ , for *both* the loss functions we consider (see Section V-B). For  $L$ , we find that setting  $L = 1$  achieves a good trade-off between

the amount of noise added for perturbing the objective, and the information loss after sample clipping across all datasets. Next, we consider *only* the synthetically generated datasets for setting the hyperparameters specific to AMP. Fixing  $\gamma = \frac{1}{n^2}$ , we find that setting  $\epsilon_2 = 0.01 \cdot \epsilon$  and  $\delta_2 = 0.01 \cdot \delta$  achieves a good trade-off between the budget for perturbing the objective, and the amount of noise that its approximate minima can tolerate. For setting  $\epsilon_3$ , we consider two separate cases:

- For  $\epsilon_1 = 0.99 \cdot \epsilon$ , and  $\epsilon_3 = f_1 \cdot \epsilon_1$ , we see that setting  $f_1 = 0.99$  for  $\epsilon_1 = 0.0099$ ,  $f_1 = 0.95$  for  $\epsilon_1 \in \{0.0313, 0.099\}$ , and  $f_1 = 0.9$  for  $\epsilon_1 \in \{0.313, 0.99, 3.13, 9.99\}$  yields a good accuracy for Synthetic-L. Hence, we observe that for very low values of  $\epsilon_1$ , a good accuracy is yielded by  $\epsilon_3$  close to  $\epsilon_1$  (i.e., most of the budget is used to reduce the scale of the noise, and the influence of regularization is kept large). As  $\epsilon_1$  increases, we see that it is more beneficial to reduce the effects of regularization. We fit a basic polynomial curve of the form  $y = a + bx^{-c}$ , where  $a, b, c > 0$ , to the above-stated values to get a dependence of  $f_1$  (the privacy budget fraction) in terms of  $\epsilon_1$ . We combine it with the lower bound imposed on  $f_1$  by Theorem 1 (for instance, we require  $f_1 \geq 0.9$  for  $\epsilon_1 = 9.99$ ) to obtain the following data-independent relationship between  $\epsilon_1$  and  $\epsilon_3$  for low-dimensional datasets:

$$\epsilon_3 = \max \left\{ \min \left\{ 0.887 + \frac{0.019}{\epsilon_1^{0.373}}, 0.99 \right\}, 1 - \frac{0.99}{\epsilon_1} \right\} \cdot \epsilon_1$$

- For Synthetic-H, we see that setting  $f_1 = 0.97$  yields a good accuracy for all the values of  $\epsilon_1$  considered. Thus, combining it with the lower bound imposed on  $f_1$  by Theorem 1, we obtain the following relationship for high-dimensional datasets:

$$\epsilon_3 = \max \left\{ 0.97, 1 - \frac{0.99}{\epsilon_1} \right\} \cdot \epsilon_1$$

Note that the results for this strategy are consistent for both loss functions across all the public and the real-world datasets considered, *none* of which were used in defining the strategy except for setting the Lipschitz constant  $L$  of the loss. They can be considered to be effectively serving as test-cases for the strategy.

## B. Loss Functions

Our evaluation considers the loss functions for two commonly used models: logistic regression and Huber SVM. This section contains results for logistic regression; results for Huber SVM are available in Appendix B.

**Logistic regression:** The  $L_2$ -regularized logistic regression loss function on a sample  $(x, y)$  with  $y \in \{1, -1\}$  is  $\ell(\theta, (x, y)) = \ln(1 + \exp(-y\langle \theta, x \rangle)) + \frac{\lambda}{2} \|\theta\|^2$ .

<sup>4</sup>The results shown are for hyperparameters tuned via the mean test set accuracy. Since all the considered algorithms aim to minimize the empirical loss, we also conducted experiments by tuning via the mean training set accuracy. Both settings provided visibly identical results.

Our experiments consider both the regularized and un-regularized (i.e.,  $\Lambda = 0$ ) settings. The un-regularized version has  $L_2$ -Lipschitz constant  $L$  when for each sample  $x$ ,  $\|x\| \leq L$ . It is also  $L^2$ -smooth. The regularized version has  $L_2$ -Lipschitz constant  $L + \Lambda C$  when for each sample  $x$ ,  $\|x\| \leq L$ , and for each model  $\theta$ ,  $\|\theta\| \leq C$ . It is also  $(L^2 + \Lambda)$ -smooth, and  $\Lambda$ -strongly convex.

### C. Experiment 1: Low-dimensional Datasets

We present the results of our experiments with logistic regression on low-dimensional data in Figure 2. All four algorithms perform better in comparison with the non-private baseline for binary classification tasks (Synthetic-L, Adult, and KDDCup99) than for multi-class problems (Covertypes and MNIST), because  $\epsilon$  and  $\delta$  must be split among the binary classifiers built for each class.

Figure 1 contains precise accuracy numbers for each dataset for reasonably low values of  $\epsilon$ . These results provide a more precise comparison between the four algorithms, and quantify the accuracy loss versus the non-private baseline for each one. Across all datasets, Approximate Minima Perturbation generally provides the most accurate models across  $\epsilon$  values.

### D. Experiment 2: High-dimensional Datasets

For this experiment, we repeat the procedure in Experiment 1 on high-dimensional data, and present the results in Figure 2. The results are somewhat different in the high-dimensional regime. We observe that although Approximate Minima Perturbation generally outperforms all the other algorithms, the private Frank-Wolfe algorithm performs the best on Synthetic-H. From prior works [11], [17], we know that both objective perturbation and the private Frank-Wolfe have near dimension-independent utility guarantees when the loss is of a GLM, and we indeed observe this expected behavior from our experiments. As in experiment 1, we present precise accuracy numbers for  $\epsilon = 0.1$  in Figure 1.

Private Frank-Wolfe works best when the optimal model is sparse (i.e., a few important features characterize the classification task well), as in the Synthetic-H dataset, which is well-characterized by just ten important features. This is because private Frank-Wolfe adds at most a single feature to the model at each iteration, and noise increases with the number of iterations. However, noise does *not* increase with the *total* number of features, since it scales with the bound on the  $\ell_\infty$ -norm of the samples. This behavior is in contrast to Approximate Minima Perturbation (and the other algorithms considered in our evaluation), for which noise scales with the

Dataset	NP baseline	AMP	H-F AMP	P-SGD	P-PSGD	P-SCPSGD	P-FW
<b>Low-Dimensional Binary Datasets (<math>\epsilon = 0.1</math>)</b>							
Synthetic-L	94.9	<b>83.1</b>	80.6	81.6	81.7	76.4	81.8
Adult	84.8	<b>79.1</b>	78.7	78.5	77.4	77.2	76.9
KDDCup99	99.1	97.5	97.4	98.0	<b>98.1</b>	95.8	96.8
<b>Low-Dimensional Multi-class Datasets (<math>\epsilon = 1^5</math>)</b>							
Covertypes	71.2	64.3	63.5	<b>65.0</b>	62.4	62.2	63.0
MNIST	91.5	<b>71.9</b>	70.5	68.6	68.0	63.2	65.0
<b>High-Dimensional Datasets (<math>\epsilon = 0.1</math>)</b>							
Synthetic-H	95.8	53.2	51.4	52.8	53.5	52.0	<b>57.6</b>
Gisette	96.6	<b>62.8</b>	59.7	61.5	62.3	61.3	58.3
Real-sim	93.3	<b>73.1</b>	71.9	66.3	66.1	65.6	69.8
RCV1	93.5	<b>64.2</b>	59.9	55.1	58.9	56.2	64.1
<b>Real-World Datasets (<math>\epsilon = 0.1</math>)</b>							
Dataset #1	75.3	<b>75.3<sup>6</sup></b>	75.3	75.3	75.3	75.3	75.3
Dataset #2	72.2	<b>70.4</b>	70.1	69.8	69.5	68.9	68.6
Dataset #3	73.6	<b>71.9</b>	71.8	71.8	71.4	71.2	71.6
Dataset #4	82.1	<b>81.7</b>	81.7	81.7	81.5	81.3	81.0

Fig. 1. Accuracy results (in %) for logistic regression. For each dataset, the result in bold represents the DP algorithm with the best accuracy for that dataset. A key for the abbreviations used for the algorithms is provided in Table III.

TABLE III  
LIST OF ABBREVIATIONS USED FOR ALGORITHMS

Abbreviation	Full-form
NP baseline	Non-private baseline
AMP	Approximate Minima Perturbation
H-F AMP	Hyperparameter-free AMP
P-SGD	Private SGD
P-PSGD	Private PSGD
P-SCPSGD	Private strongly convex PSGD
P-FW	Private Frank-Wolfe

bound on the  $\ell_2$ -norm of the samples. Private Frank-Wolfe therefore approaches the non-private baseline better than the other algorithms for high-dimensional datasets with sparse models, even at low values of  $\epsilon$ .

### E. Experiment 3: Real-world Use Cases

For this experiment, we repeat the procedure in Experiment 1 on real-world use cases, obtained in collaboration with Uber. These use cases are represented by four datasets, each of which has separately been used to train a production model deployed at Uber. The details of these datasets are listed in Table I. The results of this

<sup>5</sup>We report the accuracy for  $\epsilon = 1$  for multi-class datasets, as compared to  $\epsilon = 0.1$  for datasets with binary classification, because multi-class classification is a more difficult task than binary classification.

<sup>6</sup>For Dataset #1, AMP slightly outperforms even the NP baseline, as can be seen from Figure 2.

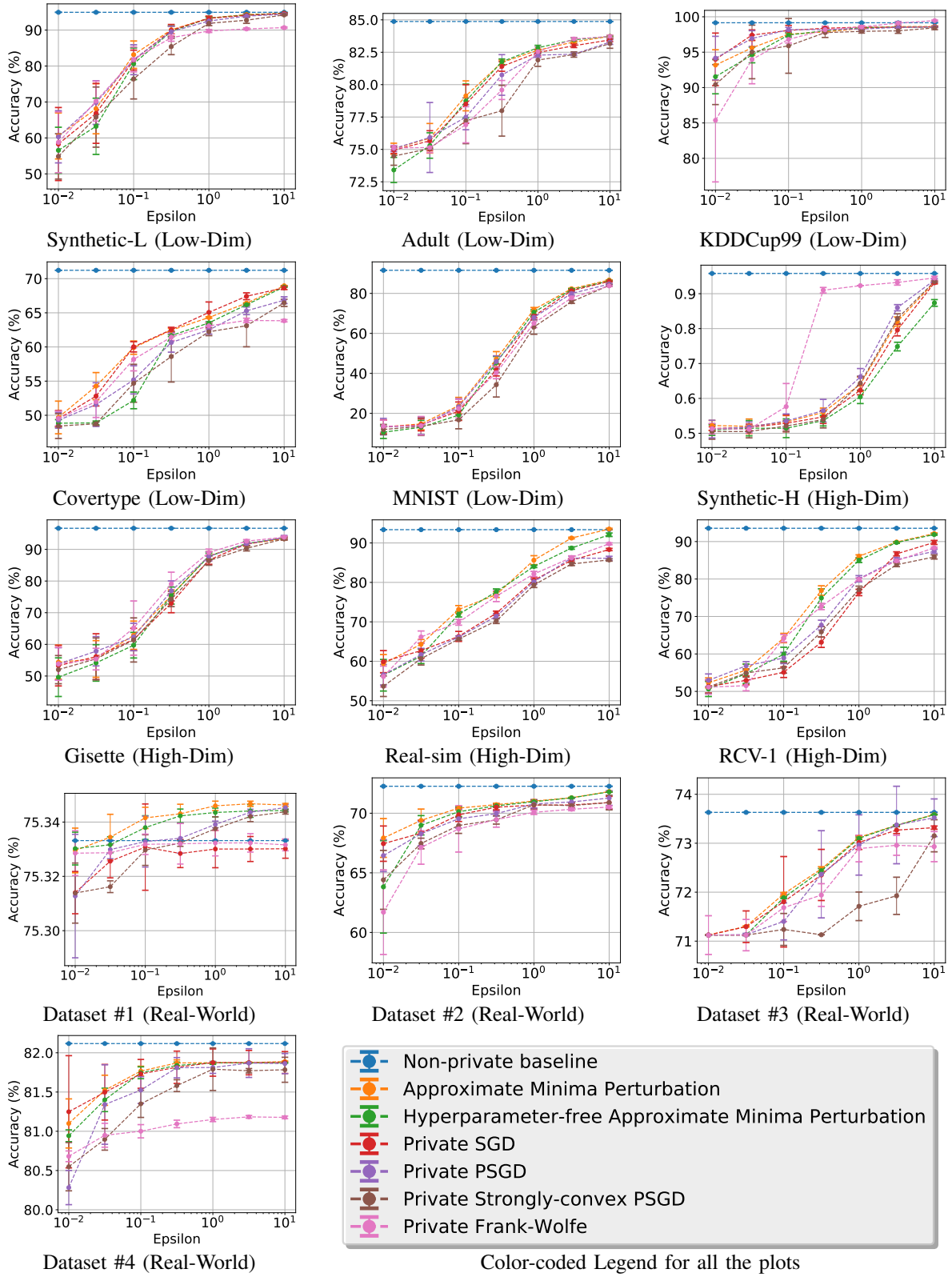


Fig. 2. Accuracy results for logistic regression on low-dimensional, high-dimensional and real-world datasets. Horizontal axis depicts varying values of  $\epsilon$ ; vertical axis shows accuracy (in %) on the testing set.

experiment are depicted in Figure 2, with more precise results for  $\epsilon = 0.1$  in Figure 1.

The real-world datasets are much larger than the datasets considered in Experiment 1. The difference in scale is reflected in the results: all of the algorithms converge to the non-private baseline for very low values of  $\epsilon$ . These results suggest that in many practical settings, the *cost of privacy is negligible*. In fact, for Dataset #1, some differentially private models exhibit a slightly *higher* accuracy than the non-private baseline for a wide range of  $\epsilon$ . For instance, even Hyperparameter-free AMP, which is end-to-end differentially private as there is no tuning involved, yields an accuracy of 75.34% for  $\epsilon = 0.1$  versus the non-private baseline of 75.33%. Some prior works [18], [19] have theorized that differential privacy could act as a type of regularization for the system, and improve the generalization error; this empirical result of ours aligns with this claim.

### F. Discussion

**For large datasets, the cost of privacy is low.** Our results confirm the expectation that very accurate differentially private models exist for large datasets. Even for relatively small datasets like Adult and KDDCup99 (where  $n < 100,000$ ), our results show that a differentially private model has accuracy within 6% of the non-private baseline even for a conservative privacy setting of  $\epsilon = 0.1$ .

For *all* the larger real-world datasets ( $n > 1\text{m}$ ), the accuracy of the best differentially private model is within 4% of the non-private baseline even for the most conservative privacy value considered ( $\epsilon = 0.01$ ). For  $\epsilon = 0.1$ , it is within 2% of the baseline for two of these datasets, essentially identical to the baseline for one of them, and even slightly higher than the baseline for one.

These results suggest that for realistic deployments on large datasets ( $n > 1\text{m}$ , and low-dimensional), a differentially private model can be deployed without much loss in accuracy.

**Approximate Minima Perturbation almost always provides the best accuracy, and is easily deployable in practice.** Our results in all the experiments demonstrate that among the available algorithms for differentially private convex optimization, our Approximate Minima Perturbation approach almost always produces models with the best accuracy. For four of the five low-dimensional datasets, and all the public high-dimensional datasets we considered, Approximate Minima Perturbation provided consistently better accuracy than the other algorithms. Under some conditions like high-dimensionality of the

datasets, and sparsity of the optimal predictive model for it, private Frank-Wolfe does give the best performance. Unlike Approximate Minima Perturbation, however, no hyperparameter-free variant of private Frank-Wolfe exists—and suboptimal hyperparameter values can reduce accuracy significantly for this algorithm.

As mentioned earlier, Approximate Minima Perturbation also has important properties that enable its practical deployment. It can leverage any off-the-shelf optimizer as a black box, allowing implementations to use existing scalable optimizers (our implementation uses Scipy’s `minimize`). None of the other evaluated algorithms have these properties.

**Hyperparameter-free Approximate Minima Perturbation provides good utility.** As demonstrated by our experimental results, AMP can be deployed without tuning hyperparameters, at little cost to accuracy. Our data-independent approach therefore enables deployment—without significant loss of accuracy—in practical settings where public data may not be available for tuning.

## VI. CONCLUSION

This paper takes two important steps towards practical differentially private convex optimization. We have presented Approximate Minima Perturbation, a novel algorithm for differentially private convex optimization that does not require the optimization process to reach the true minima. It can leverage any off-the-shelf solver, and can be employed without hyperparameter tuning. Therefore, it is amenable to be deployed in practice.

We have also performed an extensive empirical evaluation of state-of-the-art approaches for differentially private convex optimization. To encourage the further development and deployment, we have released the implementations used in our evaluation, and the benchmarking scripts used to obtain the datasets and perform the experiments. This benchmark provides a standard point of comparison for further advances in differentially private convex optimization.

## VII. ACKNOWLEDGEMENTS

The authors would like to thank Adam Smith for the discussions regarding the main privacy proof of AMP, and the anonymous reviewers for their helpful comments. This material is in part based upon work supported by NSF CCF-1740850, DARPA contract #N66001-15-C-4066, the Center for Long-Term Cybersecurity, and Berkeley Deep Drive. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the authors, and do not necessarily reflect the views of the sponsors.

## REFERENCES

- [1] M. Fredrikson, S. Jha, and T. Ristenpart, “Model inversion attacks that exploit confidence information and basic countermeasures,” in *Proceedings of the 22Nd ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS ’15. New York, NY, USA: ACM, 2015, pp. 1322–1333.
- [2] X. Wu, M. Fredrikson, S. Jha, and J. F. Naughton, “A methodology for formalizing model-inversion attacks,” in *2016 IEEE 29th Computer Security Foundations Symposium (CSF)*, June 2016, pp. 355–370.
- [3] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, “Membership inference attacks against machine learning models,” in *2017 IEEE Symposium on Security and Privacy (SP)*, May 2017, pp. 3–18.
- [4] C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor, “Our data, ourselves: Privacy via distributed noise generation,” in *EUROCRYPT*, 2006.
- [5] C. Dwork, F. McSherry, K. Nissim, and A. Smith, “Calibrating noise to sensitivity in private data analysis,” in *Theory of Cryptography Conference*. Springer, 2006, pp. 265–284.
- [6] Ú. Erlingsson, V. Pihur, and A. Korolova, “Rappor: Randomized aggregatable privacy-preserving ordinal response,” in *Proceedings of the 2014 ACM SIGSAC conference on computer and communications security*. ACM, 2014, pp. 1054–1067.
- [7] “Apple tries to peek at user habits without violating privacy,” *The Wall Street Journal*, 2016.
- [8] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate, “Differentially private empirical risk minimization,” *JMLR*, 2011.
- [9] D. Kifer, A. Smith, and A. Thakurta, “Private convex empirical risk minimization and high-dimensional regression,” *Journal of Machine Learning Research*, vol. 1, p. 41, 2012.
- [10] S. Song, K. Chaudhuri, and A. D. Sarwate, “Stochastic gradient descent with differentially private updates,” in *Global Conference on Signal and Information Processing (GlobalSIP)*, 2013 IEEE. IEEE, 2013, pp. 245–248.
- [11] P. Jain and A. Thakurta, “(near) dimension independent risk bounds for differentially private learning,” in *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32*, ser. ICML’14. JMLR.org, 2014, pp. 1–476–1–484.
- [12] X. Wu, F. Li, A. Kumar, K. Chaudhuri, S. Jha, and J. Naughton, “Bolt-on differential privacy for scalable stochastic gradient descent-based analytics,” in *Proceedings of the 2017 ACM International Conference on Management of Data*, ser. SIGMOD ’17. New York, NY, USA: ACM, 2017, pp. 1307–1322.
- [13] “Differentially Private Convex Optimization Benchmark,” 2017. [Online]. Available: <https://github.com/sunblaze-ucb/dpml-benchmark>
- [14] K. Chaudhuri and S. Vinterbo, “A stability-based validation procedure for differentially private machine learning,” in *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*, ser. NIPS’13. USA: Curran Associates Inc., 2013, pp. 2652–2660.
- [15] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, “Deep learning with differential privacy,” in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS ’16. New York, NY, USA: ACM, 2016, pp. 308–318.
- [16] R. Bassily, A. Smith, and A. Thakurta, “Private empirical risk minimization: Efficient algorithms and tight error bounds,” in *Foundations of Computer Science (FOCS)*, 2014 IEEE 55th Annual Symposium on. IEEE, 2014, pp. 464–473.
- [17] K. Talwar, A. Thakurta, and L. Zhang, “Private empirical risk minimization beyond the worst case: The effect of the constraint set geometry,” *CoRR*, vol. abs/1411.5417, 2014.
- [18] R. Bassily, A. D. Smith, and A. Thakurta, “Private empirical risk minimization, revisited,” *CoRR*, vol. abs/1405.7085, 2014.
- [19] C. Dwork, V. Feldman, M. Hardt, T. Pitassi, O. Reingold, and A. Roth, “Generalization in adaptive data analysis and holdout reuse,” in *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, ser. NIPS’15. Cambridge, MA, USA: MIT Press, 2015, pp. 2350–2358.
- [20] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [21] S. Bubeck *et al.*, “Convex optimization: Algorithms and complexity,” *Foundations and Trends® in Machine Learning*, vol. 8, no. 3–4, pp. 231–357, 2015.
- [22] L. Zhang, T. Yang, and R. Jin, “Empirical risk minimization for stochastic convex optimization:  $O(1/n)$ - and  $O(1/n^2)$ -type of risk bounds,” in *Proceedings of the 2017 Conference on Learning Theory*, ser. Proceedings of Machine Learning Research, S. Kale and O. Shamir, Eds., vol. 65. Amsterdam, Netherlands: PMLR, 07–10 Jul 2017, pp. 1954–1979.
- [23] V. Feldman, “Generalization of erm in stochastic convex optimization: The dimension strikes back,” in *Advances in Neural Information Processing Systems 29*, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, Eds. Curran Associates, Inc., 2016, pp. 3576–3584.
- [24] A. Smith and A. Thakurta, “Differentially private feature selection via stability arguments, and the robustness of the lasso,” in *COLT*, 2013.
- [25] M. Frank and P. Wolfe, “An algorithm for quadratic programming,” *Naval Research Logistics (NRL)*, vol. 3, no. 1–2, pp. 95–110, 1956.
- [26] M. Jaggi, “Revisiting frank-wolfe: projection-free sparse convex optimization,” in *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*. JMLR.org, 2013, pp. 1–427.
- [27] S. Lacoste-Julien and M. Jaggi, “An affine invariant linear convergence analysis for frank-wolfe algorithms,” *arXiv preprint arXiv:1312.7864*, 2013.
- [28] —, “On the global linear convergence of frank-wolfe optimization variants,” in *Advances in Neural Information Processing Systems*, 2015, pp. 496–504.
- [29] S. Lacoste-Julien, “Convergence Rate of Frank-Wolfe for Non-Convex Objectives,” Jun. 2016, 6 pages.
- [30] P. Jain, P. Kothari, and A. Thakurta, “Differentially private online learning,” in *COLT*, vol. 23, 2012, pp. 24–1.
- [31] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, “Local privacy and statistical minimax rates,” in *Foundations of Computer Science (FOCS)*, 2013 IEEE 54th Annual Symposium on. IEEE, 2013, pp. 429–438.
- [32] A. G. Thakurta and A. Smith, “(nearly) optimal algorithms for private online learning in full-information and bandit settings,” in *Advances in Neural Information Processing Systems 26*, C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, Eds., 2013, pp. 2733–2741.
- [33] P. Jain and A. Thakurta, “Differentially private learning with kernels,” in *ICML*, 2013.
- [34] J. Zhang, Z. Zhang, X. Xiao, Y. Yang, and M. Winslett, “Functional mechanism: Regression analysis under differential privacy,” *Proc. VLDB Endow.*, vol. 5, no. 11, pp. 1364–1375, Jul. 2012.
- [35] X. Wu, M. Fredrikson, W. Wu, S. Jha, and J. F. Naughton, “Revisiting differentially private regression: Lessons from learning theory and their consequences,” *CoRR*, vol. abs/1512.06388, 2015.
- [36] R. Shokri and V. Shmatikov, “Privacy-preserving deep learning,” in *2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, Sept 2015, pp. 909–910.
- [37] Z. Ji, Z. C. Lipton, and C. Elkan, “Differential privacy and machine learning: a survey and review,” *CoRR*, vol. abs/1412.7584, 2014.

- [38] A. Nikolov, K. Talwar, and L. Zhang, "The geometry of differential privacy: The sparse and approximate cases," in *Proceedings of the Forty-fifth Annual ACM Symposium on Theory of Computing*, ser. STOC '13. New York, NY, USA: ACM, 2013, pp. 351–360.
- [39] P. Billingsley, *Probability and Measure*, ser. Wiley Series in Probability and Statistics. Wiley, 1995.
- [40] C. Dwork, A. Roth *et al.*, "The algorithmic foundations of differential privacy," *Foundations and Trends in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2014.
- [41] R. Paulavičius and J. Žilinskas, "Analysis of different norms and corresponding lipschitz constants for global optimization," *Ukio Technologinis ir Ekonominis Vystymas*, vol. 12, no. 4, pp. 301–306, 2006.
- [42] K. Chaudhuri, C. Monteleoni, and A. D. Sarwate, "Differentially private empirical risk minimization," *Journal of Machine Learning Research*, vol. 12, no. Mar, pp. 1069–1109, 2011.
- [43] S. Dasgupta and L. Schulman, "A probabilistic analysis of em for mixtures of separated, spherical gaussians," *JMLR*, 2007.
- [44] K. Talwar, A. Thakurta, and L. Zhang, "Nearly optimal private lasso," in *NIPS*, 2015.

## APPENDIX

### A. Omitted Proofs

Here, we provide a proof for the utility guarantee of Algorithm 1, which is provided in Theorem 2. For bounding the expected risk of the algorithm, we first need to bound its empirical risk (Lemma A.1).

**Lemma A.1** (Empirical Risk). *Let  $\hat{\theta}$  be the minimizer of the objective function  $\mathcal{L}(\theta; D) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; d_i)$ , and  $\theta_{\min}$  be the minimizer of the objective function  $\mathcal{L}_{\text{priv}}(\theta; D) = \mathcal{L}(\theta; D) + \frac{\Lambda}{2n} \|\theta\|^2 + b_1^T \theta$ , where  $b_1$  is as defined in Algorithm 1. Also, let  $\theta_{\text{out}}$  be the output of Algorithm 1. We have:*

$$\begin{aligned} \mathcal{L}(\theta_{\text{out}}; D) - \mathcal{L}(\hat{\theta}; D) &\leq L \left( \frac{n\gamma}{\Lambda} + \|b_2\| \right) \\ &\quad + \frac{\Lambda \|\hat{\theta}\|^2}{2n} + \frac{2n \|b_1\|^2}{\Lambda}. \end{aligned}$$

*Proof.* We have

$$\begin{aligned} \mathcal{L}(\theta_{\text{out}}; D) - \mathcal{L}(\hat{\theta}; D) &= (\mathcal{L}(\theta_{\text{out}}; D) - \mathcal{L}(\theta_{\min}; D)) \\ &\quad + (\mathcal{L}(\theta_{\min}; D) - \mathcal{L}(\hat{\theta}; D)) \end{aligned} \quad (4)$$

First, we will bound  $(\mathcal{L}(\theta_{\text{out}}; D) - \mathcal{L}(\theta_{\min}; D))$ . We have:

$$\begin{aligned} \mathcal{L}(\theta_{\text{out}}; D) - \mathcal{L}(\theta_{\min}; D) &\leq |\mathcal{L}(\theta_{\text{out}}; D) - \mathcal{L}(\theta_{\min}; D)| \\ &\leq L \|\theta_{\text{out}} - \theta_{\min}\| \\ &= L \|\theta_{\text{approx}} - \theta_{\min} + b_2\| \\ &\leq L \|\theta_{\text{approx}} - \theta_{\min}\| \\ &\quad + L \|b_2\| \\ &\leq L \left( \frac{n\gamma}{\Lambda} + \|b_2\| \right) \end{aligned} \quad (5)$$

The second inequality above follows from the Lipschitz property of  $\mathcal{L}(\cdot; D)$ . The first equality follows as  $\theta_{\text{out}} = \theta_{\min} + (\theta_{\text{approx}} - \theta_{\min} + b_2)$ , whereas the last inequality follows from inequality 3.

Next, we bound  $(\mathcal{L}(\theta_{\min}; D) - \mathcal{L}(\hat{\theta}; D))$  on the lines of the proof of Lemma 3 in [9]. Let  $\theta^\# = \arg \min_{\theta \in \mathbb{R}^p} \mathcal{L}^\#(\theta; D)$ , where  $\mathcal{L}^\#(\theta; D) = \mathcal{L}(\theta; D) + \frac{\Lambda}{2n} \|\theta\|^2$ . As a result,  $\mathcal{L}_{\text{priv}}(\theta; D) = \mathcal{L}^\#(\theta; D) + b_1^T \theta$ . So, we have:

$$\begin{aligned} \mathcal{L}(\theta_{\min}; D) - \mathcal{L}(\hat{\theta}; D) &= \mathcal{L}^\#(\theta_{\min}; D) - \mathcal{L}^\#(\theta^\#; D) \\ &\quad + \mathcal{L}^\#(\theta^\#; D) - \mathcal{L}^\#(\hat{\theta}; D) \\ &\quad + \frac{\Lambda \|\hat{\theta}\|^2}{2n} - \frac{\Lambda \|\theta_{\min}\|^2}{2n} \\ &\leq \mathcal{L}^\#(\theta_{\min}; D) - \mathcal{L}^\#(\theta^\#; D) \\ &\quad + \frac{\Lambda \|\hat{\theta}\|^2}{2n} \end{aligned} \quad (6)$$

The inequality above follows as  $\mathcal{L}^\#(\theta^\#; D) \leq \mathcal{L}^\#(\hat{\theta}; D)$ .

Let us now bound  $\mathcal{L}^\#(\theta_{\min}; D) - \mathcal{L}^\#(\theta^\#; D)$ . To this end, we first observe that since  $\mathcal{L}_{\text{priv}}$  is  $\frac{\Lambda}{n}$ -strongly convex in  $\theta$ , we have that

$$\begin{aligned} \mathcal{L}_{\text{priv}}(\theta^\#; D) &\geq \mathcal{L}_{\text{priv}}(\theta_{\min}; D) \\ &\quad - \nabla \mathcal{L}_{\text{priv}}(\theta_{\min}; D)^T (\theta_{\min} - \theta^\#) \\ &\quad + \frac{\Lambda}{2n} \|\theta^\# - \theta_{\min}\|^2 \\ &= \mathcal{L}_{\text{priv}}(\theta_{\min}; D) + \frac{\Lambda}{2n} \|\theta^\# - \theta_{\min}\|^2 \end{aligned} \quad (7)$$

The equality above follows as  $\|\nabla \mathcal{L}_{\text{priv}}(\theta_{\min}; D)\| = 0$ .

Substituting the definition of  $\mathcal{L}_{\text{priv}}(\cdot; D)$  in equality 7, we get that

$$\begin{aligned} \mathcal{L}^\#(\theta_{\min}; D) - \mathcal{L}^\#(\theta^\#; D) &\leq b_1^T (\theta^\# - \theta_{\min}) \\ &\quad - \frac{\Lambda}{2n} \|\theta^\# - \theta_{\min}\|^2 \end{aligned} \quad (8)$$

$$\leq \|b_1\| \cdot \|\theta^\# - \theta_{\min}\| \quad (9)$$

Inequality 9 above follows by the Cauchy–Schwarz inequality.

Now, since  $\mathcal{L}^\#(\theta_{\min}; D) - \mathcal{L}^\#(\theta^\#; D) \geq 0$ , it follows from inequalities 8 and 9 that

$$\begin{aligned} \|b_1\| \cdot \|\theta^\# - \theta_{\min}\| &\geq \frac{\Lambda}{2n} \|\theta^\# - \theta_{\min}\|^2 \\ \Rightarrow \|\theta^\# - \theta_{\min}\| &\leq \frac{2n \|b_1\|}{\Lambda} \end{aligned} \quad (10)$$

We get the statement of the lemma from equation 4, and inequalities 5, 6, 9, and 10.  $\square$

Now, we are ready to prove Theorem 2.

*Proof of Theorem 2.* The proof is on the lines of the proof of Theorem 4 in [9]. First, let us get a high probability bound on  $\mathcal{L}(\theta_{out}; D) - \mathcal{L}(\hat{\theta}; D)$ . To this end, we will first bound  $\|b_1\|$  and  $\|b_2\|$  w.h.p., where  $b_s \sim \mathcal{N}(0, \sigma_s^2 I_{p \times p})$  for  $s \in \{1, 2\}$ . Using Lemma 2 from [43], we get that w.p.  $\geq 1 - \frac{\alpha}{2}$ ,

$$\|b_s\| \leq \sigma_s \sqrt{2p \log \frac{2}{\alpha}}.$$

Substituting this into Lemma A.1, we get that w.p.  $\geq 1 - \alpha$ ,

$$\begin{aligned} \mathcal{L}(\theta_{out}; D) - \mathcal{L}(\hat{\theta}; D) &\leq L \left( \frac{n\gamma}{\Lambda} + \sigma_2 \sqrt{2p \log \frac{2}{\alpha}} \right) \\ &\quad + \frac{\Lambda \|\hat{\theta}\|^2}{2n} + \frac{4n\sigma_1^2 p \log \frac{2}{\alpha}}{\Lambda}. \end{aligned}$$

It is easy to see that by making  $\epsilon_i = \frac{\epsilon}{2}$  for  $i \in \{1, 2\}$ ,  $\epsilon_3 = \max\{\frac{\epsilon_1}{2}, \epsilon_1 - 0.99\}$ ,  $\delta_j = \frac{\delta}{2}$  for  $j \in \{1, 2\}$ , and setting  $\Lambda = \Theta\left(\frac{L\sqrt{rp \log 1/\delta}}{\epsilon \|\hat{\theta}\|} + \frac{n}{\|\hat{\theta}\|} \sqrt{\frac{L\gamma \sqrt{p \log 1/\delta}}{\epsilon}}\right)$  such that it satisfies the constraint in Step 2 in Algorithm 1, we get the statement of the theorem.  $\square$

### B. Results for Huber SVM

This section reports the results of experiments with the Huber SVM loss function. The Huber SVM loss function is a differentiable and smooth approximation of the standard SVM's hinge loss. We define the loss function as in [18]. Defining  $z = y\langle x, \theta \rangle$ , the Huber SVM loss function is:

$$\ell(\theta, (x, y)) = \begin{cases} 1 - z & 1 - z > h \\ 0 & 1 - z < -h \\ \frac{(1-z)^2}{4h} + \frac{1-z}{2} + \frac{h}{4} & \text{otherwise} \end{cases}$$

As with logistic regression, the Huber SVM loss function has  $L_2$ -Lipschitz constant  $L$  when for each sample  $x$ ,  $\|x\| \leq L$ .

We repeat the experiments of Section V with the Huber SVM loss. To ensure that the experiments run to completion for Synthetic-H, we run the experiments on 2000 samples, each consisting of 2000 dimensions. For all the experiments, we obtain the non-private baseline using SciPy's `minimize` procedure with the Huber SVM loss function defined above. Following Wu et

Dataset	NP baseline	AMP	H-F AMP	P-SGD	P-PSGD	P-SCPSGD	P-FW
<b>Low-Dimensional Binary Datasets (<math>\epsilon = 0.1</math>)</b>							
Synthetic-L	94.9	<b>89.3</b>	87.8	85.6	86.2	79.4	86.8
Adult	84.8	<b>79.6</b>	77.5	79.0	76.5	76.0	77.8
KDDCup99	99.1	98.7	<b>98.7</b> <sup>7</sup>	98.5	98.5	98.1	98.0
<b>Low-Dimensional Multi-class Datasets (<math>\epsilon = 1^8</math>)</b>							
Coverttype	71.5	<b>66.4</b>	65.3	64.3	62.3	62.7	63.3
MNIST	91.5	<b>74.7</b>	73.7	69.6	72.9	70.6	65.1
<b>High-Dimensional Datasets (<math>\epsilon = 0.1</math>)</b>							
Synthetic-H <sup>9</sup>	96.5	55.2	54.3	55.0	<b>56.6</b>	55.6	56.0
Gisette	96.6	69.9	67.9	65.7	<b>70.6</b>	66.8	66.8
Real-sim	93.6	<b>78.3</b>	76.7	73.6	71.8	69.7	78.3
RCV1 <sup>9</sup>	93.8	74.5	72.9	71.3	70.1	69.7	<b>75.8</b>
<b>Real-World Datasets (<math>\epsilon = 0.1</math>)</b>							
Dataset #1	75.3	75.3	75.3	75.3	<b>75.3</b> <sup>10</sup>	75.3	75.3
Dataset #2	72.2	<b>70.8</b>	70.6	70.8	70.3	70.2	68.6
Dataset #3	73.6	<b>71.3</b>	71.2	71.2	71.1	71.1	71.1
Dataset #4 <sup>9</sup>	81.9	81.5	81.3	<b>81.7</b>	81.5	81.2	81.2

Fig. 3. Accuracy results (in %) for Huber SVM. For each dataset, the result in bold represents the DP algorithm with the best accuracy for that dataset. A key for the abbreviations used for the algorithms is provided in Table III.

al. [12], we set  $h = 0.1$ . The results are shown in Figure 4, with more precise results in Figure 3. They demonstrate a similar trend to the earlier results for logistic regression, with our Approximate Minima Perturbation approach generally providing the highest accuracy. However, the advantage of Approximate Minima Perturbation is less pronounced in this setting.

<sup>7</sup>H-F AMP can outperform AMP when the data-independent strategy provides a better value for the privacy budget fraction  $f_1$  than the specific set of values we consider for tuning in AMP.

<sup>8</sup>We report the accuracy for  $\epsilon = 1$  for multi-class datasets, as compared to  $\epsilon = 0.1$  for datasets with binary classification, as multi-class classification is a more difficult task than binary classification.

<sup>9</sup>The numbers cited here do not reflect the trend for this dataset, as can be seen from Figure 3

<sup>10</sup>Slightly outperforms even the NP baseline, as can be seen from Figure 2.



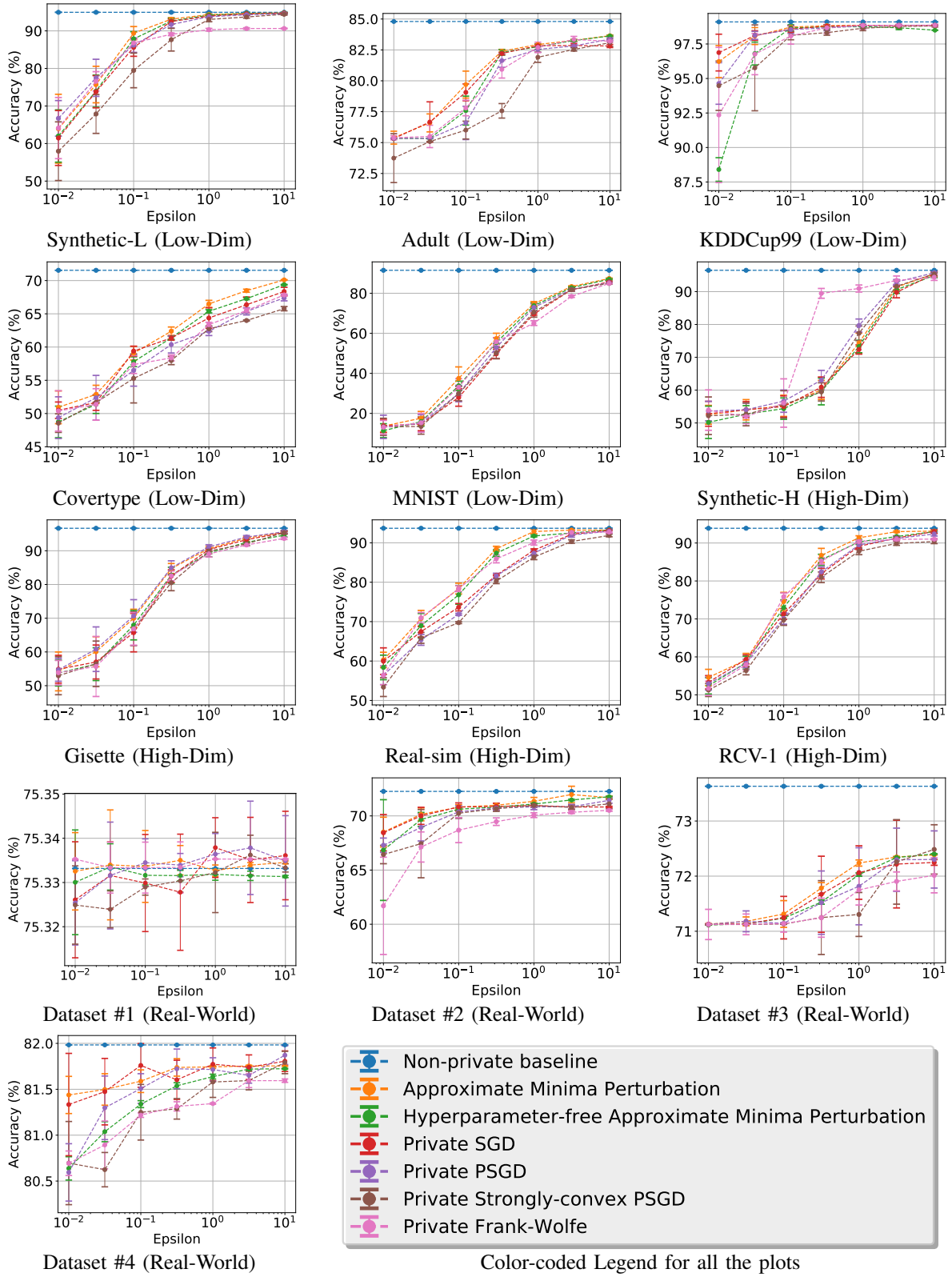


Fig. 4. Accuracy results for Huber SVM. Horizontal axis depicts varying values of  $\epsilon$ ; vertical axis shows accuracy (in %) on the testing set.

C. Pseudocodes for Algorithms evaluated in Section V

---

**Algorithm 2:** Differentially Private Minibatch Stochastic Gradient Descent [16], [15]

---

**Input:** Data set:  $D = \{d_1, \dots, d_n\}$ , loss function:  $\ell(\theta; D_i)$  with  $L_2$ -Lipschitz constant  $L$ , privacy parameters:  $(\epsilon, \delta)$ , number of iterations:  $T$ , minibatch size:  $k$ , learning rate function:  $\eta : [T] \rightarrow \mathbb{R}$ .

```

1  $\sigma^2 \leftarrow \frac{16L^2 T \log \frac{1}{\delta}}{n^2 \epsilon^2}$ 
2  $\theta_1 = 0^p$ 
3 for  $t = 1$  to  $T-1$  do
4    $s_1, \dots, s_k \leftarrow$  Sample  $k$  samples uniformly with replacement from  $D$ 
5    $b_t \sim \mathcal{N}(0, \sigma^2 I_{p \times p})$ 
6    $\theta_{t+1} = \theta_t - \eta(t) \left[ \left( \frac{1}{k} \sum_{i=1}^k \nabla \ell(\theta; s_i) \right) + b_t \right]$ 
7 end
8 Output  $\theta_T$ 
```

---



---

**Algorithm 3:** Differentially Private Permutation-based Stochastic Gradient Descent [12]

---

**Input:** Data set:  $D = \{d_1, \dots, d_n\}$ , loss function:  $\ell(\theta; D_i)$  with  $L_2$ -Lipschitz constant  $L$ , privacy parameters:  $(\epsilon, \delta)$ , number of passes:  $T$ , minibatch size:  $k$ , constant learning rate:  $\eta$ .

```

1  $\theta \leftarrow 0^p$ 
2 Let  $\tau$  be a random permutation of  $[n]$ 
3 for  $t = 1$  to  $T-1$  do
4   for  $b = 1$  to  $\frac{n}{k}$  do
5     Let  $s_1 = d_{\tau(bk)}, \dots, s_k = d_{\tau(b(k+1)-1)}$ 
6      $\theta \leftarrow \theta - \eta \left( \frac{1}{k} \sum_{i=1}^k \nabla \ell(\theta; s_i) \right)$ 
7   end
8 end
9  $\sigma^2 \leftarrow \frac{8T^2 L^2 \eta^2 \log(\frac{2}{\delta})}{k^2 \epsilon^2}$ 
10  $b \sim \mathcal{N}(0, \sigma^2 I_{p \times p})$ 
11 Output  $\theta_{priv} = \theta + b$ 
```

---



---

**Algorithm 4:** Differentially Private Strongly Convex Permutation-based Stochastic Gradient Descent [12]

---

**Input:** Data set:  $D = \{d_1, \dots, d_n\}$ , loss function:  $\ell(\theta; D_i)$  that is  $\xi$ -strongly convex and  $\beta$ -smooth with  $L_2$ -Lipschitz constant  $L$ , privacy parameters:  $(\epsilon, \delta)$ , number of passes:  $T$ , minibatch size:  $k$ .

```

1  $\theta \leftarrow 0^p$ 
2 Let  $\tau$  be a random permutation of  $[n]$ 
3 for  $t = 1$  to  $T-1$  do
4    $\eta_t \leftarrow \min \frac{1}{\beta}, \frac{1}{\xi t}$ 
5   for  $b = 1$  to  $\frac{n}{k}$  do
6     Let  $s_1 = d_{\tau(bk)}, \dots, s_k = d_{\tau(b(k+1)-1)}$ 
7      $\theta \leftarrow \theta - \eta_t \left( \frac{1}{k} \sum_{i=1}^k \nabla \ell(\theta; s_i) \right)$ 
8   end
9 end
10  $\sigma^2 \leftarrow \frac{8L^2 \log(\frac{2}{\delta})}{\xi^2 n^2 \epsilon^2}$ 
11  $b \sim \mathcal{N}(0, \sigma^2 I_{p \times p})$ 
12 Output  $\theta_{priv} = \theta + b$ 
```

---



---

**Algorithm 5:** Differentially Private Frank-Wolfe [44]

---

**Input:** Data set:  $D = \{d_1, \dots, d_n\}$ , loss function:  $L(\theta; D) = \frac{1}{n} \sum_{i=1}^n \ell(\theta; d_i)$  (with  $L_1$ -Lipshitz constant  $L$  for  $\ell$ ), privacy parameters:  $(\epsilon, \delta)$ , convex set:  $C = \text{conv}(S)$  with  $\|C\|_1$  denoting  $\max_{s \in S} \|s\|_1$  and  $S$  being the set of corners.

```

1 Choose an arbitrary  $\theta_1$  from  $C$ 
2  $\sigma^2 \leftarrow \frac{32L^2 \|C\|_1^2 T \log(1/\delta)}{n^2 \epsilon^2}$ 
3 for  $t = 1$  to  $T-1$  do
4    $\forall s \in S, \alpha_s \leftarrow \langle s, \nabla L(\theta_t; D) \rangle + \text{Lap}(\sigma)$ , where  $\text{Lap}(\lambda) \sim \frac{1}{2\lambda} e^{-|x|/\lambda}$ 
5    $\tilde{\theta}_t \leftarrow \arg \min_{s \in S} \alpha_s$ 
6    $\theta_{t+1} \leftarrow (1 - \eta_t) \theta_t + \eta_t \tilde{\theta}_t$ , where  $\eta_t = \frac{1}{t+1}$ 
7 end
8 Output  $\theta_{priv} = \theta_T$ 
```

---