

# DP-LSSGD: A Stochastic Optimization Method to Lift the Utility in Privacy-Preserving ERM

**Bao Wang**

WANGBAONJ@SGMAIL.COM

*Department of Mathematics and Scientific Computing & Imaging Institute  
University of Utah, Salt Lake City, UT, USA*

**Quanquan Gu**

QGU@CS.UCLA.EDU

*Department of Computer Science  
University of California, Los Angeles, CA, USA*

**March Boedihardjo**

MARCH@MATH.UCLA.EDU

*Department of Mathematics  
University of California, Los Angeles, CA, USA*

**Lingxiao Wang**

LINGXW@G.UCLA.EDU

*Department of Computer Science  
University of California, Los Angeles, CA, USA*

**Farzin Barekat**

FBAREKAT@MATH.UCLA.EDU

*Department of Mathematics  
University of California, Los Angeles, CA, USA*

**Stanley J. Osher**

SJO@MATH.UCLA.EDU

*Department of Mathematics  
University of California, Los Angeles, CA, USA*

## Abstract

Machine learning (ML) models trained by differentially private stochastic gradient descent (DP-SGD) have much lower utility than the non-private ones. To mitigate this degradation, we propose a DP Laplacian smoothing SGD (DP-LSSGD) to train ML models with differential privacy (DP) guarantees. At the core of DP-LSSGD is the Laplacian smoothing, which smooths out the Gaussian noise used in the Gaussian mechanism. Under the same amount of noise used in the Gaussian mechanism, DP-LSSGD attains the same DP guarantee, but in practice, DP-LSSGD makes training both convex and nonconvex ML models more stable and enables the trained models to generalize better. The proposed algorithm is simple to implement and the extra computational complexity and memory overhead compared with DP-SGD are negligible. DP-LSSGD is applicable to train a large variety of ML models, including DNNs. The code is available at <https://github.com/BaoWangMath/DP-LSSGD>.

**Keywords:** Laplacian Smoothing, Differential Privacy, Machine Learning, Optimization

## 1. Introduction

Many released machine learning (ML) models are trained on sensitive data that are often crowd-sourced or contain private information (Yuen et al., 2011; Feng et al., 2017; Liu et al., 2017). With overparameterization, deep neural nets (DNNs) can memorize the private training data, and it is possible to recover them and break the privacy by attacking the released models (Shokri et al., 2017).

For example, Fredrikson et al. demonstrated that a model-inversion attack can recover training images from a facial recognition system (Fredrikson et al., 2015). Protecting the private data is one of the most critical tasks in ML.

Differential privacy (DP) (Dwork et al., 2006) is a theoretically rigorous tool for designing algorithms on aggregated databases with a privacy guarantee. The idea is to add a certain amount of noise to randomize the output of a given algorithm such that the attackers cannot distinguish outputs of any two adjacent input datasets that differ in only one entry.

For repeated applications of additive noise based mechanisms, many tools have been invented to analyze the DP guarantee for the model obtained at the final stage. These include the basic and strong composition theorems and their refinements (Dwork et al., 2006, 2010; Kairouz et al., 2015), the moments accountant (Abadi et al., 2016), etc. Beyond the original notion of DP, there are also many other ways to define the privacy, e.g., local DP (Duchi et al., 2014), concentrated/zero-concentrated DP (Dwork and Rothblum, 2016; Bun and Steinke, 2016), and Rényi-DP (RDP) (Mironov, 2017).

Differentially private stochastic gradient descent (DP-SGD) reduces the utility of the trained models severely compared with SGD. As shown in Figure 1, the training and validation losses of the logistic regression on the MNIST dataset increase rapidly when the DP guarantee becomes stronger. The convolutional neural net (CNN)<sup>1</sup> trained by DP-SGD has much lower testing accuracy than the non-private one on the MNIST. We will discuss the detailed experimental settings in Section 4. A natural question raised from such performance degradations is:

*Can we improve DP-SGD, with negligible extra computational complexity and memory cost, such that it can be used to train general ML models with improved utility?*

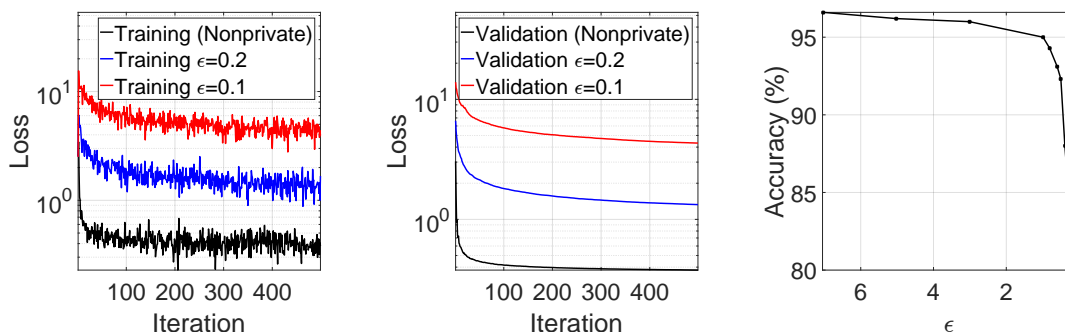


Figure 1: Training (left) and validation (middle) losses of the logistic regression on the MNIST trained by DP-SGD with  $(\epsilon, \delta = 10^{-5})$ -DP guarantee. (right): testing accuracy of a simple CNN on the MNIST trained by DP-SGD with  $(\epsilon, \delta = 10^{-5})$ -DP guarantee.

We answer the above question affirmatively by proposing differentially private Laplacian smoothing SGD (DP-LSSGD) to improve the utility in privacy-preserving empirical risk minimization (ERM). DP-LSSGD leverages the Laplacian smoothing (Osher et al., 2018) as a post-processing to smooth the injected Gaussian noise in the differentially private SGD (DP-SGD) to improve the convergence of DP-SGD in training ML models with DP guarantee.

1. [github.com/tensorflow/privacy/blob/master/tutorials/mnist\\_dpsgd\\_tutorial.py](https://github.com/tensorflow/privacy/blob/master/tutorials/mnist_dpsgd_tutorial.py)

### 1.1. Our Contributions

The main contributions of our work are highlighted as follows:

- We propose DP-LSSGD and prove its privacy and utility guarantees for convex/nonconvex optimizations. We prove that under the same privacy budget, DP-LSSGD achieves better utility, excluding a small term that is usually dominated by the other terms, than DP-SGD by a factor that is much less than one for convex optimization.
- We perform a large number of experiments logistic regression and CNN to verify the utility improvement by using DP-LSSGD. Numerical results show that DP-LSSGD remarkably reduces training and validation losses and improves the generalization of the trained private models.

In Table 1, we compare the privacy and utility guarantees of DP-LSSGD and DP-SGD. For the utility, the notation  $\tilde{O}(\cdot)$  hides the same constant and log factors for each bound. The constants  $d$  and  $n$  denote the dimension of the model’s parameters and the number of training points, respectively. The numbers  $\gamma$  and  $\beta$  are positive constants that are strictly less than one, and  $D_0, D_\sigma, G$  are positive constants, which will be defined in Section 3.

Table 1: Utility and Differential Privacy Guarantees.

Algorithm	DP	Assumption	Utility	Measurement	Reference
DP-SGD	$(\epsilon, \delta)$	convex	$\tilde{O}\left(\frac{\sqrt{(D_0+G^2)d}}{(\epsilon n)}\right)$	optimality gap	Bassily et al. (2014)
DP-SGD	$(\epsilon, \delta)$	nonconvex	$\tilde{O}\left(\sqrt{d}/(\epsilon n)\right)$	$\ell_2$ -norm of gradient	Zhang et al. (2017)
DP-LSSGD	$(\epsilon, \delta)$	convex	$\tilde{O}\left(\frac{\sqrt{\gamma(D_\sigma+G^2)d}}{(\epsilon n)}\right)$	optimality gap	<b>This Work</b>
DP-LSSGD	$(\epsilon, \delta)$	nonconvex	$\tilde{O}\left(\sqrt{\beta d}/(\epsilon n)\right)^1$	$\ell_2$ -norm of gradient	<b>This Work</b>

<sup>1</sup> Measured in the norm induced by  $\mathbf{A}_\sigma^{-1}$ , we will discuss this in detail in Section 4.

### 1.2. Related Work

There is a massive volume of research over the past decade on designing algorithms for privacy-preserving ML. Objective perturbation, output perturbation, and gradient perturbation are the three major approaches to perform ERM with a DP guarantee. Chaudhuri and Monteleoni (2008); Chaudhuri et al. (2011) considered both output and objective perturbations for privacy-preserving ERM, and gave theoretical guarantees for both privacy and utility for logistic regression and SVM. Song et al. (2013) numerically studied the effects of learning rate and batch size in DP-ERM. Wang et al. (2016) studied stability, learnability and other properties of DP-ERM. Lee and Kifer (2018) proposed an adaptive per-iteration privacy budget in concentrated DP gradient descent. The utility bound of DP-SGD has also been analyzed for both convex and nonconvex smooth objectives (Bassily et al., 2014; Zhang et al., 2017). Jayaraman et al. (2018) analyzed the excess empirical risk of DP-ERM in a distributed setting. Besides ERM, many other ML models have been made differentially private. These include: clustering (Su et al., 2015; Y. Wang and Singh, 2015; Balcan et al., 2017), matrix completion (Jain et al., 2018), online learning (Jain et al., 2012), sparse learning (Talwar et al., 2015; Wang and Gu, 2019), and topic modeling (Park et al., 2016). Gilbert and McMillan (2017) exploited the ill-conditionedness of inverse problems to design algorithms to release differentially private measurements of the physical system.

Shokri and Shmatikov (2015) proposed distributed selective SGD to train deep neural nets (DNNs) with a DP guarantee in a distributed system, however, the obtained privacy guarantee was very loose. Abadi et al. (2016) considered applying DP-SGD to train DNNs in a centralized setting. They clipped the gradient  $\ell_2$  norm to bound the sensitivity and invented the moment accountant to get better privacy loss estimation. Papernot et al. (2017) proposed Private Aggregation of Teacher Ensembles/PATE based on the semi-supervised transfer learning to train DNNs, and this framework improves both privacy and utility on top of the work by Abadi et al. (2016). Recently Papernot et al. (2018) introduced new noisy aggregation mechanisms for teacher ensembles that enable a tighter theoretical DP guarantee. The modified PATE is scalable to the large dataset and applicable to more diversified ML tasks.

Laplacian smoothing (LS) can be regarded as a denoising technique that performs post-processing on the Gaussian noise injected stochastic gradient. Denoising has been used in the DP earlier: Post-processing can enforce consistency of contingency table releases (Barak et al., 2007) and leads to accurate estimation of the degree distribution of private network (Hay et al., 2009). Nikolov et al. (2013) showed that post-processing by projecting linear regression solutions, when the ground truth solution is sparse, to a given  $\ell_1$ -ball can remarkably reduce the estimation error. Bernstein et al. (2017) used Expectation-Maximization to denoise a class of graphical models' parameters. Balle and Wang (2018) showed that in the output perturbation based differentially private algorithm design, denoising dramatically improves the accuracy of the Gaussian mechanism in the high-dimensional regime. To the best of our knowledge, we are the first to design a denoising technique on the Gaussian noise injected gradient to improve the utility of the trained private ML models.

### 1.3. Notation

We use boldface upper-case letters  $\mathbf{A}$ ,  $\mathbf{B}$  to denote matrices and boldface lower-case letters  $\mathbf{x}$ ,  $\mathbf{y}$  to denote vectors. For vectors  $\mathbf{x}$  and  $\mathbf{y}$  and positive definite matrix  $\mathbf{A}$ , we use  $\|\mathbf{x}\|_2$  and  $\|\mathbf{x}\|_{\mathbf{A}}$  to denote the  $\ell_2$ -norm and the induced norm by  $\mathbf{A}$ , respectively;  $\langle \mathbf{x}, \mathbf{y} \rangle$  denotes the inner product of  $\mathbf{x}$  and  $\mathbf{y}$ ; and  $\lambda_i(\mathbf{A})$  denotes the  $i$ -th largest eigenvalue of  $\mathbf{A}$ . We denote the set of numbers from 1 to  $n$  by  $[n]$ .  $\mathcal{N}(\mathbf{0}, \mathbf{I}_{d \times d})$  represents  $d$ -dimensional standard Gaussian.

### 1.4. Organization

This paper is organized in the following way: In Section 2, we introduce the DP-LSSGD algorithm. In Section 3, we analyze the privacy and utility guarantees of DP-LSSGD for both convex and nonconvex optimizations. We numerically verify the efficiency of DP-LSSGD in Section 4. We conclude this work and point out some future directions in Section 5.

## 2. Problem Setup and Algorithm

### 2.1. Laplacian Smoothing Stochastic Gradient Descent (LSSGD)

In this paper, we consider empirical risk minimization problem as follows. Given a training set  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  drawn from some unknown but fixed distribution, we aim to find an

empirical risk minimizer that minimizes the empirical risk as follows,

$$\min_{\mathbf{w}} F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}), \quad \mathbf{w} \in \mathbb{R}^d, \quad (1)$$

where  $F(\mathbf{w})$  is the empirical risk (a.k.a., training loss),  $f_i(\mathbf{w}) = \ell(\mathbf{w}; \mathbf{x}_i, y_i)$  is the loss function of a given ML model defined on the  $i$ -th training example  $(\mathbf{x}_i, y_i)$ , and  $\mathbf{w} \in \mathbb{R}^d$  is the model parameter we want to learn. Empirical risk minimization serves as the mathematical foundation for training many ML models that are mentioned above. The LSSGD (Osher et al., 2018) for solving (1) is given by

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \mathbf{A}_\sigma^{-1} \left( \frac{1}{b} \sum_{i_k \in \mathcal{B}_k} \nabla f_{i_k}(\mathbf{w}^k) \right), \quad (2)$$

where  $\eta$  is the learning rate,  $\nabla f_{i_k}$  denotes the stochastic gradient of  $F$  evaluated from the pair of input-output  $\{\mathbf{x}_{i_k}, y_{i_k}\}$ , and  $\mathcal{B}_k$  is a random subset of size  $b$  from  $[n]$ . Let  $\mathbf{A}_\sigma = \mathbf{I} - \sigma \mathbf{L}$  for  $\sigma \geq 0$  being a constant, where  $\mathbf{I} \in \mathbb{R}^{d \times d}$  and  $\mathbf{L} \in \mathbb{R}^{d \times d}$  are the identity and the discrete one-dimensional Laplacian matrix with periodic boundary condition, respectively. Therefore,

$$\mathbf{A}_\sigma := \begin{bmatrix} 1 + 2\sigma & -\sigma & 0 & \dots & 0 & -\sigma \\ -\sigma & 1 + 2\sigma & -\sigma & \dots & 0 & 0 \\ 0 & -\sigma & 1 + 2\sigma & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ -\sigma & 0 & 0 & \dots & -\sigma & 1 + 2\sigma \end{bmatrix} \quad (3)$$

When  $\sigma = 0$ , LSSGD reduces to SGD.

Note that  $\mathbf{A}_\sigma$  is positive definite with condition number  $1 + 4\sigma$  that is independent of  $\mathbf{A}_\sigma$ 's dimension, and LSSGD guarantees the same convergence rate as SGD in both convex and nonconvex optimization. Moreover, Laplacian smoothing (LS) can reduce the variance of SGD on-the-fly, and lead to better generalization in training many ML models including DNNs (Osher et al., 2018). For  $\mathbf{v} \in \mathbb{R}^d$ , let  $\mathbf{u} := \mathbf{A}_\sigma^{-1} \mathbf{v}$ , i.e.,  $\mathbf{v} = \mathbf{A}_\sigma \mathbf{u}$ . Note  $\mathbf{A}_\sigma$  is a convolution matrix, therefore,  $\mathbf{v} = \mathbf{A}_\sigma \mathbf{u} = \mathbf{u} - \sigma \mathbf{d} * \mathbf{u}$ , where  $\mathbf{d} = [-2, 1, 0, \dots, 0, 1]^T$  and  $*$  is the convolution operator. By the fast Fourier transform (FFT), we have

$$\mathbf{A}_\sigma^{-1} \mathbf{v} = \mathbf{u} = \text{ifft}(\text{fft}(\mathbf{v}) / (1 - \sigma \cdot \text{fft}(\mathbf{d}))),$$

where the division in the right hand side parentheses is performed in a coordinate wise way.

## 2.2. DP-LSSGD

DP ERM aims to learn a DP model,  $\mathbf{w}$ , for the problem (1). A common approach is injecting Gaussian noise into the stochastic gradient, and it resulting in the following DP-SGD

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \left( \frac{1}{b} \sum_{i_k \in \mathcal{B}_k} \nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n} \right), \quad (4)$$

where  $\mathbf{n}$  is the injected Gaussian noise for DP guarantee. Note that the LS matrix  $\mathbf{A}_\sigma^{-1}$  can remove the noise in  $\mathbf{v}$ . If we assume  $\mathbf{v}$  is the initial signal, then  $\mathbf{A}_\sigma^{-1} \mathbf{v}$  can be regarded as performing an approximate diffusion step on the initial noisy signal which removes the noise from  $\mathbf{v}$ . We will provide a detailed argument for the diffusion process in the appendix. As numerical illustrations, we consider the following two signals:

---

**Algorithm 1** DP-LSSGD

---

**Input:**  $f_i(\mathbf{w})$  is  $G$ -Lipschitz for  $i = 1, 2, \dots, n$ .  
 $\mathbf{w}^0$ : initial guess of  $\mathbf{w}$ ,  $(\epsilon, \delta)$ : the privacy budget,  $\eta$ : the step size,  $T$ : the total number of iterations.  
**Output:**  $(\epsilon, \delta)$ -differentially private classifier  $\mathbf{w}_{\text{priv}}$ .  
**for**  $k = 0, 1, \dots, T - 1$  **do**  
 $\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \mathbf{A}_\sigma^{-1} \left( \frac{1}{b} \sum_{i_k \in \mathcal{B}_k} \nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n} \right)$ , where  $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \nu^2 \mathbf{I})$  and  $\nu$  is defined in Theorem 2, and  $\mathcal{B}_k \subset [n]$ .  
**return**  $\mathbf{w}^T$

---

- 1D:  $\mathbf{v}_1 = \{\sin(2i\pi/100) + 0.1\mathcal{N}(0, 1) | i = 1, 2, \dots, 100\}$ .
- 2D:  $\mathbf{v}_2 = \{\sin(2i\pi/100) \sin(2j\pi/100) + 0.2\mathcal{N}(0, \mathbf{I}_{2 \times 2}) | i, j = 1, 2, \dots, 100\}$ .

We reshape  $\mathbf{v}_2$  into 1D with row-major ordering and then perform LS. Figure 2 shows that LS can remove noise efficiently. This noise removal property enables LSSGD to be more stable to the noise injected stochastic gradient, therefore improves training DP models with gradient perturbations.

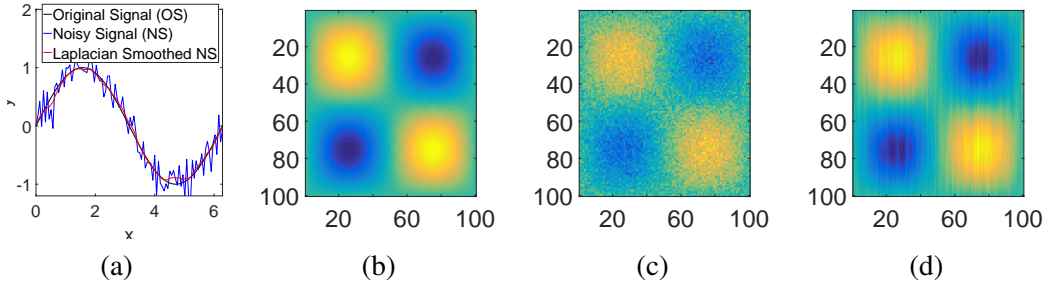


Figure 2: Illustration of LS ( $\sigma = 10$  for  $\mathbf{v}_1$  and  $\sigma = 100$  for  $\mathbf{v}_2$ ). (a): 1D signal sampled uniformly from  $\sin(x)$  for  $x \in [0, 2\pi]$ . (b), (c), (d): 2D original, noisy, and Laplacian Smoothed noisy signals sampled uniformly from  $\sin(x) \sin(y)$  for  $(x, y) \in [0, 2\pi] \times [0, 2\pi]$ .

We propose the following DP-LSSGD for solving (1) with DP guarantee

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \mathbf{A}_\sigma^{-1} \left( \frac{1}{b} \sum_{i_k \in \mathcal{B}_k} \nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n} \right). \tag{5}$$

In this scheme, we first inject the noise  $\mathbf{n}$  to the stochastic gradient  $\nabla f_{i_k}(\mathbf{w}^k)$ , and then apply the LS operator  $\mathbf{A}_\sigma^{-1}$  to denoise the noisy stochastic gradient,  $\nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n}$ , on-the-fly. We assume that each component function  $f_i$  in (1) is  $G$ -Lipschitz. The DP-LSSGD for finite-sum optimization is summarized in Algorithm 1. Compared with LSSGD, the main difference of DP-LSSGD lies in injecting Gaussian noise into the stochastic gradient, before applying the Laplacian smoothing, to guarantee the DP.

### 3. Main Theory

In this section, we present the privacy and utility guarantees for DP-LSSGD. The technical proofs are provided in the appendix.

**Definition 1** ( $(\epsilon, \delta)$ -DP) (*Dwork et al. (2006)*) A randomized mechanism  $\mathcal{M} : \mathcal{S}^N \rightarrow \mathcal{R}$  satisfies  $(\epsilon, \delta)$ -DP if for any two adjacent datasets  $S, S' \in \mathcal{S}^N$  differing by one element, and any output subset  $O \subseteq \mathcal{R}$ , it holds that

$$\mathbb{P}[\mathcal{M}(S) \in O] \leq e^\epsilon \cdot \mathbb{P}[\mathcal{M}(S') \in O] + \delta.$$

**Theorem 2 (Privacy Guarantee)** Suppose that each component function  $f_i$  is  $G$ -Lipschitz. Given the total number of iterations  $T$ , for any  $\delta > 0$  and privacy budget  $\epsilon$ , DP-LSSGD, with injected Gaussian noise  $\mathcal{N}(0, \nu^2)$  for each coordinate, satisfies  $(\epsilon, \delta)$ -DP with  $\nu^2 = 20T\alpha G^2 / (\mu n^2 \epsilon)$ , where  $\alpha = \log(1/\delta) / ((1 - \mu)\epsilon) + 1$ , if there exists  $\mu \in (0, 1)$  such that  $\alpha \leq \log(\mu n^3 \epsilon / (5b^3 T \alpha + \mu b n^2 \epsilon))$  and  $5b^2 T \alpha / (\mu n^2 \epsilon) \geq 1.5$ .

**Remark 3** It is straightforward to show that the noise in Theorem 2 is in fact also tight to guarantee the  $(\epsilon, \delta)$ -DP for DP-SGD. We will omit the dependence of  $\mu$  in our results in the rest of the paper since  $\mu$  is a constant.

For convex ERM, DP-LSSGD guarantees the following utility in terms of the gap between the ergodic average of the points along the DP-LSSGD path and the optimal solution  $\mathbf{w}^*$ .

**Theorem 4 (Utility Guarantee for convex optimization)** Suppose  $F$  is convex and each component function  $f_i$  is  $G$ -Lipschitz. Given  $\epsilon, \delta > 0$ , under the same conditions of Theorem 2 on  $\nu^2, \alpha$ , if we choose  $\eta_k = 1/\sqrt{T}$  and  $T = C_1(D_\sigma + G^2/b)n^2\epsilon^2 / (dG^2 \log(1/\delta))$ , where  $D_\sigma = \|\mathbf{w}^0 - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2$  and  $\mathbf{w}^*$  is the global minimizer of  $F$ , the DP-LSSGD output  $\tilde{\mathbf{w}} = \sum_{k=0}^{T-1} \eta_k / (\sum_{i=0}^{T-1} \eta_i) \mathbf{w}^k$  satisfies the following utility

$$\mathbb{E}(F(\tilde{\mathbf{w}}) - F(\mathbf{w}^*)) \leq \frac{C_2 G \sqrt{6\gamma(D_\sigma + G^2/b)d \log(1/\delta)}}{n\epsilon},$$

where  $\gamma = 1/d \sum_{i=1}^d 1/[1 + 2\sigma - 2\sigma \cos(2\pi i/d)]$ ,  $C_1, C_2$  are universal constants.

**Proposition 5** In Theorem 4,  $\gamma = \frac{1+\omega^d}{(1-\omega^d)\sqrt{4\sigma+1}}$ , where  $\omega = \frac{2\sigma+1-\sqrt{4\sigma+1}}{2\sigma} < 1$ . That is,  $\gamma$  converge to 0 almost exponentially as the dimension,  $d$ , increases.

**Remark 6** In the above utility bound for convex optimization, for different  $\sigma$  ( $\sigma = 0$  corresponds to DP-SGD), the only difference lies in the term  $\gamma(D_\sigma + G^2)$ . The first part  $\gamma D_\sigma$  depends on the gap between initialization  $\mathbf{w}^0$  and the optimal solution  $\mathbf{w}^*$ . The second part  $\gamma G^2$  decrease monotonically as  $\sigma$  increases.  $\sigma$  should be selected to get an optimal trade-off between these two parts. Based on our test on multi-class logistic regression for MNIST classification,  $\sigma \neq 0$  always outperforms the case when  $\sigma = 0$ .

For nonconvex ERM, DP-LSSGD has the following utility bound measured in gradient norm.



**Theorem 7 (Utility Guarantee for nonconvex optimization)** *Suppose that  $F$  is nonconvex and each component function  $f_i$  is  $G$ -Lipschitz and has  $L$ -Lipschitz continuous gradient. Given  $\epsilon, \delta > 0$ , under the same conditions of Theorem 2 on  $\nu^2, \alpha$ , if we choose  $\eta = 1/\sqrt{T}$  and  $T = C_1(D_F + LG^2/b)n^2\epsilon^2/(dLG^2 \log(1/\delta))$ , where  $D_F = F(\mathbf{w}^0) - F(\mathbf{w}^*)$  with  $\mathbf{w}^*$  being the global minimum of  $F$ , then the DP-LSSGD output  $\tilde{\mathbf{w}} = \sum_{k=0}^{T-1} \mathbf{w}^k/T$  satisfies the following utility*

$$\mathbb{E}\|\nabla F(\tilde{\mathbf{w}})\|_{\mathbf{A}_\sigma^{-1}}^2 \leq C_2 \frac{G\sqrt{\beta dL(2D_F + LG^2/b) \log(1/\delta)}}{n\epsilon},$$

where  $\beta = 1/d \sum_{i=1}^d 1/[1 + 2\sigma - 2\sigma \cos(2\pi i/d)]^2$ ,  $C_1, C_2$  are universal constants.

**Proposition 8** *In Theorem 7,  $\beta = \frac{2\omega^{2d+1} - \xi\omega^{2d} + 2\xi d\omega^d - 2\omega + \xi}{\sigma^2 \xi^3 (1 - \omega^d)^2}$ , where  $\omega = \frac{2\sigma + 1 - \sqrt{4\sigma + 1}}{2\sigma}$  and  $\xi = -\frac{\sqrt{1+4\sigma}}{\sigma}$ . Therefore,  $\beta \in (0, 1)$ .*

It is worth noting that if we use the  $\ell_2$ -norm instead of the induced norm, we have the following utility guarantee

$$\mathbb{E}\|\nabla F(\tilde{\mathbf{w}})\|_2^2 \leq \frac{\mathbb{E}\|\nabla F(\tilde{\mathbf{w}})\|_{\mathbf{A}_\sigma^{-1}}^2}{\lambda_{\min}(\mathbf{A}_\sigma^{-1})} \leq (1 + 4\sigma)\mathbb{E}\|\nabla F(\tilde{\mathbf{w}})\|_{\mathbf{A}_\sigma^{-1}}^2 \leq 4\zeta \frac{G\sqrt{6dL(2D_F + LG^2) \log(1/\delta)}}{n\epsilon}$$

where  $\zeta = \sqrt{\frac{1}{d} \sum_{i=1}^d \frac{(1+4\sigma)^2}{(1+2\sigma-2\sigma \cos(2\pi i/d))^2}} > 1$ . In the  $\ell_2$ -norm, DP-LSSGD has a bigger utility upper bound than DP-SGD (set  $\sigma = 0$  in  $\zeta$ ). However, this does not mean that DP-LSSGD has worse performance. We provide an example to support this claim in the appendix.

## 4. Experiments

In this section, we verify the efficiency of DP-LSSGD in training multi-class logistic regression and CNNs for MNIST and CIFAR10 classification. We use  $\mathbf{v} \leftarrow \mathbf{v}/\max(1, \|\mathbf{v}\|_2/C)$  (Abadi et al., 2016) to clip the gradient  $\ell_2$ -norms of the CNNs to  $C$ . The gradient clipping guarantee the Lipschitz condition for the objective functions. We train all the models below with  $(\epsilon, 10^{-5})$ -DP guarantee for different  $\epsilon$ . For Logistic regression we use the privacy budget given by Theorem 2, and for CNNs we use the privacy budget in the Tensorflow privacy (Andrew and et al., 2019). We checked that these two privacy budgets are consistent.

### 4.1. Logistic Regression for MNIST Classification

We ran 50 epochs of DP-LSSGD with learning rate scheduled as  $1/t$  with  $t$  being the index of the iteration to train the  $\ell_2$ -regularized (regularization constant  $10^{-4}$ ) multi-class logistic regression. We split the training data into 50K/10K with batch size 128 for cross-validation. We plot the evolution of training and validation loss over iterations for privacy budgets  $(0.2, 10^{-5})$  and  $(0.1, 10^{-5})$  in Figure 3. We see that the training loss curve of DP-SGD ( $\sigma = 0$ ) is much higher and more oscillatory (log-scale on the  $y$ -axis) than that of DP-LSSGD ( $\sigma = 1, 3$ ). Also, the validation loss of the model trained by DP-LSSGD decays faster and has a much smaller loss value than that of the model trained by DP-SGD. Moreover, when the privacy guarantee gets stronger, the utility improvement by DP-LSSGD becomes more significant.



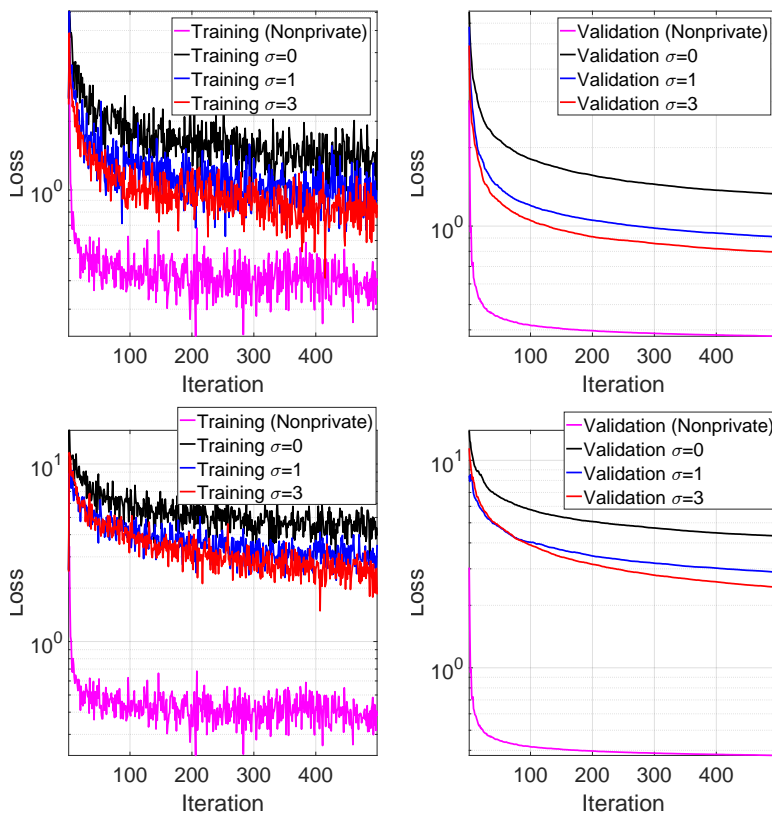


Figure 3: Training and validation losses of the multi-class logistic regression by DP-LSSGD. (a) and (b): training and validation curves with  $(0.2, 10^{-5})$ -DP guarantee; (c) and (d): training and validation curves with  $(0.1, 10^{-5})$ -DP guarantee. (Average over 5 runs)

Next, consider the testing accuracy of the multi-class logistic regression trained with  $(\epsilon, 10^{-5})$ -DP guarantee by DP-LSSGD includes  $\sigma = 0$ , i.e., DP-SGD. We list the test accuracy of logistic regression trained in different settings in Table 2. These results reveal that DP-LSSGD with  $\sigma = 1, 2, 3$  can improve the accuracy of the trained private model and also reduce the variance, especially when the privacy guarantee is very strong, e.g.,  $(0.1, 10^{-5})$ .

Table 2: Testing accuracy of the multi-class logistic regression trained by DP-LSSGD with  $(\epsilon, \delta = 10^{-5})$ -DP guarantee and different LS parameter  $\sigma$ . Unit: %. (5 runs)

$\epsilon$	0.30	0.25	0.20	0.15	0.10
$\sigma = 0$	81.74 $\pm$ 0.96	81.45 $\pm$ 1.59	78.92 $\pm$ 1.14	77.03 $\pm$ 0.69	73.49 $\pm$ 1.60
$\sigma = 1$	84.21 $\pm$ 0.51	83.27 $\pm$ 0.35	81.56 $\pm$ 0.79	79.46 $\pm$ 1.33	76.29 $\pm$ 0.53
$\sigma = 2$	84.23 $\pm$ 0.65	<b>83.65 <math>\pm</math> 0.76</b>	82.15 $\pm$ 0.59	80.77 $\pm$ 1.26	76.31 $\pm$ 0.93
$\sigma = 3$	<b>85.11 <math>\pm</math> 0.45</b>	82.97 $\pm$ 0.48	<b>82.22 <math>\pm</math> 0.28</b>	<b>80.81 <math>\pm</math> 1.03</b>	<b>77.13 <math>\pm</math> 0.77</b>

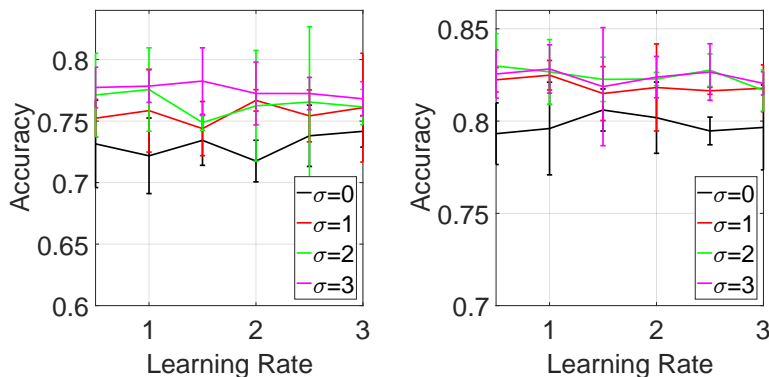


Figure 4: Accuracy of the logistic regression on MNIST when different learning rates are used to train the model. Left:  $(0.1, 10^{-5})$ -DP; Right:  $(0.2, 10^{-5})$ -DP.

#### 4.1.1. THE EFFECTS OF STEP SIZE

We know that the step size in DP-SGD/DP-LSSGD may affect the accuracy of the trained private models. We try different step size scheduling of the form  $\{a/t|a = 0.5, 1.0, 1.5, 2.0, 2.5, 3.0\}$ , where  $t$  is again the index of iteration, and all the other hyper-parameters are used the same as before. Figure. 4 plots the test accuracy of the logistic regression model trained with different learning rate scheduling and different privacy budget. We see that the private logistic regression model trained by DP-LSSGD always outperforms DP-SGD.

#### 4.2. CNN for MNIST and CIFAR10 Classification

In this subsection, we consider training a small CNN<sup>2</sup> with DP-guarantee for MNIST classification. We implement DP-LSSGD and DP-LSAdam (Kingma and Ba, 2014) (simply replace the noisy gradient in DP-Adam in the Tensorflow privacy with the Laplacian smoothed surrogate) into the Tensorflow privacy framework (Andrew and et al., 2019). We use the default learning rate 0.15 for DP-(LS)SGD and 0.001 for DP-(LS)Adam and decay them by a factor of 10 at the 10K-th iteration, norm clipping (1), batch size (256), and micro-batches (256). We vary the noise multiplier (NM), and larger NM guarantees stronger DP. As shown in Figure 5, the privacy budget increases at exactly the same speed (dashed red line) for four optimization algorithms. When the NM is large, i.e., DP-guarantee is strong, DP-SGD performs very well in the initial period. However, after a few epochs, the validation accuracy gets highly oscillatory and decays. DP-LSSGD can mitigate the training instability issue of DP-SGD. DP-Adam outperforms DP-LSSGD, and DP-LSAdam can further improve validation accuracy on top of DP-Adam.

Next, we consider the effects of the LS constant ( $\sigma$ ) and the learning rate in training the DP-CNN for MNIST classification. We fixed the NM to be 10, and run 60 epochs of DP-SGD and DP-LSSGD with different  $\sigma$  and different learning rate. We show the comparison of DP-SGD with DP-LSSGD with different  $\sigma$  in the left panel of Figure 7, and we see that as  $\sigma$  increases it becomes more stable in training CNNs with DP-guarantee even though initially it becomes slightly slower. In the middle

<sup>2</sup> [github.com/tensorflow/privacy/blob/master/tutorials/mnist\\_dpsgd\\_tutorial.py](https://github.com/tensorflow/privacy/blob/master/tutorials/mnist_dpsgd_tutorial.py)

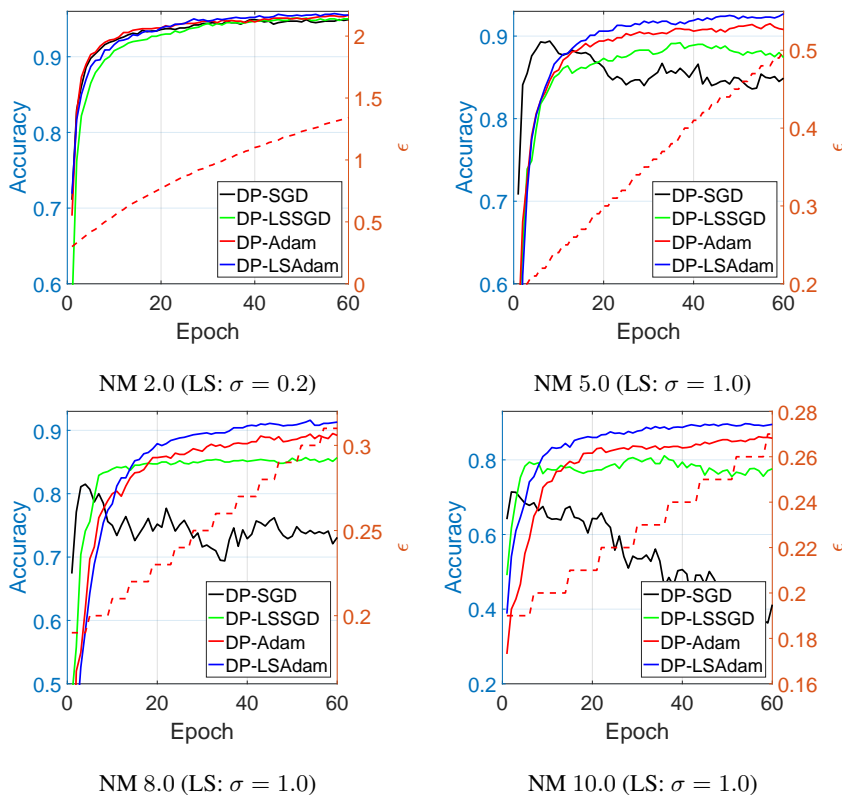


Figure 5: Performance comparison (validation accuracy) between different DP optimization algorithms in training CNN for MNIST classification with a fixed  $\delta = 10^{-5}$ .

panel of Figure 7, we plot the evolution of validation accuracy curves of the DP-CNN trained by DP-SGD and DP-LSSGD with different learning rate, where the solid lines represent results for DP-LSSGD and dashed lines for DP-SGD. DP-LSSGD outperforms DP-SGD in all learning rates tested, and DP-LSSGD is much more stable than DP-SGD when a larger learning rate is used.

Finally, we go back to the accuracy degradation problem raised in Figure 1. As shown in Figure 3, LS can efficiently reduce both training and validation losses in training multi-class logistic regression for MNIST classification. Moreover, as shown in the right panel of Figure 7, DP-LSSGD can improve the testing accuracy of the CNN used above significantly. In particular, DP-LSSGD improves the testing accuracy of CNN by 3.2% and 5.0% for  $(0.4, 10^{-5})$  and  $(0.2, 10^{-5})$ , respectively, on top of DP-SGD. DP-LSAdam can further boost test accuracy. All the accuracies associated with any given privacy budget in Figure 7 (right panel), are the optimal ones searched over the results obtained in the above experiments with different learning rate, number of epochs, and NM.

### 4.3. CNN for CIFAR10 Classification

In this section, we will show that LS can also improve the utility of the DP-CNN trained by DP-SGD and DP-Adam for CIFAR10 classification. We simply replace the CNN architecture used

above for MNIST classification with the benchmark architecture in the Tensorflow tutorial <sup>3</sup> for CIFAR10 classification. Also, we use the same set of parameters as that used for training DP-CNN for MNIST classification except we fixed the noise multiplier to be 2.0 and clip the gradient  $\ell_2$  norm to 3. As shown in Figure 6, LS can significantly improve the validation accuracy of the model trained by DP-SGD and DP-Adam, and the DP guarantee for all these algorithms are the same (dashed line in Figure 6).

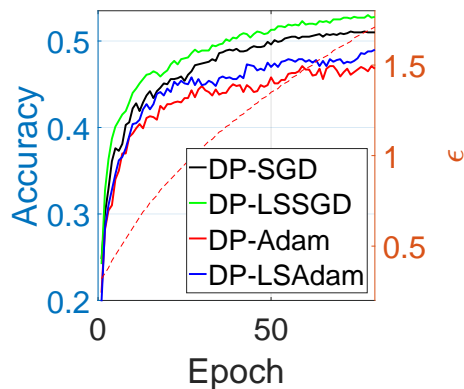


Figure 6: Performance comparison between different differentially private optimization algorithms in training CNN for CIFAR10 classification with a fixed  $\delta = 10^{-5}$ .

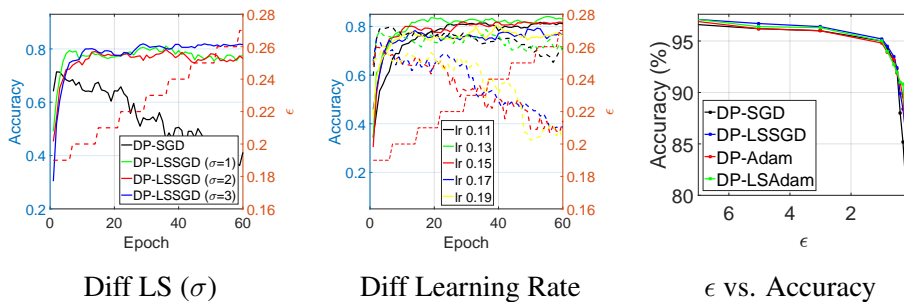


Figure 7: Left & middle panels: Contrasting performance (validation acc) of DP-SGD and DP-LSSGD with different  $\sigma$  and different learning rate. Right panel:  $\epsilon$  vs. Testing accuracy of the private models trained by different DP-optimization algorithms with a fixed  $\delta = 10^{-5}$ .

### 5. Conclusions

In this paper, we integrated Laplacian smoothing with DP-SGD for privacy-preserving ERM. The resulting algorithm is simple to implement and the extra computational cost compared with the DP-

<sup>3</sup> [github.com/tensorflow/models/tree/master/tutorials/image/cifar10](https://github.com/tensorflow/models/tree/master/tutorials/image/cifar10)

SGD is almost negligible. We show that DP-LSSGD can improve the utility of the trained private ML models both numerically and theoretically.

## Acknowledgments

This material is based on research sponsored by the National Science Foundation under grant number DMS-1924935 and DMS-1554564 (STROBE). The Air Force Research Laboratory under grant numbers FA9550-18-0167 and MURI FA9550-18-1-0502, the Office of Naval Research under grant number N00014-18-1-2527. QG is partially supported by the National Science Foundation under grant number SaTC-1717950.

## References

- M. Abadi, A. Chu, I. Goodfellow, H. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep Learning with Differential Privacy. In *23rd ACM Conference on Computer and Communications Security (CCS 2016)*, 2016.
- G. Andrew and et al. TensorFlow Privacy. <https://github.com/tensorflow/privacy>, 2019.
- M. Balcan, T. Dick, Y. Liang, W. Mou, and H. Zhang. Differentially Private Clustering in High-Dimensional Euclidean Spaces. In *34th International Conference on Machine Learning (ICML 2017)*, 2017.
- B. Balle and Y. Wang. Improving the Gaussian Mechanism for Differential Privacy: Analytical Calibration and Optimal Denoising. In *International Conference on Machine Learning*, pages 403–412, 2018.
- B. Barak, K. Chaudhuri, C. Dwork, S. Kale, F. McSherry, and K. Talwar. Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 273–282. ACM, 2007.
- R. Bassily, A. Smith, and A. Thakurta. Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds. In *55th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2014)*, 2014.
- G. Bernstein, R. McKenna, T. Sun, D. Sheldon, M. Hay, and G. Miklau. Differentially private learning of undirected graphical models using collective graphical models. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 478–487. JMLR. org, 2017.
- M. Bun and T. Steinke. Concentrated Differential Privacy: Simplifications, Extensions, and Lower Bounds. *ArXiv:1605.02065*, 2016.
- K. Chaudhuri and C. Monteleoni. Privacy-Preserving Logistic Regression. In *Advances in Neural Information Processing Systems (NIPS 2008)*, 2008.

- K. Chaudhuri, C. Monteleoni, and A. Sarwate. Differentially Private Empirical Risk Minimization. *Journal of Machine Learning Research*, 12, 2011.
- J. Duchi, M. Jordan, and M. Wainwright. Privacy Aware Learning. *Journal of the Association for Computing Machinery*, 61(6), 2014.
- C. Dwork and G. Rothblum. Concentrated Differentially Privacy. *ArXiv:1603.01887*, 2016.
- C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor. Our data, ourselves: Privacy via distributed noise generation. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pages 486–503. Springer, 2006.
- C. Dwork, G. Rothblum, and S. Vadhan. Boosting and Differential Privacy. In *51th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2010)*, 2010.
- W. Feng, Z. Yan, H. Zhang, K. Zeng, Y. Xiao, and Y. T. Hou. A Survey on Security, Privacy and Trust in Mobile Crowdsourcing. *IEEE Internet of Things Journal*, 2017.
- M. Fredrikson, S. Jha, and T. Ristenpart. Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures. In *22nd ACM SIGSAC Conference on Computer and Communications Security (CCS 2015)*, 2015.
- A. Gilbert and A. McMillan. Local Differential Privacy for Physical Sensor Data and Sparse Recovery. *arXiv:1706.05916*, 2017.
- M. Hay, C. Li, G. Miklau, and D. Jensen. Accurate estimation of the degree distribution of private networks. In *2009 Ninth IEEE International Conference on Data Mining*, pages 169–178. IEEE, 2009.
- P. Jain, P. Kothari, and A. Thakurta. Differentially Private Online Learning. In *25th Conference on Learning Theory (COLT 2012)*, 2012.
- P. Jain, O. Thakkar, and A. Thakurta. Differentially Private Matrix Completion. In *35th International Conference on Machine Learning (ICML 2018)*, 2018.
- B. Jayaraman, L. Wang, D. Evans, and Q. Gu. Distributed Learning without Distress: Privacy-Preserving Empirical Risk Minimization. In *Advances in Neural Information Processing Systems (NIPS 2018)*, 2018.
- P. Kairouz, S. Oh, and P. Viswanath. The Composition Theorem for Differential Privacy. In *32nd International Conference on Machine Learning (ICML 2015)*, 2015.
- D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- J. Lee and D. Kifer. Concentrated Differentially Private Gradient Descent with Adaptive per-iteration Privacy Budget. *ArXiv:1808.09501*, 2018.
- Y. Liu, K. Gadepalli, M. Norouzi, G. Dahl, T. Kohlberger, A. Boyko, S. Venugopalan, A. Timofeev, P. Nelson, G. Corrado, and et al. Detecting Cancer Metastases on Gigapixel Pathology Images. *arXiv:1703.02442*, 2017.

- I. Mironov. Renyi Differential Privacy. In *Computer Security Foundations Symposium (CSF), 2017 IEEE 30th*, pages 263–275. IEEE, 2017.
- 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*, pages 351–360. ACM, 2013.
- S. Osher, B. Wang, P. Yin, X. Luo, M. Pham, and A. Lin. Laplacian Smoothing Gradient Descent. *ArXiv:1806.06317*, 2018.
- N. Papernot, M. Abadi, U. Erlingsson, I. Goodfellow, and K. Talwar. Semisupervised Knowledge Transfer for Deep Learning from Private Training Data. In *5th International Conference on Learning Representation (ICLR 2017)*, 2017.
- N. Papernot, S. Song, I. Mironov, A. Raghunathan, K. Talwar, and U. Erlingsson. Scalable Private Learning with PATE. In *International Conference on Learning Representations (ICLR 2018)*, 2018.
- M. Park, J. Foulds, K. Chaudhuri, and M. Welling. Private Topic Modeling. *arXiv:1609:04120*, 2016.
- R. Shokri and V. Shmatikov. Privacy-Preserving Deep Learning. In *22nd ACM SIGSAC Conference on Computer and Communications Security (CCS 2015)*, 2015.
- R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership Inference Attacks Against Machine Learning Models. *Proceedings of the 2017 IEEE Symposium on Security and Privacy*, 2017.
- S. Song, K. Chaudhuri, and A. Sarwate. Stochastic Gradient Descent with Differentially Private Updates. In *GlobalSIP Conference*, 2013.
- D. Su, J. Cao, N. Li, E. Bertino, and H. Jin. Differentially Private k-Means Clustering. *arXiv:1504.05998*, 2015.
- K. Talwar, A. Thakurta, and L. Zhang. Nearly optimal private lasso. In *Advances in Neural Information Processing Systems*, pages 3025–3033, 2015.
- L. Wang and Q. Gu. Differentially Private Iterative Gradient Hard Thresholding for Sparse Learning. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, 2019.
- Lingxiao Wang, Bargav Jayaraman, David Evans, and Quanquan Gu. Efficient privacy-preserving nonconvex optimization. *arXiv preprint arXiv:1910.13659*, 2019.
- Y. Wang, J. Lei, and S. Fienberg. Learning with Differential Privacy: Stability, Learnability and the Sufficiency and Necessity of ERM Principle. *ArXiv:1502.06309*, 2016.
- Y. Wang Y. Wang and A. Singh. Differentially Private Subspace Clustering. In *Advances in Neural Information Processing Systems (NIPS 2015)*, 2015.
- M. Yuen, I. King, and K. Leung. A Survey of Crowdsourcing Systems. In *Proceedings of the IEEE international conference on social computing (Socialcom 2011)*, 2011.
- J. Zhang, K. Zheng, W. Mou, and L. Wang. Efficient Private ERM for Smooth Objectives. In *The Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI 2017)*, 2017.



## Appendix A. Proof of the Main Theorems

### A.1. Privacy Guarantee

To prove the privacy guarantee in Theorem 2, we first introduce the following  $\ell_2$ -sensitivity.

**Definition 9 ( $\ell_2$ -Sensitivity)** For any given function  $f(\cdot)$ , the  $\ell_2$ -sensitivity of  $f$  is defined by

$$\Delta(f) = \max_{\|S-S'\|_1=1} \|f(S) - f(S')\|_2,$$

where  $\|S - S'\|_1 = 1$  means the data sets  $S$  and  $S'$  differ in only one entry.

We will adapt the concepts and techniques of Rényi DP (RDP) to prove the DP-guarantee of the proposed DP-LSSGD.

**Definition 10 (RDP)** For  $\alpha > 1$  and  $\rho > 0$ , a randomized mechanism  $\mathcal{M} : \mathcal{S}^n \rightarrow \mathcal{R}$  satisfies  $(\alpha, \rho)$ -Rényi DP, i.e.,  $(\alpha, \rho)$ -RDP, if for all adjacent datasets  $S, S' \in \mathcal{S}^n$  differing by one element, we have

$$D_\alpha(\mathcal{M}(S) \parallel \mathcal{M}(S')) := \frac{1}{\alpha - 1} \log \mathbb{E} \left( \frac{\mathcal{M}(S)}{\mathcal{M}(S')} \right)^\alpha \leq \rho,$$

where the expectation is taken over  $\mathcal{M}(S')$ .

**Lemma 11 (Wang et al., 2019)** Given a function  $q : \mathcal{S}^n \rightarrow \mathcal{R}$ , the Gaussian Mechanism  $\mathcal{M} = q(S) + \mathbf{u}$ , where  $\mathbf{u} \sim N(0, \sigma^2 \mathbf{I})$ , satisfies  $(\alpha, \alpha \Delta^2(q)/(2\sigma^2))$ -RDP. In addition, if we apply the mechanism  $\mathcal{M}$  to a subset of samples using uniform sampling without replacement,  $\mathcal{M}$  satisfies  $(\alpha, 5\tau^2 \Delta^2(q)\alpha/\sigma^2)$ -RDP given  $\sigma'^2 = \sigma^2/\Delta^2(q) \geq 1.5$ ,  $\alpha \leq \log(1/\tau(1 + \sigma'^2))$ , where  $\tau$  is the subsample rate.

**Lemma 12 (Mironov, 2017)** If  $k$  randomized mechanisms  $\mathcal{M}_i : \mathcal{S}^n \rightarrow \mathcal{R}$ , for  $i \in [k]$ , satisfy  $(\alpha, \rho_i)$ -RDP, then their composition  $(\mathcal{M}_1(S), \dots, \mathcal{M}_k(S))$  satisfies  $(\alpha, \sum_{i=1}^k \rho_i)$ -RDP. Moreover, the input of the  $i$ -th mechanism can be based on outputs of the previous  $(i - 1)$  mechanisms.

**Lemma 13** If a randomized mechanism  $\mathcal{M} : \mathcal{S}^n \rightarrow \mathcal{R}$  satisfies  $(\alpha, \rho)$ -RDP, then  $\mathcal{M}$  satisfies  $(\rho + \log(1/\delta)/(\alpha - 1), \delta)$ -DP for all  $\delta \in (0, 1)$ .

With the definition (Def. 10) and guarantees of RDP (Lemmas 11 and 12), and the connection between RDP and  $(\epsilon, \delta)$ -DP (Lemma 13), we can prove the following DP-guarantee for DP-LSSGD.

**Proof** [Proof of Theorem 2] Let us denote the update of DP-SGD and DP-LSSGD at the  $k$ -th iteration starting from any given points  $\mathbf{w}^k$  and  $\tilde{\mathbf{w}}^k$ , respectively, as

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta_k \left( \frac{1}{b} \sum_{i_k \in \mathcal{B}_k} \nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n} \right), \quad (6)$$

and

$$\tilde{\mathbf{w}}^{k+1} = \tilde{\mathbf{w}}^k - \eta_k \mathbf{A}_\sigma^{-1} \left( \frac{1}{b} \sum_{i_k \in \mathcal{B}_k} \nabla f_{i_k}(\tilde{\mathbf{w}}^k) + \mathbf{n} \right), \quad (7)$$

where  $\mathcal{B}_k$  is a mini batch that are drawn uniformly from  $[n]$ , and  $|\mathcal{B}_k| = b$  is the mini batch size.

We will show that with the aforementioned Gaussian noise  $\mathcal{N}(0, \nu^2)$  for each coordinate of  $\mathbf{n}$ , the output of DP-SGD,  $\tilde{\mathbf{w}}$ , after  $T$  iterations is  $(\epsilon, \delta)$ -DP. Let us consider the mechanism  $\hat{\mathcal{M}}_k = \frac{1}{b} \sum_{i_k \in \mathcal{B}_k} \nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n}$ , and  $\mathcal{M}_k = \frac{n}{b} \nabla F(\mathbf{w}^k) + \mathbf{n}$  with the query  $\mathbf{q}_k = \frac{n}{b} \nabla F(\mathbf{w}^k)$ . We have the  $\ell_2$ -sensitivity of  $\mathbf{q}_k$  as  $\Delta(\mathbf{q}_k) = \|\nabla f_{i_k}(\mathbf{w}^k) - \nabla f_{i'_k}(\mathbf{w}^k)\|_2 \leq \frac{2G}{b}$ . According to Lemma 11, if we add noise with variance

$$\nu^2 = \frac{20T\alpha G^2}{n^2\epsilon\mu},$$

the mechanism  $\mathcal{M}_k$  will satisfy  $(\alpha, (n^2\epsilon\mu/b^2)/(10T))$ -RDP. By post-processing theorem, we immediately have that under the same noise,  $\tilde{\mathcal{M}}_k = \mathbf{A}_\sigma^{-1}(\nabla F(\mathbf{w}^k) + \mathbf{n})$  also satisfies  $(\alpha, (n^2\epsilon\mu/b^2)/(10T))$ -RDP. According to Lemma 11,  $\hat{\mathcal{M}}_k$  will satisfy  $(\alpha, \mu\epsilon/T)$ -RDP provided that  $\nu^2/\Delta(\mathbf{q}_k)^2 \geq 1.5$ , because  $\tau = b/n$ . Let  $\alpha = \log(1/\delta)/((1-\mu)\epsilon) + 1$ , we obtain that  $\hat{\mathcal{M}}_k$  satisfies  $(\log(1/\delta)/((1-\mu)\epsilon) + 1, \mu\epsilon/T)$ -RDP as long as we have

$$\frac{\nu^2}{\Delta(\mathbf{q}_k)^2} = \frac{5T\alpha b^2}{n^2\epsilon\mu} \geq 1.5.$$

In addition, we have

$$\frac{1}{\tau(1 + \nu^2/\Delta(\mathbf{q}_k)^2)} = \frac{\mu n^3 \epsilon}{5b^3 T \alpha + \mu b n^2 \epsilon},$$

which implies that  $\alpha = \log(1/\delta)/((1-\mu)\epsilon) + 1 \leq \log(\mu n^3 \epsilon / (5b^3 T \alpha + \mu b n^2 \epsilon))$ . Therefore, according to Lemma 12, we have  $\mathbf{w}^k$  satisfies  $(\log(1/\delta)/((1-\mu)\epsilon) + 1, k\mu\epsilon/T)$ -RDP. Finally, by Lemma 13, we have  $\mathbf{w}^k$  satisfies  $(k\mu\epsilon/T + (1-\mu)\epsilon, \delta)$ -DP. Therefore, the output of DP-SGD,  $\tilde{\mathbf{w}}$ , is  $(\epsilon, \delta)$ -DP.  $\blacksquare$

**Remark 14** *In the above proof, we used the following estimate of the  $\ell_2$  sensitivity*

$$\Delta(\mathbf{q}_k) = \|\mathbf{A}_\sigma^{-1} \nabla f_i(\mathbf{w}^k) - \mathbf{A}_\sigma^{-1} \nabla f_{i'}(\mathbf{w}^k)\|_2/n \leq 2G/n.$$

*Indeed, let  $\mathbf{g} = \nabla f_i(\mathbf{w}^k) - \nabla f_{i'}(\mathbf{w}^k)$  and  $\mathbf{d} = \mathbf{A}_\sigma^{-1} \mathbf{g}$ , then according to Osher et al. (2018) we have*

$$\|\mathbf{d}\|_2 + 2\sigma \frac{\|\mathbf{D}_+ \mathbf{d}\|_2^2}{d} + \sigma^2 \frac{\|\mathbf{L} \mathbf{d}\|_2^2}{d} = \|\mathbf{g}\|_2,$$

where  $d$  is the dimension of  $\mathbf{d}$ , and

$$\mathbf{D}_+ = \begin{bmatrix} -1 & 1 & 0 & \dots & 0 & 0 \\ 0 & -1 & 1 & \dots & 0 & 0 \\ 0 & 0 & -1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & \dots & 0 & -1 \end{bmatrix}.$$

*Moreover, if we assume the  $\mathbf{g}$  is randomly sampled from a unit ball in a high dimensional space, then a high probability estimation of the compression ratio of the  $\ell_2$  norm can be derived from Lemma. 16.*

*Numerical experiments show that  $\|\mathbf{A}_\sigma^{-1} \nabla f_i(\mathbf{w}^k) - \mathbf{A}_\sigma^{-1} \nabla f_{i'}(\mathbf{w}^k)\|_2$  is much less than  $\|\nabla f_i(\mathbf{w}^k) - \nabla f_{i'}(\mathbf{w}^k)\|_2$ , so for the above noise, it can give much stronger privacy guarantee.*

## A.2. Utility Guarantee – Convex Optimization

To prove the utility guarantee for convex optimization, we first show that the LS operator compresses the  $\ell_2$  norm of any given Gaussian random vector with a specific ratio in expectation.

**Lemma 15** *Let  $\mathbf{x} \in \mathbb{R}^d$  be the standard Gaussian random vector. Then*

$$\mathbb{E}\|\mathbf{x}\|_{\mathbf{A}_\sigma^{-1}}^2 = \sum_{i=1}^d \frac{1}{1 + 2\sigma - 2\sigma \cos(2\pi i/d)},$$

where  $\|\mathbf{x}\|_{\mathbf{A}_\sigma^{-1}}^2 \doteq \langle \mathbf{x}, \mathbf{A}_\sigma^{-1} \mathbf{x} \rangle$  is the square of the induced norm of  $\mathbf{x}$  by the matrix  $\mathbf{A}_\sigma^{-1}$ .

**Proof** [Proof of Lemma 15] Let the eigenvalue decomposition of  $\mathbf{A}_\sigma^{-1}$  be  $\mathbf{A}_\sigma^{-1} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ , where  $\mathbf{\Lambda}$  is a diagonal matrix with  $\Lambda_{ii} = \frac{1}{1+2\sigma-2\sigma \cos(2\pi i/d)}$ . We have

$$\begin{aligned} \mathbb{E}\|\mathbf{x}\|_{\mathbf{A}_\sigma^{-1}}^2 &= \mathbb{E}[\text{Tr}(\mathbf{x}^T \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T \mathbf{x})] \\ &= \sum_{i=1}^d \Lambda_{ii} \\ &= \sum_{i=1}^d \frac{1}{1 + 2\sigma - 2\sigma \cos(2\pi i/d)} = \gamma. \end{aligned}$$

■

**Proof** [Proof of Theorem 4] Recall that we have the following update rule  $\mathbf{w}^{k+1} = \mathbf{w}^k - \eta_k \mathbf{A}_\sigma^{-1} (\nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n})$ , where  $i_k$  are drawn uniformly from  $[n]$ , and  $\mathbf{n} \sim \mathcal{N}(0, \nu^2 \mathbf{I})$ . Let  $\nabla f_{\mathcal{B}_k} = \sum_{i_k \in \mathcal{B}_k} \nabla f_{i_k}(\mathbf{w}^k)/b$ , observe that

$$\begin{aligned} \|\mathbf{w}^{k+1} - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2 &= \|\mathbf{w}^k - \eta_k \mathbf{A}_\sigma^{-1} (\nabla f_{\mathcal{B}_k}(\mathbf{w}^k) + \mathbf{n}) - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2 \\ &= \|\mathbf{w}^k - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2 + \eta_k^2 (\|\mathbf{A}_\sigma^{-1} (\nabla f_{\mathcal{B}_k}(\mathbf{w}^k) - \nabla F(\mathbf{w}^k) + \nabla F(\mathbf{w}^k))\|_{\mathbf{A}_\sigma}^2 + \|\mathbf{A}_\sigma^{-1} \mathbf{n}\|_{\mathbf{A}_\sigma}^2 \\ &\quad + 2\langle \mathbf{A}_\sigma^{-1} \nabla f_{\mathcal{B}_k}(\mathbf{w}^k), \mathbf{n} \rangle - 2\eta_k \langle \nabla f_{\mathcal{B}_k}(\mathbf{w}^k) + \mathbf{n}, \mathbf{w}^k - \mathbf{w}^* \rangle. \end{aligned}$$

Taking expectation with respect to  $\mathcal{B}_k$  and  $\mathbf{n}$  given  $\mathbf{w}^k$ , we have

$$\begin{aligned} \mathbb{E}\|\mathbf{w}^{k+1} - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2 &= \mathbb{E}\|\mathbf{w}^k - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2 - 2\eta_k \mathbb{E}\langle \nabla F(\mathbf{w}^k), \mathbf{w}^k - \mathbf{w}^* \rangle + \eta_k^2 \mathbb{E}\|\nabla f_{\mathcal{B}_k}(\mathbf{w}^k) - \nabla F(\mathbf{w}^k)\|_{\mathbf{A}_\sigma^{-1}}^2 \\ &\quad + \eta_k^2 \mathbb{E}\|\nabla F(\mathbf{w}^k)\|_{\mathbf{A}_\sigma^{-1}}^2 + \eta_k^2 \mathbb{E}\|\mathbf{n}\|_{\mathbf{A}_\sigma^{-1}}^2. \end{aligned}$$

In addition, we have

$$\mathbb{E}\|\nabla f_{\mathcal{B}_k}(\mathbf{w}^k) - \nabla F(\mathbf{w}^k)\|_{\mathbf{A}_\sigma^{-1}}^2 \leq \mathbb{E}\|\nabla f_{\mathcal{B}_k}(\mathbf{w}^k) - \nabla F(\mathbf{w}^k)\|_2^2 \leq \frac{G^2}{b}, \quad (8)$$

and

$$\left(1 - \frac{L\eta_k}{2}\right) \eta_k \|\nabla F(\mathbf{w}^k)\|_2^2 \leq F(\mathbf{w}^k) - F(\mathbf{w}^*), \quad (9)$$

which implies that

$$\eta_k^2 \mathbb{E} \|\nabla F(\mathbf{w}^k)\|_{\mathbf{A}_\sigma^{-1}}^2 \leq \eta_k^2 \mathbb{E} \|\nabla F(\mathbf{w}^k)\|_2^2 \leq \left( \frac{2}{2 - L\eta_k} \right) \eta_k \mathbb{E}(F(\mathbf{w}^k) - F(\mathbf{w}^*)) \leq \frac{4}{3} \eta_k \mathbb{E}(F(\mathbf{w}^k) - F(\mathbf{w}^*)),$$

where the last inequality is due to the fact that  $\eta_t \leq 1/(2L)$ . Therefore, we have

$$\mathbb{E} \|\mathbf{w}^{k+1} - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2 \leq \mathbb{E} \|\mathbf{w}^k - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2 - \frac{2}{3} \eta_k \mathbb{E}(F(\mathbf{w}^k) - F(\mathbf{w}^*)) + \eta_k^2 (G^2/b + \gamma d\nu^2),$$

where the inequality is due to the convexity of  $F$ , and Lemma 15. It implies that

$$\frac{2}{3} \eta_k \mathbb{E}(F(\mathbf{w}^k) - F(\mathbf{w}^*)) \leq (\mathbb{E} \|\mathbf{w}^k - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2 - \mathbb{E} \|\mathbf{w}^{k+1} - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2) + \eta_k^2 (G^2/b + \gamma d\nu^2).$$

Now taking the full expectation and summing up over  $T$  iterations, we have

$$\sum_{k=0}^{T-1} \frac{2}{3} \eta_k \mathbb{E}(F(\mathbf{w}^k) - F(\mathbf{w}^*)) \leq D_\sigma + \sum_{k=0}^{T-1} \eta_k^2 (G^2/b + \gamma d\nu^2),$$

where  $D_\sigma = \|\mathbf{w}^0 - \mathbf{w}^*\|_{\mathbf{A}_\sigma}^2$ . Let  $v_k = \eta_k / (\sum_{k=0}^{T-1} \eta_k)$ , we have

$$\sum_{k=0}^{T-1} v_k \mathbb{E}(F(\mathbf{w}^k) - F(\mathbf{w}^*)) \leq \frac{D_\sigma + \sum_{k=0}^{T-1} \eta_k^2 (G^2/b + \gamma d\nu^2)}{2 \sum_{k=0}^{T-1} \eta_k / 3}.$$

According to the definition of  $\tilde{\mathbf{w}}$  and the convexity of  $F$ , we obtain

$$\begin{aligned} \mathbb{E}(F(\tilde{\mathbf{w}}) - F(\mathbf{w}^*)) &\leq \frac{D_\sigma + \sum_{k=0}^{T-1} \eta_k^2 (G^2/b + \gamma d\nu^2)}{2 \sum_{k=0}^{T-1} \eta_k / 3} \\ &\leq \frac{D_\sigma + \sum_{k=0}^{T-1} \eta_k^2 G^2/b}{2 \sum_{k=0}^{T-1} \eta_k / 3} + \frac{\sum_{k=0}^{T-1} \eta_k^2}{2 \sum_{k=0}^{T-1} \eta_k / 3} \cdot \frac{20\gamma d T G^2 \log(1/\delta)}{n^2 \epsilon^2 \mu(1 - \mu)}. \end{aligned}$$

Let  $\eta = 1/\sqrt{T}$  and  $T = C_1(D_\sigma + G^2/b)n^2\epsilon^2/(\gamma d G^2 \log(1/\delta))$ , we can obtain that

$$\mathbb{E}(F(\tilde{\mathbf{w}}) - F(\mathbf{w}^*)) \leq \frac{C_2 G \sqrt{\gamma(D_\sigma + G^2/b)d \log(1/\delta)}}{n\epsilon},$$

where  $C_1, C_2$  are universal constants. ■

### A.3. Utility Guarantee – Nonconvex Optimization

To prove the utility guarantee for nonconvex optimization, we need the following lemma, which shows that the LS operator compresses the  $\ell_2$  norms of any given Gaussian random vector with a specific ratio in expectation.

**Lemma 16** *Let  $\mathbf{x} \in \mathbb{R}^d$  be the standard Gaussian random vector. Then*

$$\mathbb{E} \|\mathbf{A}_\sigma^{-1} \mathbf{x}\|_2^2 = \sum_{i=1}^d \frac{1}{(1 + 2\sigma - 2\sigma \cos(2\pi i/d))^2}.$$

**Proof** [Proof of Lemma 16] Let the eigenvalue decomposition of  $\mathbf{A}_\sigma^{-1}$  be  $\mathbf{A}_\sigma^{-1} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ , where  $\mathbf{\Lambda}$  is a diagonal matrix with  $\Lambda_{ii} = \frac{1}{1+2\sigma-2\sigma\cos(2\pi i/n)}$ . We have

$$\begin{aligned}\mathbb{E}\|\mathbf{A}_\sigma^{-1}\mathbf{x}\|_2^2 &= \mathbb{E}[\text{Tr}(\mathbf{x}^\top\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top\mathbf{x})] \\ &= \mathbb{E}[\text{Tr}(\mathbf{x}^\top\mathbf{U}\mathbf{\Lambda}^2\mathbf{U}^\top\mathbf{x})] \\ &= \sum_{i=1}^d \Lambda_{ii}^2 \\ &= \sum_{i=1}^d \frac{1}{(1+2\sigma-2\sigma\cos(2\pi i/d))^2} = \beta.\end{aligned}$$

■

**Proof** [Proof of Theorem 7] Recall that we have the following update rule  $\mathbf{w}^{t+1} = \mathbf{w}^k - \eta_k \mathbf{A}_\sigma^{-1}(\nabla f_{i_k}(\mathbf{w}^k) + \mathbf{n})$ , where  $i_k$  are drawn uniformly from  $[n]$ , and  $\mathbf{n} \sim \mathcal{N}(0, \nu^2 \mathbf{I})$ . Let  $\nabla f_{\mathcal{B}_k} = \sum_{i_k \in \mathcal{B}_k} \nabla f_{i_k}(\mathbf{w}^k)/b$ , since  $F$  is  $L$ -smooth, we have

$$\begin{aligned}F(\mathbf{w}^{k+1}) &\leq F(\mathbf{w}^k) + \langle \nabla F(\mathbf{w}^k), \mathbf{w}^{k+1} - \mathbf{w}^k \rangle + \frac{L}{2} \|\mathbf{w}^{k+1} - \mathbf{w}^k\|_2^2 \\ &= F(\mathbf{w}^k) - \eta_k \langle \nabla F(\mathbf{w}^k), \mathbf{A}_\sigma^{-1}(\nabla f_{\mathcal{B}_k}(\mathbf{w}^k) + \mathbf{n}) \rangle \\ &\quad + \frac{\eta_k^2 L}{2} \left( \|\mathbf{A}_\sigma^{-1}(\nabla f_{\mathcal{B}_k}(\mathbf{w}^k) - \nabla F(\mathbf{w}^k))\|_2^2 + \|\mathbf{A}_\sigma^{-1}\mathbf{n}\|_2^2 + 2\langle \mathbf{A}_\sigma^{-1}\nabla f_{\mathcal{B}_k}(\mathbf{w}^k), \mathbf{A}_\sigma^{-1}\mathbf{n} \rangle \right).\end{aligned}$$

Taking expectation with respect to  $\mathcal{B}_k$  and  $\mathbf{n}$  given  $\mathbf{w}^k$ , we have

$$\begin{aligned}\mathbb{E}F(\mathbf{w}^{k+1}) &\leq \mathbb{E}F(\mathbf{w}^k) - \eta_k \mathbb{E}\langle \nabla F(\mathbf{w}^k), \mathbf{A}_\sigma^{-1}\nabla f_{\mathcal{B}_k}(\mathbf{w}^k) \rangle + \frac{\eta_k^2 L}{2} \left( \mathbb{E}\|\mathbf{A}_\sigma^{-1}(\nabla f_{\mathcal{B}_k}(\mathbf{w}^k) - \nabla F(\mathbf{w}^k))\|_2^2 \right. \\ &\quad \left. + \mathbb{E}\|\mathbf{A}_\sigma^{-1}\nabla F(\mathbf{w}^k)\|_2^2 + \mathbb{E}\|\mathbf{A}_\sigma^{-1}\mathbf{n}\|_2^2 \right) \\ &\leq \mathbb{E}F(\mathbf{w}^k) - \eta_k \left( 1 - \frac{\eta_k L}{2} \right) \mathbb{E}\|\nabla F(\mathbf{w}^k)\|_{\mathbf{A}_\sigma^{-1}}^2 + \frac{\eta_k^2 L}{2} (G^2/b + d\beta\nu^2) \\ &\leq \mathbb{E}F(\mathbf{w}^k) - \frac{\eta_k}{2} \mathbb{E}\|\nabla F(\mathbf{w}^k)\|_{\mathbf{A}_\sigma^{-1}}^2 + \frac{\eta_k^2 L(G^2 + d\beta\nu^2)}{2},\end{aligned}$$

where the second inequality uses Lemma 16, the inequality (8), and the last inequality is due to  $1 - \eta_k L/2 > 1/2$ . Now taking the full expectation and summing up over  $T$  iterations, we have

$$\mathbb{E}F(\mathbf{w}^T) \leq F(\mathbf{w}^0) - \sum_{k=1}^{T-1} \frac{\eta_k}{2} \mathbb{E}\|\nabla F(\mathbf{w}^k)\|_{\mathbf{A}_\sigma^{-1}}^2 + \sum_{k=1}^{T-1} \frac{\eta_k^2 L(G^2/b + d\beta\nu^2)}{2}.$$

If we choose fix step size, i.e.,  $\eta_k = \eta$ , and rearranging the above inequality, and using  $F(\mathbf{w}^0) - \mathbb{E}F(\mathbf{w}^T) \leq F(\mathbf{w}^0) - F(\mathbf{w}^*)$ , we get

$$\frac{1}{T} \sum_{k=1}^{T-1} \mathbb{E}\|\nabla F(\mathbf{w}^k)\|_{\mathbf{A}_\sigma^{-1}}^2 \leq \frac{2}{\eta T} (F(\mathbf{w}^0) - F(\mathbf{w}^*)) + \eta L(G^2/b + d\beta\nu^2),$$

which implies that

$$\begin{aligned}\mathbb{E}\|\nabla F(\tilde{\mathbf{w}})\|_{\mathbf{A}_\sigma^{-1}}^2 &\leq \frac{2D_F}{\eta T} + \eta L(G^2/b + d\beta\nu^2) \\ &\leq \frac{2D_F}{\eta T} + \eta L\left(G^2/b + \frac{20d\beta TG^2 \log(1/\delta)}{n^2\epsilon^2\mu(1-\mu)}\right).\end{aligned}$$

Let  $\eta = 1/\sqrt{T}$  and  $T = C_1(2D_F + LG^2/b)n^2\epsilon^2/(dL\beta G^2 \log(1/\delta))$ , where  $D_F = F(\mathbf{w}^0) - F(\mathbf{w}^*)$ , we obtain

$$\mathbb{E}\|\nabla F(\tilde{\mathbf{w}})\|_{\mathbf{A}_\sigma^{-1}}^2 \leq C_2 \frac{G\sqrt{\beta dL(2D_F + LG^2/b)\log(1/\delta)}}{n\epsilon},$$

where  $C_1, C_2$  are universal constants. ■

It is worth noting that if we use the  $\ell_2$ -norm instead of the induced norm, we have the following utility guarantee

$$\mathbb{E}\|\nabla F(\tilde{\mathbf{w}})\|_2^2 \leq \frac{\mathbb{E}\|\nabla F(\tilde{\mathbf{w}})\|_{\mathbf{A}_\sigma^{-1}}^2}{\lambda_{\min}(\mathbf{A}_\sigma^{-1})} \leq (1+4\sigma)\mathbb{E}\|\nabla F(\tilde{\mathbf{w}})\|_{\mathbf{A}_\sigma^{-1}}^2 \leq 4\zeta \frac{G\sqrt{6dL(2D_F + LG^2)\log(1/\delta)}}{n\epsilon}$$

where  $\zeta = \sqrt{\frac{1}{d} \sum_{i=1}^d \frac{(1+4\sigma)^2}{(1+2\sigma-2\sigma \cos(2\pi i/d))^2}} > 1$ . In the  $\ell_2$ -norm, DP-LSSGD has a bigger utility upper bound than DP-SGD (set  $\sigma = 0$  in  $\zeta$ ). However, this does not mean that DP-LSSGD has worse performance. To see this point, let us consider the following simple nonconvex function

$$f(x, y) = \begin{cases} \frac{x^2}{4} + y^2, & \text{for } \frac{x^2}{4} + y^2 \leq 1 \\ \sin\left(\frac{\pi}{2}\left(\frac{x^2}{4} + y^2\right)\right), & \text{for } \frac{x^2}{4} + y^2 > 1. \end{cases} \quad (10)$$

For two points  $\mathbf{a}_1 = (2, 0)$  and  $\mathbf{a}_2 = (1, \sqrt{3}/2)$ , the distance to the local minima  $\mathbf{a}^* = (0, 0)$  are 2 and  $\sqrt{7}/2$ , while  $\|\nabla f(\mathbf{a}_1)\|_2 = 1$  and  $\|\nabla f(\mathbf{a}_2)\|_2 = \sqrt{13}/2$ . So  $\mathbf{a}_2$  is closer to the local minima  $\mathbf{a}^*$  than  $\mathbf{a}_1$  while its gradient has a larger  $\ell_2$ -norm.

## Appendix B. Calculations of $\beta$ and $\gamma$

### B.1. Calculation of $\gamma$

To prove Proposition 5, we need the following two lemmas.

**Lemma 17 (Residue Theorem)** *Let  $f(z)$  be a complex function defined on  $\mathbb{C}$ , then the residue of  $f$  around the pole  $z = c$  can be computed by the formula*

$$\text{Res}(f, c) = \frac{1}{(n-1)!} \lim_{z \rightarrow c} \frac{d^{n-1}}{dz^{n-1}} ((z-c)^n f(z)). \quad (11)$$

where the order of the pole  $c$  is  $n$ . Moreover,

$$\oint f(z)dz = 2\pi i \sum_{c_i} \text{Res}(f, c_i), \quad (12)$$

where  $\{c_i\}$  be the set of pole(s) of  $f(z)$  inside  $\{z \mid |z| < 1\}$ .

The proof of Lemma 17 can be found in any complex analysis textbook.

**Lemma 18** For  $0 \leq \theta \leq 2\pi$ , suppose

$$F(\theta) = \frac{1}{1 + 2\sigma(1 - \cos(\theta))},$$

has the discrete-time Fourier transform of series  $f[k]$ . Then, for integer  $k$ ,

$$f[k] = \frac{\alpha^{|k|}}{\sqrt{4\sigma + 1}}$$

where

$$\alpha = \frac{2\sigma + 1 - \sqrt{4\sigma + 1}}{2\sigma}$$

**Proof** By definition,

$$f[k] = \frac{1}{2\pi} \int_0^{2\pi} F(\theta) e^{ik\theta} d\theta = \frac{1}{2\pi} \int_0^{2\pi} \frac{e^{ik\theta}}{1 + 2\sigma(1 - \cos(\theta))} d\theta. \quad (13)$$

We compute (13) by using Residue theorem. First, note that because  $F(\theta)$  is real valued,  $f[k] = f[-k]$ ; therefore, it suffices to compute (13) for nonnegative  $k$ . Set  $z = e^{i\theta}$ . Observe that  $\cos(\theta) = 0.5(z + 1/z)$  and  $dz = izd\theta$ . Substituting in (13) and simplifying yields that

$$f[k] = \frac{-1}{2\pi i \sigma} \oint \frac{z^k}{(z - \alpha_-)(z - \alpha_+)} dz, \quad (14)$$

where the integral is taken around the unit circle, and  $\alpha_{\pm} = \frac{2\sigma+1 \pm \sqrt{4\sigma+1}}{2\sigma}$  are the roots of quadratic  $-\sigma z^2 + (2\sigma + 1)z - \sigma$ . Note that  $\alpha_-$  lies within the unit circle; whereas,  $\alpha_+$  lies outside of the unit circle. Therefore, because  $k$  is nonnegative,  $\alpha_-$  is the only singularity of the integrand in (14) within the unit circle. A straightforward application of the Residue Theorem, i.e., Lemma 17, yields that

$$f[k] = \frac{-\alpha_-^k}{\sigma(\alpha_- - \alpha_+)} = \frac{\alpha^k}{\sqrt{4\sigma + 1}}.$$

This completes the proof. ■

**Proof** [Proof of Proposition 5] First observe that we can re-write  $\gamma$  as

$$\frac{1}{d} \sum_{j=0}^{d-1} \frac{1}{1 + 2\sigma(1 - \cos(\frac{2\pi j}{d}))}. \quad (15)$$

It remains to show that the above summation is equal to  $\frac{1+\alpha^d}{(1-\alpha^d)\sqrt{4\sigma+1}}$ . This follows by lemmas 18 and standard sampling results in Fourier analysis (i.e. sampling  $\theta$  at points  $\{2\pi j/d\}_{j=0}^{d-1}$ ). Nevertheless, we provide the details here for completeness: Observe that that the inverse discrete-time Fourier transform of

$$G(\theta) = \sum_{j=0}^{d-1} \delta\left(\theta - \frac{2\pi j}{d}\right).$$



is given by

$$g[k] = \begin{cases} d/2\pi & \text{if } k \text{ divides } d, \\ 0 & \text{otherwise.} \end{cases}$$

Furthermore, let

$$F(\theta) = \frac{1}{1 + 2\sigma(1 - \cos(\theta))},$$

and use  $f[k]$  to denote its inverse discrete-time Fourier transform. Now,

$$\begin{aligned} \frac{1}{d} \sum_{j=0}^{d-1} \frac{1}{1 + 2\sigma(1 - \cos(\frac{2\pi j}{d}))} &= \frac{1}{d} \int_0^{2\pi} F(\theta)G(\theta) \\ &= \frac{2\pi}{d} \text{DTFT}^{-1}[F \cdot G][0] \\ &= \frac{2\pi}{d} (\text{DTFT}^{-1}[F] * \text{DTFT}^{-1}[G])[0] \\ &= \frac{2\pi}{d} \sum_{r=-\infty}^{\infty} f[-r]g[r] \\ &= \frac{2\pi}{d} \sum_{\ell=-\infty}^{\infty} f[-\ell d] \frac{d}{2\pi} \\ &= \sum_{\ell=-\infty}^{\infty} f[-\ell d]. \end{aligned}$$

The proof is completed by substituting the result of lemma 18 in the above sum and simplifying. ■

We list some typical values of  $\gamma$  in Table 5.

Table 3: The values of  $\gamma$  corresponding to some  $\sigma$  and  $d$ .

$\sigma$	1	2	3	4	5
$d = 1000$	0.447	0.333	0.277	0.243	0.218
$d = 10000$	0.447	0.333	0.277	0.243	0.218
$d = 100000$	0.447	0.333	0.277	0.243	0.218

## B.2. Calculation of $\beta$

The proof of Proposition 8 is similar as the proof of Proposition 5. The only difference is that we need to compute

$$f[k] = \frac{1}{2\pi} \int_0^{2\pi} \frac{e^{ik\theta}}{(1 + 2\sigma(1 - \cos \theta))^2} d\theta. \quad (16)$$

By Residue theorem, for  $k > 0$  (note that  $f[-k] = f[k]$ ), we have

$$\begin{aligned}
 f[k] &= \frac{1}{2\pi} \int_0^{2\pi} \frac{e^{ik\theta}}{(1 + 2\sigma(1 - \cos \theta))^2} d\theta \\
 &= \frac{1}{2\pi i} \oint \frac{z^{k+1}}{(z + \sigma(2z - z^2 - 1))^2} dz \\
 &= \lim_{z \rightarrow \alpha^-} \frac{d}{dz} \left( \frac{z^{k+1}}{(z + \sigma(2z - z^2 - 1))^2} \right) \\
 &= \lim_{z \rightarrow \alpha^-} \frac{d}{dz} \left( \frac{z^{k+1}}{\sigma^2(z - \alpha^+)^2} \right) \\
 &= \frac{(k+1)\alpha^k}{4\sigma+1} + \frac{2\sigma\alpha^{k+1}}{(4\sigma+1)^{3/2}},
 \end{aligned}$$

where  $\alpha_- = \frac{2\sigma+1-\sqrt{4\sigma+1}}{2\sigma}$ . Therefore, we have

$$\beta = \frac{2\alpha^{2d+1} - \xi\alpha^{2d} + 2\xi d\alpha^d - 2\alpha + \xi}{\sigma^2\xi^3(1 - \alpha^d)^2}.$$

We list some typical values of  $\beta$  in Table 8.

Table 4: The values of  $\beta$  corresponding to some  $\sigma$  and  $d$ .

$\sigma$	1	2	3	4	5
$d = 1000$	0.268	0.185	0.149	0.128	0.114
$d = 10000$	0.268	0.185	0.149	0.128	0.114
$d = 100000$	0.268	0.185	0.149	0.128	0.114

### Appendix C. Laplacian Smoothing and Diffusion Equation

Let  $u(x, t)$  be a function defined on the space-time domain  $[0, 1] \times [0, +\infty)$ , suppose it satisfies the following diffusion equation with the Neumann boundary condition

$$\begin{cases} \frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2}, & (x, t) \in [0, 1] \times [0, +\infty), \\ \frac{\partial u(0, t)}{\partial x} = \frac{\partial u(1, t)}{\partial x} = 0, & t \in [0, +\infty) \\ u(x, 0) = f(x), & x \in [0, 1] \end{cases} \quad (17)$$

If we apply the backward Euler in time and central finite difference in space to discretize the governing equation in (17), we get

$$\mathbf{v}^{\Delta t} - \mathbf{v}^0 = \Delta t \mathbf{L} \mathbf{v}^{\Delta t},$$

where  $\mathbf{v}^0$  is the discretization of  $f(x)$ , and  $\mathbf{v}^{\Delta t}$  is the numerical solution of (17) at time  $\Delta t$ . Therefore, we have

$$\mathbf{v}^{\Delta t} = (I - \Delta t \mathbf{L})^{-1} \mathbf{v}^0,$$

which is the LS with  $\sigma = \Delta t$ .