

Differentially Private Neural Network Training under Hidden State Assumption

Ding Chen * Chen Liu*

July 12, 2024

Abstract

We present a novel approach called differentially private stochastic block coordinate descent (DP-SBCD) for training neural networks with provable guarantees of differential privacy under the hidden state assumption. Our methodology incorporates Lipschitz neural networks and decomposes the training process of the neural network into sub-problems, each corresponding to the training of a specific layer. By doing so, we extend the analysis of differential privacy under the hidden state assumption to encompass non-convex problems and algorithms employing proximal gradient descent. Furthermore, in contrast to existing methods, we adopt a novel approach by utilizing calibrated noise sampled from adaptive distributions, yielding improved empirical trade-offs between utility and privacy.

1 Introduction

Machine learning models, especially deep neural networks, have exhibited remarkable progress in the last decade across diverse fields. Their applications, such as face recognition [11] and large language models [20], have been integrated into people’s daily lives. However, the increasing demand for large amounts of training data in the training process has given rise to growing concerns regarding the potential privacy vulnerabilities [14, 28] associated with these models. For example, deep neural networks are shown to memorize the training data [34] so that we can even reconstruct part of the training data from the learned model parameters [19]. To address these issues, differential privacy [10] has become the gold standard for making formal and quantitative guarantees on model’s privacy and has been widely applied in learning problems [1, 18].

The predominate approach to ensuring differential privacy in the context of machine learning involves the incorporation of calibrated noise during each update step in the training phase, which leads to a trade-off between utility and privacy loss [2, 3]. Excessive noise often leads to a significant loss in utility, while insufficient noise may result in privacy leakage. Moreover, many privacy accountant methods, such as moment accountant [1], are typically based on the composition properties of differential privacy. Specifically, they usually assume that the calibrated noise in each training iteration follows the same distribution and that the algorithm’s internal states, i.e., the model’s parameters in each training step, can be revealed to the adversaries. As a result, the total privacy loss of the algorithm significantly increases with the number of training iterations, because more internal states will be revealed with more training iterations. Such assumptions lead to a very loose privacy estimation of the algorithm, because, in practice, most internal states are usually not even recorded during training.

The gap between theory and practice motivates researchers to consider more practical assumptions. Recently, the hidden state assumption was proposed to narrow down this gap [9, 13, 32]. This assumption posits that the internal states of the training phase are hidden, and that the adversaries only have access to the last iterate. Under this assumption, Chourasia et al. [9], Ye and Shokri [32] use Langevin diffusion to track the change rate of *Rényi Differential Privacy* (RDP) in each epoch and bound the privacy loss for *Differential Privacy Gradient Descent* (DP-GD) and *Differential Privacy Stochastic mini-batch Gradient Descent* (DP-SGD) [1]. Recently, Asoodeh and Diaz [4] also considered directly using hockey-stick divergence instead of Rényi divergence. These works claim a converged and small

*City University of Hong Kong, Hong Kong, China/ emails: ding.chen@my.cityu.edu.hk;chen.liu@cityu.edu.hk

privacy loss (usually less than 1) under some strict assumption such as strongly convex loss function and gradient Lipschitzness.

From a theoretical aspect, existing theorems indicate a valuable property that privacy loss could *converge* under the hidden state assumption in learning problems with strongly convex and β -smooth loss function. Otherwise, the privacy loss will *increase exponentially*. Consequently, existing methods cannot be directly applicable to training deep neural networks, whose loss functions are non-convex and may exhibit non-smoothness when using activation functions like ReLU [21].

To extend existing analyses of differential privacy under the hidden state assumptions to deep learning problem, we introduce deep neural networks with a Lipschitz constraint. To this end, we propose *differentially private stochastic block coordinate descent* (DP-SBCD). In our proposed approach, we formulate the task of training neural networks as a constrained optimization problem and subsequently construct the associated Lagrangian function. The DP-SBCD algorithm then considers each layer individually, decomposing the original learning problem into multiple sub-problems. By observing that the Lagrangian function exhibits convexity and smoothness concerning the model parameters of each layer, our theoretical analyses yield the conclusion that the DP-SBCD algorithm can yield models with guaranteed differential privacy under the hidden state assumption. Besides, our algorithm’s privacy loss is significantly smaller than that of composition theory.

In addition, this work investigates the connections between the calibrated noise and the smoothness of the loss function in the context of differential privacy. Intuitively speaking, we need noise of larger variance to achieve the same level of differential privacy if the Lipschitz constant of the loss gradient is bigger, but noise of large variance harms the utility when the algorithm approaches the optimal solution. Due to this, we study the impact of the variance of the calibrated noise on differential privacy and propose to adaptively adjust it during training. We also explain why the privacy loss will converge under the hidden state assumption. Empirically, our proposed algorithm with adaptive calibrated noise achieves a better trade-off between the model’s utility and privacy.

In summary, we highlight our contributions as follows:

- We propose *differentially private stochastic block coordinate descent* (DP-SBCD), which, to the best of our knowledge, is the first feasible method to solve non-convex problems with differential privacy guarantees under the hidden state assumption.
- The algorithms and theorems in this work possess a generic nature, rendering them compatible with proximal gradient descent and adaptive calibrated noise.
- Empirically, the adaptive calibrated noise proposed in this work is able to achieve better trade-offs between the model’s utility and privacy under the hidden state assumption.

Terminology and Notation The privacy loss refers to the *Rényi differential privacy* (RDP). When using the diffusion process to analyze the update scheme, we use t to represent the time in the process. We use $\|\cdot\|_2$ to represent l_2 norm for vectors and the spectral norm for matrices. In addition, $\|\cdot\|_F$ denotes Frobenius norm. Other notations are explicitly defined before usage.

2 Related Works

Differential Privacy Differential Privacy (DP) [10] was initially introduced to quantitatively ensure the privacy guarantees within the realm of databases. Nevertheless, its rigorous mathematical foundations and nice properties, such as composition theorem and post-processing immunity, have raised more and more interest in its application in machine learning. Among the practical methods for achieving DP-guaranteed models, DP-SGD [29] is the most popular one, which involves the addition of calibrated noise to the clipped per-sample gradients in each iterative update.

Besides the DP-guaranteed training algorithm, the proper accounting of privacy loss is also critical. Most accountant methods rely on the composition theorem of differential privacy. In this regard, Abadi et al. [1] introduces the momentum accountant, Mironov [25] further proposes Rényi differential privacy (RDP), which turns out a better tool to analyze the algorithm’s privacy loss. Specifically, a randomized mechanism $\mathcal{M} : \mathcal{D} \rightarrow \mathcal{R}$ is said to have (α, ϵ) -RDP if for any adjacent datasets $d, d' \in \mathcal{D}$, it holds that $D_\alpha(\mathcal{M}(d) || \mathcal{M}(d')) \leq \epsilon$ where D_α denotes the Rényi divergence of order α .

Recently, the hidden state assumption has been explored in many works [9, 13, 32]. While the composition theory assumes that adversaries have access to the intermediate states of the algorithm,

which requires the summation of the privacy loss across each training iteration [27], the hidden state assumption posits that these intermediate states remain concealed. In real-life applications, the intermediate states usually cannot be tracked, so the hidden state assumption is more practical. It is first proposed in [13], which also proves that a sequence of contractive maps can exhibit bounded (α, ε) -RDP. Subsequently, Chourasia et al. [9] consider the update scheme as a Langevin diffusion process and analyze the privacy loss for DP-GD using the log-Sobolev inequality [30]. Expanding upon this, Ye and Shokri [32] analyze DP-SGD by leveraging the convergent analysis of unadjusted Langevin algorithm proposed by Wibisono [31].

Block Coordinate Descent The block coordinate descent algorithm is designed to decompose the original problem into several sub-problems which allow for independent and efficient solutions. It not only helps prevent the vanishing gradient problem [35] but also facilitates distributed or parallel implementations [23]. In the era of deep learning, many works focus on employing block coordinate descent to train neural networks [16, 33, 35]. Zhang and Brand [35] utilize a lifting trick to solve the learning problem of deep neural networks with ReLU activation function, which is then extended by Gu et al. [16], Mangold et al. [24]. Furthermore, Zeng et al. [33] comprehensively analyzes two general types of optimization formulas: Two-splitting and Three-splitting. They also investigate the loss function, activation function, and architecture that these algorithms require.

Lipschitz Neural Networks The Lipschitz constant of a neural network holds significant importance, especially in the context of robustness [12, 22] and generalization [7, 26]. Recently, some works [8] remove the computationally expensive per-sample gradient clipping in DP-SGD by controlling the Lipschitz constants of neural networks. However, it is challenging to accurately enforce the Lipschitz constant of neural networks. Gouk et al. [15] address this challenge by using power iteration to approximate the layerwise Lipschitz constant.

3 Solving Neural Network using Block Coordinate Descent

In this section, we train Lipschitz neural networks with block coordinate descent and analyze the properties of the decomposed sub-problems.

3.1 Lipschitz Neural Network

We consider learning an D -layer neural network parameterized by $\boldsymbol{\theta} \stackrel{\text{def}}{=} \{\theta_d\}_{d=0}^D$:

$$\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) \stackrel{\text{def}}{=} \mathcal{R}(\theta_D, \mathbf{x}_D; y) \quad \text{s.t.} \quad \mathbf{x}_{d+1} = \sigma_d(\theta_d \mathbf{x}_d) \quad d = 0, \dots, D-1 \quad (1)$$

Here, $\mathbf{x} \stackrel{\text{def}}{=} \mathbf{x}_0$ is the input, and $\{\mathbf{x}_d\}_{d=1}^D$ are vectors representing the intermediate activations of each layer. \mathcal{R} is the loss function, such as softmax cross-entropy function, which is convex w.r.t. θ_D . Moreover, $\{\sigma_d\}_{d=0}^{D-1}$ are activation functions, such as ReLU function, that are Lipschitz continuous. Finally, $\{\theta_d\}_{d=0}^D$ are matrices representing the weights of linear layers, including fully-connected layers and convolutional layers. This formulation is generic and consistent with practices.

As formulated in (1), the neural network is the composition of D layers, so the global Lipschitz constant is the product of the Lipschitz constant of each layer [26]. We control the Lipschitz constant of each layer by normalizing the weight parameters. Specifically, following the practice of Behrmann et al. [6], we use power iteration [15] to approximate the spectral norm $\tilde{\lambda}_d$ of each layer and normalize the corresponding parameter θ_d given the designated bound of the Lipschitz constant ρ_d by:

$$\tilde{\theta}_d = \theta_d \cdot \min(\rho_d / \tilde{\lambda}_d, 1) \quad (2)$$

3.2 Block Coordinate Descent Algorithm

In this work, we use *Three-splitting* scheme [33] to solve problem (1). First, we introduce the variables $\{\mathbf{U}_d\}_{d=0}^{D-1}$ and rewrite the problem as follows:

$$\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{x}) := \mathcal{R}(\theta_D, \mathbf{x}_D; y) \quad \text{s.t.} \quad \mathbf{x}_{d+1} = \sigma_d(\mathbf{U}_d), \mathbf{U}_d = \theta_d \mathbf{x}_d \quad d = 0, \dots, D-1 \quad (3)$$

We then consider the Lagrangian function of the problem (3) with a multiplier co-efficient γ , which is set 1 in this work.

$$\mathcal{F}(\boldsymbol{\theta}, \mathbf{x}, \mathbf{U}) = \mathcal{R}(\theta_D, \mathbf{x}_D; y) + \frac{\gamma}{2} \sum_{d=0}^{D-1} (\|\mathbf{x}_{d+1} - \sigma_d(\mathbf{U}_d)\|_2^2 + \|\mathbf{U}_d - \theta_d \mathbf{x}_d\|_F^2) \quad (4)$$

\mathcal{F} is a function of $\{\theta_d\}_{d=0}^D$, $\{\mathbf{x}_d\}_{d=0}^D$ and $\{\mathbf{U}_d\}_{d=0}^{D-1}$. For notation simplicity, we only explicitly highlight the parameter we consider for \mathcal{F} if there is no ambiguity. For example $\mathcal{F}(\theta'_d)$ means we update the parameter θ_d in \mathcal{F} while keeping the other parameters fixed. We use block coordinate descent to update parameters of \mathcal{F} where we treat each layer as a sub-problem. The algorithm is summarized in Algorithm 1. Considering the loss function \mathcal{R} is a convex function w.r.t. both θ_D and \mathbf{x}_D , we have $\forall d$, $\mathcal{F}(\theta_d)$ is convex w.r.t. θ_d , and $\forall d$, $\mathcal{F}(\mathbf{x}_d)$ is convex w.r.t. \mathbf{x}_d .

To update θ_d , Algorithm 1 also considers the damping term $\frac{1}{2\eta}\|\theta_d - \theta'_d\|_F^2$ and the regularization term $r_d(\theta'_d)$, which can be optimized by the proximal gradient descent elaborated in the next section. To update \mathbf{x}_d , we have the analytical solution due to the convexity of $\mathcal{F}(\mathbf{x}_d)$.

$$\mathbf{x}_d \leftarrow \begin{cases} (\theta_d^T \theta_d + \mathbf{I})^{-1} (\theta_d^T \mathbf{U}_d + \sigma(\mathbf{U}_{d-1})) & \text{if } d = 0, 1, \dots, D-1. \\ \text{Prox}_{\frac{1}{\gamma}, \mathcal{R}}(\sigma_{D-1}(\mathbf{U}_{D-1})). & \text{if } d = D \end{cases} \quad (5)$$

To update \mathbf{U}_d , we utilize Zeng et al. [33, Lemma 13] as follows to obtain the analytical solution.

$$\mathbf{U}_d \leftarrow \begin{cases} \theta_d \mathbf{x}_d & \text{if } -\mathbf{x}_{d+1} \leq \theta_d \mathbf{x}_d \leq -(\sqrt{2}-1)\mathbf{x}_{d+1} \leq 0. \\ \min\{0, \theta_d \mathbf{x}_d\} & \text{if } \theta_d \mathbf{x}_d + \mathbf{x}_{d+1} \leq 0. \\ \frac{1}{2}(\theta_d \mathbf{x}_d + \mathbf{x}_{d+1}) & \text{otherwise.} \end{cases} \quad (6)$$

Algorithm 1 Stochastic Block Coordinate Descent For Lipschitz Neural Network

Input: training set, regularization schemes $\{r_d\}_{d=0}^D$, parameter η , batch size b , Lipschitz constant ρ .
Initialization: $\{\theta_d\}_{d=0}^D, \{\mathbf{x}_d\}_{d=0}^D, \{\mathbf{U}_d\}_{d=0}^{D-1}$.
for epoch $k = 0, 1, \dots, K-1$ **do**
 for each mini-batch of size b **do**
 for layer $d = 0, 1, \dots, D$ **do**
 Update $\theta_d, \mathbf{U}_d, \mathbf{x}_d$ simultaneously as follows:

 $\mathbf{x}_d \leftarrow \arg \min_{\mathbf{x}'_d} \mathcal{F}(\mathbf{x}'_d)$
 $\mathbf{U}_d \leftarrow \arg \min_{\mathbf{U}'_d} \mathcal{F}(\mathbf{U}'_d)$
 Normalize θ_d by (2) with coefficient ρ .
 $\theta_d \leftarrow \arg \min_{\theta'_d} \mathcal{F}(\theta'_d) + \frac{1}{2\eta}\|\theta_d - \theta'_d\|_F^2 + r_d(\theta'_d)$
 end for
 end for
end for

Algorithm 2 Differentially Private Stochastic Block Coordinate Descent For Lipschitz Neural Network

Input: training set, regularization schemes $\{r_d\}_{d=0}^D$, step size η , batch size b , noise control function $o(\eta, k, j)$, Lipschitz constant ρ .
Initialization: $\{\theta_d\}_{d=0}^D, \{\mathbf{x}_d\}_{d=0}^D, \{\mathbf{U}_d\}_{d=0}^{D-1}$.
for epoch $k = 0, 1, \dots, K-1$ **do**
 for each mini-batch of size b **do**
 for layer $d = 0, 1, \dots, D$ **do**
 Update $\theta_d, \mathbf{U}_d, \mathbf{x}_d$ simultaneously as follows:

 $\mathbf{x}_d \leftarrow \arg \min_{\mathbf{x}'_d} \mathcal{F}(\mathbf{x}'_d)$
 $\mathbf{U}_d \leftarrow \arg \min_{\mathbf{U}'_d} \mathcal{F}(\mathbf{U}'_d)$
 Normalize θ_d by (2) with coefficient ρ .
 Update θ_d based on (8)
 end for
 end for
end for

Now, we assume the activations $\{\mathbf{x}_d\}_{d=0}^D$ are bounded.

Assumption 3.1 (Bounded layer input) $\forall d, \exists X_d < +\infty$ such that \mathbf{x}_d always satisfies $\|\mathbf{x}_d\|_2^2 \leq X_d$.

Bounded input assumption is common in the machine learning community. Since we loosen $\{\mathbf{x}_d\}_{d=1}^D$ in Equation (4), we bound these intermediate activations as well. This assumption is benign, because \mathbf{x}_0 is anchored as the input, and $\{\mathbf{x}_d\}_{d=1}^D$ are softly constrained by several quadratic terms in Equation (4) that we aim to jointly minimize.

The following lemma confirms the smoothness of the function $\mathcal{F}(\theta_d)$ as a function of θ_d , the loss function for layerwise parameters.

Lemma 3.2 If Assumption 3.1 is satisfied, and the activation function is ReLU, then the function $\mathcal{F}(\theta_d)$ for any layer $0 \leq d \leq D$ is β_d -smooth and the smoothness constant β_d is γX_d^2 .

We have to highlight that the $\mathcal{F}(\theta_d)$ **assumes that other parameters except θ_d are fixed**. In addition, Lemma 3.2 indicates that the smoothness constant β_d does not depend on other parameters. Based on this, we can calculate the Hessian matrix of $\mathcal{F}(\theta_d)$ as $\nabla^2 \mathcal{F}(\theta_d) = \mathbf{x}_d \mathbf{x}_d^T \geq 0$. Therefore, $\mathcal{F}(\theta_d)$ is convex w.r.t. θ_d .

4 Differentially Private Stochastic Block Coordinate Descent

In this section, we propose *Differentially Private Stochastic Block Coordinate Descent* (DP-SBCD), i.e., the differentially private version of Algorithm 1. We then calculate the algorithm’s privacy loss under the hidden state assumption. We discuss the privacy loss in a generic form, especially the case when using adaptive calibrated noise.

4.1 Differentially Private Update Scheme

We first consider the update scheme for θ_d in Algorithm 1 and use first order Taylor expansion of $\mathcal{F}(\theta'_d)$ to write it in the proximal gradient descent format.

$$\theta_d \leftarrow \arg \min_{\theta'_d} \langle \nabla \mathcal{F}(\theta_d), \theta'_d - \theta_d \rangle + \frac{1}{2\eta} \|\theta'_d - \theta_d\|_F^2 + r_d(\theta'_d) = \text{Prox}_{\eta, r_d}(\theta_d - \eta \nabla \mathcal{F}(\theta_d)) \quad (7)$$

Owing to the post-processing immunity of differential privacy, the privacy guarantee is required for each sub-problem. Furthermore, as the algorithm converges gradually, the gradient $\nabla \mathcal{F}(\theta_d)$ diminishes, prompting the use of calibrated noise with an adaptive variance instead of a fixed one. As a result, we propose the following update scheme:

$$\theta_d \leftarrow \text{Prox}_{\eta, r_d}(\theta_d - \eta \nabla \mathcal{F}(\theta_d) + \mathcal{N}(0, 2\eta \cdot o(\eta, k, j) \mathbf{I})) \quad (8)$$

where $\mathcal{N}(0, \mathbf{I})$ is the standard Gaussian distribution, $o(\eta, k, j)$ is a function of the step size η , epoch index k and iteration index j to control the magnitude of the calibrated noise. By incorporating the update scheme (8) into Algorithm 1, we obtain Algorithm 2, which will be proved to guarantee differential privacy in the following section.

4.2 Privacy Loss of Sub-problems

We now study the privacy loss of Algorithm 2. Leveraging the post-processing immunity and composition property, we can establish an upper bound of the privacy loss of Algorithm 2 by summing the privacy losses of its individual sub-problems. Consequently, our primary focus here is to estimate the privacy loss associated with each sub-problem. For notation simplicity, we omit the subscript d in this subsection, as our analyses apply to any sub-problem.

In each iteration of Algorithm 2, the value of θ is updated by (8) and then normalized by (2). Among them, it is clear that the scaling operation of θ in (2) does not contribute to the privacy loss, so we focus on the update scheme (8).

We regard the update scheme (8) in Algorithm 2 as a diffusion process [5, 9, 32]. Specifically, the update scheme consists of three parts: The gradient descent part $\theta - \eta \nabla \mathcal{F}(\theta)$; The noise part $o(\eta, k, j)$; The proximal operator associated with a convex regularization function r . From a distributional perspective, let Θ be the distribution of the parameter θ , then the distribution of the parameter after one-step update in (8) can be represented as follows:

$$\tilde{\Theta} = T_{\#}(F_{\#}(\Theta) * \mathcal{N}(0, 2t \cdot o(\eta, k, j) \mathbb{I}_d)) \quad (9)$$

where $F_{\#}$ and $T_{\#}$ are two push-forward mappings. F represents the gradient descent update, T represents the proximal operator, and $*$ represents the convolution operator between two distributions.

Based on the smoothness property of $\mathcal{F}(\theta)$ indicated in Lemma 3.2, we prove the Lipschitz continuity of the first part of the update scheme (8).

Lemma 4.1 (Lipschitz continuity for F) *If $\mathcal{F}(\theta)$ is a β -smooth function and $\nabla^2 \mathcal{F}(\theta) \geq \omega$, then the update function $F(\theta) = \theta - \eta \nabla \mathcal{F}(\theta)$ is Lipschitz continuous with a constant $L_F \leq \max\{|1 - \eta\omega|, |1 - \eta\beta|\}$.*

We also need to highlight that the $\mathcal{F}(\theta)$ in the Lemma 4.1 **assumes that other parameters except θ are fixed**. In addition, the bound of the Lipschitz constant derived in Lemma 4.1 does not depend on other parameters. That is to say, the bound in Lemma 4.1 is valid for arbitrary values of other parameters. Since $\mathcal{F}(\theta)$ is proven convex in Section 3, $0 \leq \omega \leq \beta$. As a result, the Lipschitz constant will be dominated by the convexity term $|1 - \eta\omega|$ when the step size $\eta \leq \frac{2}{\omega + \beta}$ and otherwise the smoothness term $|1 - \eta\beta|$.

Similarly, we can prove the Lipschitz continuity of the third part of the update scheme (8).

Lemma 4.2 (Lipschitz continuity for T) *If $\eta > 0$ and $r(\theta)$ is a convex function, then the proximal operator function $T(\theta) = \text{Prox}_{\eta,r}(\theta) = \arg \min_{\tilde{\theta}} \left\{ r(\tilde{\theta}) + \frac{1}{2\eta} \|\tilde{\theta} - \theta\|_2^2 \right\}$ is Lipschitz continuous with a constant $L_T \leq 2$.*

Typical examples of the regularization function r include 1) no regularization: $r(\theta) = 0$, and then $\text{Prox}_{\eta,r}(\theta) = \theta$; 2) l_2 regularization in weight decay: $r(\theta) = \frac{1}{2} \|\theta\|_2^2$, and then $\text{Prox}_{\eta,r}(\theta) = \frac{\eta}{1+\eta} \theta$; 3) l_1 regularization in LASSO: $r(\theta) = \|\theta\|_1$, and then $\text{Prox}_{\eta,r}(\theta) = \text{sign}(\theta) \cdot \max(0, |\theta| - \eta)$. In all these three cases, the proximal function is Lipschitz continuous and the Lipschitz constant is 1.

Now we assume that θ for each subproblem satisfies the log-Sobolev inequality (LSI), which is a benign assumption used in Vempala and Wibisono [30], Ye and Shokri [32]. The formal definition of log-Sobolev inequality is provided in Definition 4.3. Based on the log-Sobolev inequality assumptions and Lipschitzness (Lemma 4.1,4.2), we consider the recursive privacy dynamics for Equation (9) and bound the change rate of RDP during one step of noisy mini-batch proximal gradient descent.

Definition 4.3 (Log-Sobolev Inequality [30]) *A distribution ν over \mathbb{R}^d satisfies log-Sobolev inequality (LSI) with a constant c if \forall smooth function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ with $\mathbb{E}_{\theta \sim \nu}[g^2(\theta)] < +\infty$,*

$$\mathbb{E}_{\theta \sim \nu}[g^2(\theta) \log(g^2(\theta))] - \mathbb{E}_{\theta \sim \nu}[g^2(\theta)] \log \mathbb{E}_{\theta \sim \nu}[g^2(\theta)] \leq \frac{2}{c} \mathbb{E}_{\theta \sim \nu}[\|\nabla g(\theta)\|_2^2] \quad (10)$$

It is an extension of Ye and Shokri [32, Lemma 3.2] to the cases of non-convex loss functions, proximal gradient descent and adaptive calibrated noise.

To derive the privacy loss of Algorithm 2, we assume a bounded sensitivity of the total gradient.

Assumption 4.4 (Sensitivity of the Total Gradient) *The l_2 sensitivity of the total gradient $\mathbb{E}_D \nabla \mathcal{F}(\theta)$ is finite. That is to say, $\exists S_g < +\infty$ such that for any dataset D and its neighbouring dataset D' that only differs in one instance, we have $S_g = \max_{D, D', \theta} \|\mathbb{E}_D \nabla \mathcal{F}(\theta) - \mathbb{E}_{D'} \nabla \mathcal{F}(\theta)\|_2$.*

Under Assumption 4.4, we obtain the privacy loss of the Algorithm 2 for each iteration.

Corollary 4.5 *Under Assumption 4.4 where the sensitivity of the total gradient is $S_g < +\infty$, let D, D' be an arbitrary pair of the neighboring datasets that only differ in the i_0 -th data point (i.e. $x_{i_0} \neq x'_{i_0}$). Let B_k^j be a fixed mini-batch used in the j -th iteration of the k -th epoch in Algorithm 2, which contains b different training instances whose indices are sampled from $\{0, 1, \dots, n-1\}$. We denote θ_k^j and $\theta'_k{}^j$ as the intermediate parameters in Algorithm 2 when using datasets D and D' , respectively. If the distributions of θ_k^j and $\theta'_k{}^j$ satisfy log-Sobolev inequality with a constant c , the update function $F(\theta) = \theta - \eta \nabla \mathcal{F}(\theta)$ and the proximal operator $T(\theta) = \text{Prox}_{\eta,r}(\theta)$ are Lipschitz continuous with Lipschitz constant L_F and L_T , respectively, then the following recursive bound for Rényi divergence holds for any order $\alpha > 1$:*

$$\begin{aligned} 1) \text{ If } i_0 \notin B_k^j, \text{ then } \frac{R_\alpha(\theta_k^{j+1} \|\theta'^{j+1}\|)}{\alpha} \text{ is upper bounded by } \frac{R_{\alpha'}(\theta_k^j \|\theta'^j\|)}{\alpha'} \cdot \left(1 + \frac{c \cdot 2\eta \cdot \sigma(\eta, k, j)}{L_F^2}\right)^{-1/L_T^2} \text{ where} \\ \text{the order } \alpha' = (\alpha - 1) \left(1 + \frac{c \cdot 2\eta \cdot \sigma(\eta, k, j)}{L_F^2}\right)^{-1} + 1. \\ 2) \text{ If } i_0 \in B_k^j, \text{ then } \frac{R_\alpha(\theta_k^{j+1} \|\theta'^{j+1}\|)}{\alpha} \text{ is upper bounded by } \frac{R_\alpha(\theta_k^j \|\theta'^j\|)}{\alpha} + \frac{\eta S_g^2}{4b^2 \cdot \sigma(\eta, k, j)}. \end{aligned}$$

Compared with the results discussed in Ye and Shokri [32] which only study the case without proximal operator, the privacy loss decay rate in the first case of Corollary 4.5 is powered by $-1/L_T^2$ instead of -1 , corresponding to the factor $1/L_T^2$ in the formulation of c_t in the change rate of RDP. This indicates that the Lipschitz constant L_T of the proximal operator also affects the privacy loss decay, a.k.a. privacy amplification when we run Algorithm 2 in the hidden state assumption. The

smaller L_T^2 is, the better privacy guarantee will be obtained. Corollary 4.5 concludes the change of the privacy loss for one iteration in Algorithm 2.

By applying Corollary 4.5 iteratively, we can obtain the algorithm's privacy loss for the whole training phase. The formal theorem is demonstrated below.

Theorem 4.6 *Under Assumption 4.4 where the sensitivity of the total gradient is $S_g < +\infty$, the distribution of θ satisfies log-Sobolev inequality with a constant c . In addition, the update function $F(\theta) = \theta - \eta \nabla \mathcal{F}(\theta)$ and the proximal operator $T(\theta) = \text{Prox}_{\eta, r}(\theta)$ are Lipschitz continuous with Lipschitz constant L_F and L_T , respectively. If we use Algorithm 2 to train model parameters θ for $K \geq 1$ epochs, then the algorithm satisfies $(\alpha, \varepsilon(\alpha))$ -Rényi differential privacy with the constant:*

$$\varepsilon_K(\alpha) \leq \frac{1}{\alpha - 1} \log \left(\sum_{j_0=0}^{n/b-1} \frac{b}{n} \cdot e^{(\alpha-1)(\varepsilon_K(\alpha, j_0))} \right)$$

where:

$$\varepsilon_K(\alpha, j_0) \leq \alpha \sum_{k=0}^{K-1} \frac{\eta S_g^2}{4b^2 \cdot o(\eta, k, j_0)} \cdot \left(\frac{c_k^{j_0+1}}{c_K^{n/b-1}} \left(\frac{1}{L_F^2 L_T^2} \right)^{(n/b-1)(K-k)-j_0} \prod_{l=k}^K \frac{c_l^{j_0+1}}{c_l^{j_0}} \right)^{-1/L_T^2} \quad (11)$$

In Inequality (11), c_k^j is the log-Sobolev inequality constant for the distribution of the model parameters in the j -th iteration of the k -th epoch. The value of c_k^j is calculated based on Lemma 4.7.

Proof Sketch. For each iteration of the Algorithm 2, the update scheme (8) fixed parameters except the θ . Hence, although the $\mathcal{F}(\theta)$ in the update scheme differs among each iteration because of different fixed parameters, it maintains Lipschitzness in each iteration. Hence, we could use Corollary 4.5 to analyze the algorithm's privacy loss. In each epoch, there is one and only one different mini-batch for two neighboring datasets, so we assume it is the j_0 -th batch w.l.o.g. and apply case 2) of Corollary 4.5 for this mini-batch update and case 1) of Corollary 4.5 for the other updates. For any α , the original privacy loss is 0 and the recursive bound in Corollary 4.5 holds. Therefore, we can uniformly bounded the value of $\frac{\varepsilon_K(\alpha, j_0)}{\alpha}$ for any α . Finally, since j_0 is uniformly distributed among the index set $\{0, 1, \dots, n/b - 1\}$, we use the joint convexity of scaled exponentiation of Rényi divergence to bound the final privacy loss $\varepsilon_K(\alpha)$. \square

Lemma 4.7 *(LSI constant sequence in Algorithm 2) For each layer's update scheme in Algorithm 2 with a batch size of b , if the update function $F(\theta) = \theta - \eta \nabla f(\theta)$ and the proximal operator $T(\theta) = \text{Prox}_{\eta, r}$ are Lipschitz continuous with Lipschitz constant L_F and L_T , respectively, then the distribution of parameter θ_k^j in the j -th iteration of the k -th epoch satisfies c_k^j log-Sobolev inequality (LSI) and the constant c_k^j is calculated by:*

$$c_k^j = \frac{1}{2\eta L_T^2} \left(\sum_{k'=0}^{k-1} \sum_{j'=0}^{n/b-1} o(\eta, k', j') (L_F L_T)^{2((k-k')(n/b)-j'+j-1)} + \sum_{j'=0}^{j-1} o(\eta, k, j') (L_F L_T)^{2(j-j'-1)} \right)^{-1}$$

$\varepsilon_K(\alpha, j_0)$ in the Theorem 4.6 represents the privacy loss when the only different instance of the two neighboring datasets is in the j_0 -th mini-batch of each epoch. The $\varepsilon_K(\alpha, j_0)$ shows that the overall privacy loss is the summation of each epoch's privacy loss term $\frac{\eta S_g^2}{4b^2 \cdot o(\eta, k, j)}$ times a decay rate term $\left(\frac{c_k^{j_0+1}}{c_K^{n/b-1}} \left(\frac{1}{L_F^2 L_T^2} \right)^{(n/b-1)(K-k)-j_0} \prod_{l=k}^K \frac{c_l^{j_0+1}}{c_l^{j_0}} \right)^{-1/L_T^2}$. More importantly, Algorithm 2 maintains the decay rate term smaller than 1 and decreases with the increase in K .

Theorem 4.6 improves the results in existing works from many aspects for estimating the differential privacy under hidden state assumptions. Firstly, it is more generally applicable to different algorithms. Theorem 4.6 could easily extend the analyses and privacy guarantees from the gradient descent algorithm with calibrated noise from a fixed distribution to proximal gradient descent with adaptive calibrated noise. Second, it provides a feasible privacy loss accountant for non-convex problems with adaptive noise. In contrast to the assumption of strong convexity and β -smoothness on

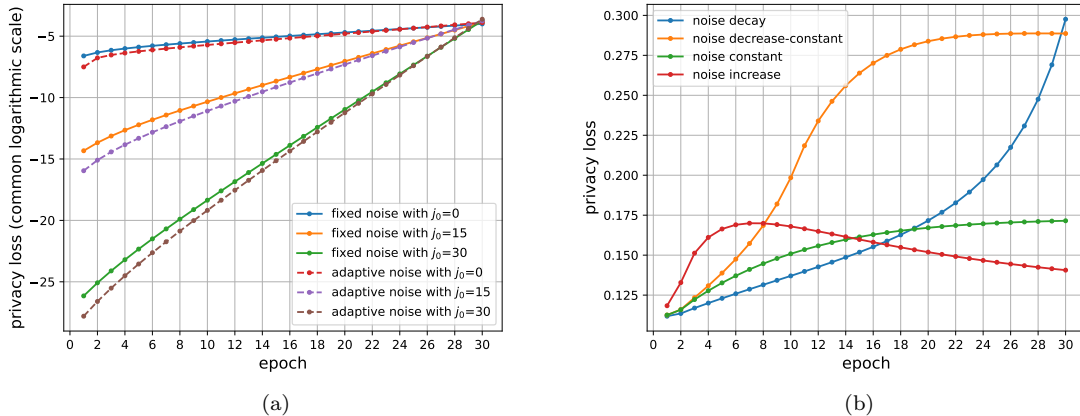


Figure 1: Common hyper-parameter settings: Lipschitz constant $\rho = 0.99$, batch size $b = 100$, stepsize $\eta = 0.01$, data set $n = 2100$. (a) Each epoch’s privacy loss contribution for the Algorithm 2 for $K = 30$ epochs. We study the scenarios of both fixed and adaptive noise. Each scenario contains three situations ($j_0 = 0, 15, 30$). In the fixed noise scenario, the $o(\eta, k, j) = 0.01$, while in the adaptive noise scenario, the $o(\eta, k, j) = 0.01 - 0.003k$. (b) Privacy loss for different numbers of epochs K for Algorithm 2. We show four scenarios: noise decay means $o(\eta, k, j) = 0.001 - 0.0003k$, noise constant means $o(\eta, k, j) = 0.0005$, noise increase means $o(\eta, k, j) = 0.0001 + 0.0003k$, and noise decrease-constant means that $o(\eta, k, j) = 0.0005 - 0.00003k$ where $k \in [1, 10]$ and remain $o(\eta, k, j) = 0.0002$ where $k \in [10, 30]$.

the objective function in existing results, Theorem 4.6 shows that Algorithm 2 applicable to general neural networks with Lipschitz constraints. Finally, even when downgrading to the case of gradient descent with calibrated noise sampled from fixed distributions, Theorem 4.6 has a tighter bound than previous work [32, Theorem 3.3]. This is because we directly bound $\varepsilon_K(\alpha, j_0)$ by recursively applying Corollary 4.5. By contrast, Ye and Shokri [32] approximate the bound $\varepsilon_K(\alpha, j_0)$ by part of iterations in one epoch rather than the whole epoch.

4.3 Privacy Loss variation under the Adaptive Noise

As discussed above, the bound of $\varepsilon_K(\alpha, j_0)$ in Theorem 4.6 elucidates how each epoch influences the cumulative privacy loss, exhibiting a decay rate. Specifically, these contributions are inversely proportional to $o(\eta, j, k)$, i.e., the variance of the noise, and decrease exponentially¹ as the number of iterations and epochs increase. That is to say, under the hidden state assumption, the calibrated noise in the last few epochs primarily contributes to the total privacy loss.

We illustrate this phenomenon in Figure 1(a), which encompasses scenarios with both constant and non-constant $o(\eta, j, k)$ values. The numerical results align with the analysis, as indicated by the near-linear curves observed in the logarithmic scale graph. Furthermore, we observe that a larger value of j_0 , signifying a delayed occurrence of the mini-batch containing the unique instance, leads to a smaller privacy loss during the initial stages of training. However, as the training progresses, this disparity diminishes. This observation is consistent with Theorem 4.6, as the privacy loss is zero for the first j_0 mini-batches of the first epoch.

Theorem 4.6 also enhances our understanding of the *convergence of the privacy loss*, initially proposed in Feldman et al. [13]: when $o(\eta, k, j)$ is a constant, an equilibrium solution emerges between the privacy loss and the decay rate after several epochs, resulting in the convergence of privacy loss under the hidden state assumption. In the case of adaptive calibrated noise, Figure 1(b) shows the tendency of privacy loss under four distinct noise settings. Empirical findings suggest that privacy loss converges when a constant magnitude of noise is utilized in the later stages of training. Conversely, the privacy loss diverges if the noise magnitude continues to decrease in the late phase of training. In contrast, increasing the noise magnitude can even lead to a decrease in privacy loss. However, it is

¹While not strictly exponential due to variations in the log-Sobolev inequality constant, the behavior closely resembles that of an exponential function curve in simulations.

worth noting that noise with a smaller variance tends to yield better utility for the model compared to noise with a larger variance, which is consistent with the utility-privacy trade-off [2, 3].

5 Numerical Experiments

We run numerical simulations in this section to investigate the model’s utility and privacy loss in different training phases when we use different distributions to sample calibrated noise.

We mainly utilize the Madelon dataset in this section. The Madelon dataset was originally introduced as a challenging classification problem in the NIPS 2003 feature selection challenge [17]. This synthetic dataset comprises 6000 instances, each with 20 features and belonging to one of the five classes. To ensure an unbiased evaluation, we divided the dataset into a training set, accounting for 80% of the data, and a testing set, containing the remaining 20%. The batch size is 960.

We employ a multilayer perceptron (MLP) model with four layers, each of which has 200 neurons. The activation function is the commonly used ReLU function, and we use squared loss as the loss function. In addition, we set the layerwise Lipschitz constant $\rho = 3$ and the step size $\eta = 0.01$. To guarantee privacy, we applied Theorem 4.6 and employed the privacy loss calculation with $\alpha = 100$. All the experiments can be efficiently executed on a single NVIDIA RTX 5000 Ada GPU.

The experiment compares the algorithm’s utility under different noise strategies: the constant strategy with $o(\eta, k, j) = 0.01$ and the decrease strategy with a linear decay rate of 0.0008 per epoch and final noise variance $o(\eta, K, j) = 0.0075$. The decrease strategy is designed in a manner so that the privacy loss of both strategies will be approximately the same for the same number of epochs K . Based on the aforementioned setup, we train the model 50 times where the total number of epochs K varies from 10 to 50. We run the whole experiment 5 times and report both the average performance and the standard deviation. The results are demonstrated in Figure 2 and Table 1.

Table 1: The accuracy(in % under 95% confidential interval) and the privacy loss of Algorithm 2 when we use different noise strategies and train the model for different numbers of epochs. D means noise decrease scenario and C means noise constant scenario

EPOCH	NOISE	PRIVACY LOSS	AVG ACC.
10	D	0.040556	20.13(± 0.32)
20	D	0.040648	42.15(± 9.35)
30	D	0.040688	92.10(± 2.20)
40	D	0.040704	93.98(± 1.38)
50	D	0.040708	94.37(± 1.46)
10	C	0.041408	20.03(± 0.59)
20	C	0.040652	23.57(± 6.89)
30	C	0.040696	53.35(± 17.68)
40	C	0.040712	82.35(± 7.98)
50	C	0.040716	91.16(± 3.20)

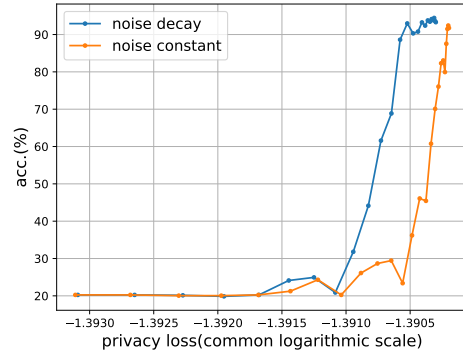


Figure 2: Relationship between model’s utility and privacy loss under different noise strategies. For each noise strategy, we run Algorithm 2 for different numbers of epochs to plot the curve.

The experiment results validate the effectiveness of Algorithm 2. What’s more, with proper settings of adaptive calibrated noise, the algorithm can demonstrate a better trade-off between the model’s utility and privacy. For the examples in Table 1, we see higher utilities and lower privacy loss when we use adaptive calibrated noise. In Figure 2, we see the curve of adaptive calibrated noise above one of its counterparts in most cases.

6 Conclusion

We propose *differentially private stochastic block coordinate descent* (DP-SBCD) algorithm, which includes proximal gradient descent and adaptive noise, to train neural networks. As far as we are aware, DP-SBCD is the first algorithm capable of addressing non-convex training problems while ensuring differential privacy guarantees under the hidden state assumption. The adaptive noise proposed in

our method can provide adjustable trade-offs between the model’s utility and privacy. Moreover, our numerical experiments show that under proper settings, our method could provide a better trade-off between utility and privacy. Going forward, our future research will concentrate on further refining the convergence and performance of the algorithm and extending its applicability to more generic settings.

References

- [1] 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*, pages 308–318, 2016.
- [2] Y. Allouah, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. On the privacy-robustness-utility trilemma in distributed learning. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 569–626. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/allouah23a.html>.
- [3] M. S. Alvim, M. E. Andrés, K. Chatzikokolakis, P. Degano, and C. Palamidessi. Differential privacy: on the trade-off between utility and information leakage. In *Formal Aspects of Security and Trust: 8th International Workshop, FAST 2011, Leuven, Belgium, September 12-14, 2011. Revised Selected Papers 8*, pages 39–54. Springer, 2012.
- [4] S. Asoodeh and M. Diaz. Privacy loss of noisy stochastic gradient descent might converge even for non-convex losses. *arXiv preprint arXiv:2305.09903*, 2023.
- [5] B. Balle, G. Barthe, M. Gaboardi, and J. Geumlek. Privacy amplification by mixing and diffusion mechanisms. *Advances in neural information processing systems*, 32, 2019.
- [6] J. Behrmann, W. Grathwohl, R. T. Chen, D. Duvenaud, and J.-H. Jacobsen. Invertible residual networks. In *International conference on machine learning*, pages 573–582. PMLR, 2019.
- [7] L. Béthune, T. Boissin, M. Serrurier, F. Mamalet, C. Friedrich, and A. Gonzalez Sanz. Pay attention to your loss: understanding misconceptions about lipschitz neural networks. *Advances in Neural Information Processing Systems*, 35:20077–20091, 2022.
- [8] L. Béthune, T. Masséna, T. Boissin, Y. Prudent, C. Friedrich, F. Mamalet, A. Bellet, M. Serrurier, and D. Vigouroux. Dp-sgd without clipping: The lipschitz neural network way. *arXiv preprint arXiv:2305.16202*, 2023.
- [9] R. Chourasia, J. Ye, and R. Shokri. Differential privacy dynamics of langevin diffusion and noisy gradient descent. *Advances in Neural Information Processing Systems*, 34:14771–14781, 2021.
- [10] C. Dwork. Differential privacy. In *International colloquium on automata, languages, and programming*, pages 1–12. Springer, 2006.
- [11] Z. Erkin, M. Franz, J. Guajardo, S. Katzenbeisser, I. Lagendijk, and T. Toft. Privacy-preserving face recognition. In *Privacy Enhancing Technologies: 9th International Symposium, PETS 2009, Seattle, WA, USA, August 5-7, 2009. Proceedings 9*, pages 235–253. Springer, 2009.
- [12] M. Fazlyab, A. Robey, H. Hassani, M. Morari, and G. Pappas. Efficient and accurate estimation of lipschitz constants for deep neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [13] V. Feldman, I. Mironov, K. Talwar, and A. Thakurta. Privacy amplification by iteration. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 521–532. IEEE, 2018.
- [14] 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*, pages 1322–1333, 2015.

- [15] H. Gouk, E. Frank, B. Pfahringer, and M. J. Cree. Regularisation of neural networks by enforcing lipschitz continuity. *Machine Learning*, 110:393–416, 2021.
- [16] F. Gu, A. Askari, and L. El Ghaoui. Fenchel lifted networks: A lagrange relaxation of neural network training. In *International Conference on Artificial Intelligence and Statistics*, pages 3362–3371. PMLR, 2020.
- [17] I. Guyon, S. Gunn, A. Ben-Hur, and G. Dror. Result analysis of the nips 2003 feature selection challenge. *Advances in neural information processing systems*, 17, 2004.
- [18] T. Ha, T. K. Dang, T. T. Dang, T. A. Truong, and M. T. Nguyen. Differential privacy in deep learning: an overview. In *2019 International Conference on Advanced Computing and Applications (ACOMP)*, pages 97–102. IEEE, 2019.
- [19] N. Haim, G. Vardi, G. Yehudai, O. Shamir, and M. Irani. Reconstructing training data from trained neural networks. *Advances in Neural Information Processing Systems*, 35:22911–22924, 2022.
- [20] N. Kandpal, E. Wallace, and C. Raffel. Deduplicating training data mitigates privacy risks in language models. In *International Conference on Machine Learning*, pages 10697–10707. PMLR, 2022.
- [21] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [22] Q. Li, S. Haque, C. Anil, J. Lucas, R. B. Grosse, and J.-H. Jacobsen. Preventing gradient attenuation in lipschitz constrained convolutional networks. *Advances in neural information processing systems*, 32, 2019.
- [23] D. Mahajan, S. S. Keerthi, and S. Sundararajan. A distributed block coordinate descent method for training l_1 regularized linear classifiers. *The Journal of Machine Learning Research*, 18(1): 3167–3201, 2017.
- [24] P. Mangold, A. Bellet, J. Salmon, and M. Tommasi. Differentially private coordinate descent for composite empirical risk minimization. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 14948–14978. PMLR, 17–23 Jul 2022.
- [25] I. Mironov. Rényi differential privacy. In *2017 IEEE 30th computer security foundations symposium (CSF)*, pages 263–275. IEEE, 2017.
- [26] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro. Exploring generalization in deep learning. *Advances in neural information processing systems*, 30, 2017.
- [27] N. Ponomareva, H. Hazimeh, A. Kurakin, Z. Xu, C. Denison, H. B. McMahan, S. Vassilvitskii, S. Chien, and A. G. Thakurta. How to dp-fy ml: A practical guide to machine learning with differential privacy. *Journal of Artificial Intelligence Research*, 77:1113–1201, 2023.
- [28] 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)*, pages 3–18. IEEE, 2017.
- [29] S. Song, K. Chaudhuri, and A. D. Sarwate. Stochastic gradient descent with differentially private updates. In *2013 IEEE global conference on signal and information processing*, pages 245–248. IEEE, 2013.
- [30] S. Vempala and A. Wibisono. Rapid convergence of the unadjusted langevin algorithm: Isoperimetry suffices. *Advances in neural information processing systems*, 32, 2019.
- [31] A. Wibisono. Proximal langevin algorithm: Rapid convergence under isoperimetry. *arXiv preprint arXiv:1911.01469*, 2019.

- [32] J. Ye and R. Shokri. Differentially private learning needs hidden state (or much faster convergence). *Advances in Neural Information Processing Systems*, 35:703–715, 2022.
- [33] J. Zeng, T. T.-K. Lau, S. Lin, and Y. Yao. Global convergence of block coordinate descent in deep learning. In *International conference on machine learning*, pages 7313–7323. PMLR, 2019.
- [34] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=Sy8gdB9xx>.
- [35] Z. Zhang and M. Brand. Convergent block coordinate descent for training tikhonov regularized deep neural networks. *Advances in Neural Information Processing Systems*, 30, 2017.