

Differentially Private ADMM for Convex Distributed Learning: Improve Accuracy via Multi-Step Approximation

Zonghao Huang

*School of Electrical and Computer Engineering
Oklahoma State University
Stillwater, OK 74075, USA*

ZONGHAO.HUANG@OKSTATE.EDU

Editor: Zonghao Huang

Abstract

Alternating Direction Method of Multipliers (ADMM) is a popular algorithm for distributed learning, where a network of agents collaboratively solve a regularized empirical risk minimization by iterative local computation and iterate exchanges. When the training data is sensitive, the exchanged iterates will cause serious privacy concern. In this paper, we propose a new differentially private distributed ADMM algorithm for a wide range of convex learning problems. In our proposed algorithm, we adopt the approximation of the augmented Lagrangian function to introduce calibrated noise into ADMM robustly, and allow l approximate primal variable updates per node in each iteration to improve the utility. We demonstrate the privacy guarantee of our approach and analyze the utility by the excess empirical risk with feasibility violation. Our theoretical results demonstrate that with limited communications our approach can obtain higher utility if we allow more primal variable updates per node in each iteration, and with sufficiently large iteration number our approach can achieve the error bounds, which are comparable to the state-of-art error bounds for differentially private empirical risk minimization.

1. Introduction

The advances in machine learning are due to the abundance of data, which could be collected over network but cannot be handled by a single processor. This motivates distributed learning, where data is distributed and possessed by multiple agents or nodes. In distributed learning frameworks, a network of agents collaboratively solve an optimization problem which is usually formulated as a regularized empirical risk minimization associated with the distributed data. Distributed learning has been widely applied in a variety of areas such as vehicle networks (Han et al., 2017) and biomedical sensing (Gong et al., 2016).

There exist approaches for distributed optimization including distributed subgradient descent algorithms (Nedic et al., 2008; Nedic and Ozdaglar, 2009), dual averaging methods (Duchi et al., 2011; Tsianos et al., 2012), and Alternating Direction Method of Multipliers (ADMM) (Boyd et al., 2011; Ling and Ribeiro, 2014; Shi et al., 2014; Zhang and Kwok, 2014). Among these algorithms, ADMM demonstrates fast convergence in many applications by both numerical and theoretical results. Prior works (Shi et al., 2014; Makhdoomi and Ozdaglar, 2017) have proved that in distributed ADMM, the iterates can converge linearly to the optimal solution while the objective value with feasibility violation can converge

to the optimum at a rate of $O(1/t)$, where t is the number of iterations. In this paper, we mainly focus on ADMM-based distributed learning.

In ADMM-based distributed learning, agents solve the regularized empirical risk minimization by iterative local computation and information exchanges. Local computation requires each agent to solve a local minimization associated with his local dataset while the information exchanges need agents to share the updated iterates with their neighbours. In this learning process, the exchanged information would cause serious privacy concern if the training data is sensitive. Therefore, additional methods are required to control privacy leakage. In this paper, we use differential privacy (Dwork et al., 2014, 2006b,a), and in ADMM-based distributed learning, differential privacy can be guaranteed by adding calibrated noise into iterates. Recently, there are a few of works (Zhang et al., 2018b; Zhang and Zhu, 2017; Zhang et al., 2018a; Huang et al., 2019) focusing on designing differentially private ADMM-based distributed learning and demonstrating their privacy-utility trade-off by numerical results. Zhang and Zhu (Zhang and Zhu, 2017) propose primal variable perturbation and dual variable perturbation to achieve dynamic differential privacy in ADMM. Zhang et al. (Zhang et al., 2018a) consider adaptive penalty parameters and propose to perturb the penalty. Huang et al. (Huang et al., 2019) propose DP-ADMM by adopting first-order approximation with time-varying Gaussian noise addition, and theoretically demonstrate that their approach can converge at a rate of $O(1/\sqrt{t})$, where t is the number of iterations. Zhang et al. (Zhang et al., 2018b) propose Recycled ADMM where the information from odd iterations can be reutilized in even iterations to save privacy budget. Hu et al. (Hu et al., 2019) consider distributed data feature and propose an ADMM-based algorithm with differential privacy guarantee.

In this paper, we propose a new differentially private distributed ADMM algorithm. The key algorithmic features of our approach are to adopt the approximation in local computation, and to allow l primal variable updates per node in each iteration. We adopt the approximate function in order to combine the calibrated noise ensuring differential privacy and the ADMM learning process robustly. This idea is similar to that one used in (Huang et al., 2019). In addition, our approach allows l primal variable updates per node in each iteration to provide a better privacy-utility trade-off. We demonstrate that by adding calibrated noise our approach can achieve differential privacy, and we analyze the utility of our proposed algorithm by the excess empirical risk with feasibility violation. Our theoretical results demonstrate that with limited communications our approach can have higher utility if we perform more primal variable updates per node in each iteration, and with sufficiently large iteration number our approach can achieve the error bounds, which are comparable to the state-of-art error bounds for differentially private empirical risk minimization.

The main contributions of this paper are summarized as follows:

1. We propose a new differentially private ADMM-based distributed learning algorithm, where the approximation is adopted and l primal variable updates are performed per node in each ADMM iteration to improve the privacy-utility trade-off.
2. We analyze the privacy guarantee of our proposed algorithm by using the moments accountant method. By properly setting the noise magnitude, we demonstrate that our approach can achieve (ϵ, δ) -differential privacy.

3. We analyze the utility of our approach theoretically by the excess empirical risk with feasibility violation by assuming the objective function is convex and Lipschitz. Our theoretical results show that with limited communications our approach can have higher utility if we perform more primal variable updates per node in each iteration.
4. We conduct numerical experiments based on real-world datasets to show the improved privacy-utility trade-off of our approach and demonstrate our theoretical results.

2. Problem Statement

In this section, we first introduce our problem setting. Then, we describe the ADMM-based distributed learning algorithm, and discuss the associated privacy concern.

2.1 Problem Setting

We consider a connected network given by a undirected graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, which consists of a set of nodes or agents $\mathcal{V} = \{1, 2, \dots, n\}$, and a set of edges \mathcal{E} . In this connected network, each agent can only exchange information with his connected neighbours, and we use \mathcal{N}_i to denote the neighbour set of agent i . Each agent $i \in \mathcal{V}$ possesses a private training dataset $\{(\mathbf{a}_{i,j}, b_{i,j}), j \in \mathcal{D}_i\}$, where $\mathbf{a}_{i,j} \in \mathbb{R}^d$ represents the d -dimensional data feature vector of the j -th training sample, and $b_{i,j} \in \{+1, -1\}$ is the corresponding data label.

The goal of our problem is to train a supervised learning model on the aggregated dataset $\{\mathcal{D}_i\}_{i \in \mathcal{V}}$, which enables predicting a label for any new data feature vector. The learning objective can be formulated as the following regularized empirical risk minimization problem:

$$\min_{\mathbf{w}} \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{D}_i} \frac{1}{|\mathcal{D}_i|} \ell(\mathbf{a}_{i,j}, b_{i,j}, \mathbf{w}) + \lambda \phi(\mathbf{w}), \quad (1)$$

where $\mathbf{w} \in \mathcal{W} \subseteq \mathbb{R}^d$ is the trained machine learning model, $\ell(\cdot)$ is the loss function used to measure the quality of the trained model, e.g., the loss function $\ell(\mathbf{a}, b, \mathbf{w})$ can be defined by $\log(1 + \exp(-b\mathbf{w}^\top \mathbf{a}))$ when we consider logistic regression, $\phi(\cdot)$ refers to the regularizer function introduced to prevent overfitting, and $\lambda > 0$ is the regularizer parameter controlling the impact of regularizer.

In this paper, we assume that the loss function $\ell(\cdot)$ and the regularizer function $\phi(\cdot)$ are both convex and Lipschitz. We use $\nabla \ell(\cdot)$ and $\nabla \phi(\cdot)$ to denote their gradient if they are differentiable or subgradient if not differentiable. We use $\|\cdot\|$ to denote the Euclidean norm.

2.2 ADMM-Based Distributed Learning Algorithm

To solve problem (1) with ADMM in distributed manner, we need to reformulate it as:

$$\min_{\{\mathbf{w}_i\}} \sum_{i \in \mathcal{V}} \left(\sum_{j \in \mathcal{D}_i} \frac{1}{|\mathcal{D}_i|} \ell(\mathbf{a}_{i,j}, b_{i,j}, \mathbf{w}_i) + \frac{\lambda}{n} \phi(\mathbf{w}_i) \right), \quad (2a)$$

$$\text{s.t. } \mathbf{w}_i = \mathbf{z}_{i,j}, \mathbf{z}_{i,j} = \mathbf{w}_j, \forall i \in \mathcal{V}, \forall j \in \mathcal{N}_i \quad (2b)$$

where $\mathbf{w}_i \in \mathcal{W} \subseteq \mathbb{R}^d$ is the local model solved by agent i , and $\mathbf{z}_{i,j} \in \mathbb{R}^d$ is an auxiliary variable imposing the consensus constraint on neighboring agents. The objective function (2a) is decoupled and constraints (2b) enforce that all the local models reach consensus finally.

Let $\{\mathbf{w}_i\}$, $\{\mathbf{z}_{i,j}\}$, and $\{\gamma_{i,j}^a\}$ be the shorthand for $\{\mathbf{w}_i\}_{i \in \mathcal{V}}$, $\{\mathbf{z}_{i,j}\}_{i \in \mathcal{V}, j \in \mathcal{N}_i}$, and $\{\gamma_{i,j}^a\}_{i \in \mathcal{V}, j \in \mathcal{N}_i, a=1,2}$, respectively. The augmented Lagrangian function associated with the problem (2) is:

$$\begin{aligned} \mathcal{L}(\{\mathbf{w}_i\}, \{\mathbf{z}_{i,j}\}, \{\gamma_{i,j}^a\}) = & \sum_{i \in \mathcal{V}} \left(\sum_{j \in \mathcal{D}_i} \frac{1}{|\mathcal{D}_i|} \ell(\mathbf{a}_{i,j}, b_{i,j}, \mathbf{w}_i) + \frac{\lambda}{n} \phi(\mathbf{w}_i) \right. \\ & - \sum_{j \in \mathcal{N}_i} \langle \gamma_{i,j}^1, \mathbf{w}_i - \mathbf{z}_{i,j} \rangle - \sum_{j \in \mathcal{N}_i} \langle \gamma_{i,j}^2, \mathbf{z}_{i,j} - \mathbf{w}_j \rangle \\ & \left. + \frac{\rho}{2} \sum_{j \in \mathcal{N}_i} (\|\mathbf{w}_i - \mathbf{z}_{i,j}\|^2 + \|\mathbf{w}_j - \mathbf{z}_{i,j}\|^2) \right), \end{aligned} \quad (3)$$

where $\{\gamma_{i,j}^a\}$ are the dual variables associated with constraints (2b) and $\rho > 0$ is the penalty parameter. The ADMM solves the problem (2) in a Gauss-Seidel manner by minimizing (3) w.r.t. $\{\mathbf{w}_i\}$ and $\{\mathbf{z}_{i,j}\}$ alternatively followed by dual updates of $\{\gamma_{i,j}^a\}$:

$$\mathbf{w}_i^{k+1} = \underset{\mathbf{w}_i}{\operatorname{argmin}} \mathcal{L}(\{\mathbf{w}_i\}, \{\mathbf{z}_{i,j}^{k-1}\}, \{\gamma_{i,j}^{a,k}\}), \quad (4a)$$

$$\mathbf{z}_{i,j}^{k+1} = \underset{\mathbf{z}_{i,j}}{\operatorname{argmin}} \mathcal{L}(\{\mathbf{w}_i^{k+1}\}, \{\mathbf{z}_{i,j}\}, \{\gamma_{i,j}^{a,k}\}), \quad (4b)$$

$$\gamma_{i,j}^{1,k+1} = \gamma_{i,j}^{1,k} - \rho(\mathbf{w}_i^{k+1} - \mathbf{z}_{i,j}^{k+1}), \quad (4c)$$

$$\gamma_{i,j}^{2,k+1} = \gamma_{i,j}^{2,k} - \rho(\mathbf{z}_{i,j}^{k+1} - \mathbf{w}_j^{k+1}). \quad (4d)$$

According to the previous works (Forero et al., 2010), the above iterate updates could be simplified by initializing $\gamma_{i,j}^{1,0} = \gamma_{i,j}^{2,0} = \mathbf{0}$, which can enforce $\gamma_{i,j}^{1,k} = \gamma_{i,j}^{2,k}$ and $\mathbf{z}_{i,j}^k = \frac{1}{2}(\mathbf{w}_i^k + \mathbf{w}_j^k)$. Let $\gamma_i^k = \sum_{j \in \mathcal{N}_i} \gamma_{i,j}^{1,k} = \sum_{j \in \mathcal{N}_i} \gamma_{i,j}^{2,k}$, let $\{\mathbf{w}_j^k\}$ be the shorthand of $\{\mathbf{w}_j^k\}_{j \in \mathcal{N}_i}$, and define $\mathcal{L}_i^s(\mathbf{w}_i, \mathbf{w}_i^k, \{\mathbf{w}_j^k\}, \gamma_i^k)$ as:

$$\begin{aligned} \mathcal{L}_i^s(\mathbf{w}_i, \mathbf{w}_i^k, \{\mathbf{w}_j^k\}, \gamma_i^k) = & \sum_{j \in \mathcal{D}_i} \frac{1}{|\mathcal{D}_i|} \ell(\mathbf{a}_{i,j}, b_{i,j}, \mathbf{w}_i) + \frac{\lambda}{n} \phi(\mathbf{w}_i) \\ & - 2 \langle \gamma_i^k, \mathbf{w}_i \rangle + \rho \sum_{j \in \mathcal{N}_i} \left\| \mathbf{w}_i - \frac{1}{2}(\mathbf{w}_i^k + \mathbf{w}_j^k) \right\|^2. \end{aligned} \quad (5)$$

Iterate updates (4) can be simplified as:

$$\mathbf{w}_i^{k+1} = \underset{\mathbf{w}_i}{\operatorname{argmin}} \mathcal{L}_i^s(\mathbf{w}_i, \mathbf{w}_i^k, \{\mathbf{w}_j^k\}, \gamma_i^k), \quad (6a)$$

$$\gamma_i^{k+1} = \gamma_i^k - \frac{\rho}{2} \sum_{j \in \mathcal{N}_i} (\mathbf{w}_i^{k+1} - \mathbf{w}_j^{k+1}), \quad (6b)$$

where Eq. (6a) is regarded as the primal variable update while Eq. (6b) is known as the dual variable update.

2.3 Privacy Concern

In ADMM, iterates are updated by solving a minimization associated with the local dataset (Eq. (6a)), and they are needed to be shared with neighbours. If the local training data is sensitive, the shared iterates would cause privacy leakage.

The main goal of this paper is to provide privacy protection against inference attacks from an adversary, who tries to infer sensitive information about the agents' private datasets from the shared messages.

In order to provide privacy guarantee against such attacks, we define our privacy model formally by the notion of differential privacy (Dwork et al., 2006b, 2014). Specifically, we adopt the (ϵ, δ) -differential privacy defined as follows:

Definition 1 ((ϵ, δ) -Differential Privacy) *A randomized mechanism \mathcal{M} is (ϵ, δ) -differentially private if for any two neighbouring datasets \mathcal{D} and \mathcal{D}' differing in only one tuple, and for any output subset $\mathcal{O} \subseteq \text{range}(\mathcal{M})$:*

$$\Pr [\mathcal{M}(\mathcal{D}) \in \mathcal{O}] \leq e^\epsilon \cdot \Pr [\mathcal{M}(\mathcal{D}') \in \mathcal{O}] + \delta, \quad (7)$$

which means, with probability of at least $1 - \delta$, the ratio of the probability distributions for two neighboring datasets is bounded by e^ϵ .

In Definition 1, δ and ϵ indicate the strength of privacy protection from the mechanism (a smaller ϵ or a smaller δ gives better privacy protection). Gaussian mechanism is a widely used randomization method used to guarantee (ϵ, δ) -differential privacy, where calibrated noise sampled from normal (Gaussian) distribution is added to the output.

To analyze the privacy guarantee from a class of differentially private algorithms under t -fold adaptive composition, we use the following advanced composition theorem to analyze the privacy guarantee.

Theorem 2 (Advanced Composition) *Let $\epsilon, \delta \geq 0$. The class of (ϵ, δ) -differentially private algorithms satisfies (ϵ', δ) -differential privacy under t -fold adaptive composition, where $\epsilon' = c_0 \sqrt{t\epsilon}$ for some constant c_0 .*

Proof The proof of the advanced composition theorem is based on the moments accountant method proposed in (Abadi et al., 2016). In moments accountant method, the τ -th log moments of privacy loss from each (ϵ, δ) -differentially private algorithm can be given by $\frac{t\tau(\tau+1)\epsilon^2}{4\ln(1.25/\delta)}$. According to the linear composability of the log moments, we obtain the log moment of the total privacy loss from the class of private algorithms by $\frac{t\tau(\tau+1)\epsilon^2}{4\ln(1.25/\delta)}$. By using the tail bound property of the log moment, we can obtain the relationship between ϵ and ϵ' , which is $\epsilon' \geq \sqrt{\frac{t\ln(1/\delta)}{\ln(1.25/\delta)}}\epsilon$. Thus, there exists a constant c_0 so that $\epsilon' = c_0 \sqrt{t\epsilon}$. Due to the limited space here, we suggest readers to refer to the previous works (Abadi et al., 2016) for the details. \blacksquare

3. Distributed ADMM with Differential Privacy

In this section, we will briefly describe the main idea of our approach, then introduce our algorithm, and give the privacy analysis.

3.1 Main Idea

As discussed in the last section, the privacy concern in ADMM-based distributed learning comes from the exchanged iterates. Calibrated noise is added to the iterates to control the privacy leakage and guarantee differential privacy. However, as pointed in (Huang et al., 2019), directly combining noise addition and conventional ADMM is not a good option since the added noise would disrupt the convergence property of ADMM. Inspired by the previous work (Huang et al., 2019), we adopt the approximation to introduce the calibrated noise into ADMM robustly. Such approximation is also widely used in linearized ADMM (Ling and Ribeiro, 2014) to reduce computation cost. In addition, instead of performing just one primal variable update per node in each iteration, our approach allows to perform l primal variable updates based on the approximate function. This idea is similar to the federated learning framework proposed by (McMahan et al., 2016).

3.2 Our Approach

Our proposed algorithm adopts the approximation when updating the primal variable with noise, and allows performing l updates per node in each iteration. Here we use $\tilde{\mathbf{w}}_i^{k,r}$ to denote the r -th noisy primal variable from agent i in the k -th ADMM iteration.

In each iteration, each node does not perform the exact minimization to update the primal variable by (6a). Instead, each node performs the inexact minimization by adopting the approximation of $\mathcal{L}_i^s(\mathbf{w}_i, \tilde{\mathbf{w}}_i^k, \{\tilde{\mathbf{w}}_j^k\}, \gamma_i^k)$ at $\tilde{\mathbf{w}}_i^{k+1,r}$ defined by:

$$\begin{aligned} \hat{\mathcal{L}}_i^s(\mathbf{w}_i, \tilde{\mathbf{w}}_i^{k+1,r}, \tilde{\mathbf{w}}_i^k, \{\tilde{\mathbf{w}}_j^k\}, \gamma_i^k) &= \sum_{j \in \mathcal{D}_j} \frac{1}{|\mathcal{D}_j|} \ell(\mathbf{a}_{i,j}, b_{i,j}, \tilde{\mathbf{w}}_i^{k+1,r}) + \frac{\lambda}{n} \phi(\tilde{\mathbf{w}}_i^{k+1,r}) \\ &+ \langle \sum_{j \in \mathcal{D}_j} \frac{1}{|\mathcal{D}_j|} \nabla \ell(\mathbf{a}_{i,j}, b_{i,j}, \tilde{\mathbf{w}}_i^{k+1,r}) + \frac{\lambda}{n} \nabla \phi(\tilde{\mathbf{w}}_i^{k+1,r}), \mathbf{w}_i - \tilde{\mathbf{w}}_i^{k+1,r} \rangle \\ &- 2 \langle \gamma_i^k, \mathbf{w}_i \rangle + \rho \sum_{j \in \mathcal{N}_i} \left\| \mathbf{w}_i - \frac{1}{2}(\tilde{\mathbf{w}}_i^k + \tilde{\mathbf{w}}_j^k) \right\|^2 + \frac{\eta_i^{k+1,r+1}}{2} \|\mathbf{w}_i - \tilde{\mathbf{w}}_i^{k+1,r}\|^2, \end{aligned} \quad (8)$$

where $\eta_i^{k+1,r+1}$ is an approximation parameter to control the distance between the updated primal variable and the previous one. Compare with $\mathcal{L}_i^s(\mathbf{w}_i, \tilde{\mathbf{w}}_i^k, \{\tilde{\mathbf{w}}_j^k\}, \gamma_i^k)$, its approximation replaces the objective function with its first-order approximation and a scalar l_2 -norm prox-function, which can lead to a closed-form solution to the primal variable update.

In addition, our approach allows l primal variable updates per node in each iteration. Thus, step (6a) is replaced by an inner l -iterative process:

$$\mathbf{w}_i^{k+1,r+1} = \min_{\mathbf{w}_i} \hat{\mathcal{L}}_i^s(\mathbf{w}_i, \tilde{\mathbf{w}}_i^{k+1,r}, \tilde{\mathbf{w}}_i^k, \{\tilde{\mathbf{w}}_j^k\}, \gamma_i^k), \quad (9a)$$

$$\tilde{\mathbf{w}}_i^{k+1,r+1} = \mathbf{w}_i^{k+1,r+1} + \boldsymbol{\xi}_i^{k+1,r+1}, \quad (9b)$$

where $\boldsymbol{\xi}_i^{k+1,r+1}$ is the sampled noise from normal (Gaussian) distribution to ensure differential privacy:

$$\boldsymbol{\xi}_i^{k+1,r+1} \sim \mathcal{N}(0, s_i^{k+1,r+1} \sigma^2 \mathbf{I}_d). \quad (10)$$

After the l -iterative primal variable updates, the dual variable update follows as:

$$\gamma_i^{k+1} = \gamma_i^k - \frac{\rho}{2} \sum_{j \in \mathcal{N}_i} (\tilde{\mathbf{w}}_i^{k+1} - \tilde{\mathbf{w}}_j^{k+1}), \quad (11)$$

where $\tilde{\mathbf{w}}_i^{k+1} = \frac{1}{l} \sum_{r=1}^l \tilde{\mathbf{w}}_i^{k+1,r}$.

The details of our approach are given in Algorithm 1. Each agent i firstly initializes his noisy primal variables $\tilde{\mathbf{w}}_i^{0,l}$ and $\tilde{\mathbf{w}}_i^0$, and dual variables γ_i^0 . Then each agent i updates his noisy primal variables by an inner l -iterative process, where $\tilde{\mathbf{w}}_i^{k+1,r+1}$ is updated by (9a) and (9b) in each inner iteration. After l iterations of the inner process, agent i obtain a noisy primal variable $\tilde{\mathbf{w}}_i^{k+1,l}$ and $\tilde{\mathbf{w}}_i^{k+1}$, and broadcast $\tilde{\mathbf{w}}_i^{k+1}$ to his neighbours $j \in \mathcal{N}_i$. After receiving the noisy primal variables $\{\tilde{\mathbf{w}}_j^{k+1}\}_{j \in \mathcal{N}_i}$ from his neighbours, agent i continues to update his dual variable γ_i^{k+1} by (11). The iterative process will continue until reaching t iterations.

Compared with the previous work (Huang et al., 2019), although we consider a different network setting, our approach has similar algorithmic feature to their approach by adopting the approximation when updating the primal variable. However, our approach allows l primal variable updates per node in each iteration. In the next section, we can prove that by performing more computation, we improve the utility of our approach.

Algorithm 1 l -Approx ADMM with Differential Privacy

```

1: Initialize  $\{\tilde{\mathbf{w}}_i^{0,l}\}_{i \in \mathcal{V}}$ ,  $\{\tilde{\mathbf{w}}_i^0\}_{i \in \mathcal{V}}$  and  $\{\gamma_i^0\}_{i \in \mathcal{V}}$ ;
2: for  $k = 0, 1, \dots, t - 1$  do
3:   for  $i \in \mathcal{V}$  do
4:     Let  $\tilde{\mathbf{w}}_i^{k+1,0} = \tilde{\mathbf{w}}_i^{k,l}$ .
5:     for  $r = 0, 1, \dots, l - 1$  do
6:       Sample  $\boldsymbol{\xi}_i^{k+1,r+1} \sim \mathcal{N}(0, s_i^{k+1,r+1^2} \sigma^2 \mathbf{I}_d)$ ;
7:       Compute  $\mathbf{w}_i^{k+1,r+1}$  by Eq. (9a);
8:       Compute  $\tilde{\mathbf{w}}_i^{k+1,r+1}$  by Eq. (9b);
9:     end for
10:    Compute  $\tilde{\mathbf{w}}_i^{k+1} = \frac{1}{l} \sum_{r=1}^l \tilde{\mathbf{w}}_i^{k+1,r}$ .
11:  end for
12:  for  $i \in \mathcal{V}$  do
13:    Broadcast  $\tilde{\mathbf{w}}_i^{k+1}$  to all neighbours  $j \in \mathcal{N}_i$ ;
14:  end for
15:  for  $i \in \mathcal{V}$  do
16:    Compute  $\gamma_i^{k+1}$  by Eq. (11).
17:  end for
18: end for

```

3.3 Privacy Analysis

In this section, we define the l_2 norm sensitivity $s_i^{k,r}$ and noise magnitude σ to achieve (ϵ, δ) -differential privacy in Algorithm 1.

Lemma 3 (L_2 -Norm Sensitivity) Assume that the loss function $\ell(\cdot)$ is c_1 -Lipschitz. The l_2 norm sensitivity of the primal variable update function (Eq. (9a)) is given by:

$$s_i^{k,r} = \frac{2c_1}{(2\rho|\mathcal{N}_i| + \eta_i^{k,r})|\mathcal{D}_i|}. \quad (12)$$

Proof The l_2 norm sensitivity of the primal variable update function (Eq. (9a)) is defined by:

$$s_i^{k,r} = \max_{\mathcal{D}_i, \mathcal{D}'_i} \|\mathbf{w}_{i,\mathcal{D}_i}^{k,r} - \mathbf{w}_{i,\mathcal{D}'_i}^{k,r}\|. \quad (13)$$

According to Eq. (8) and Eq. (9a), we obtain a closed-form solution to $\mathbf{w}_i^{k,r}$:

$$\begin{aligned} \mathbf{w}_i^{k,r} = & \frac{1}{2\rho|\mathcal{N}_i| + \eta_i^{k,r}} \left(- \sum_{j \in \mathcal{D}_j} \frac{1}{|\mathcal{D}_j|} \nabla \ell(\mathbf{a}_{i,j}, b_{i,j}, \tilde{\mathbf{w}}_i^{k,r-1}) - \frac{\lambda}{n} \nabla \phi(\tilde{\mathbf{w}}_i^{k,r-1}) \right. \\ & \left. + \gamma_i^k + \rho \sum_{j \in \mathcal{N}_i} \tilde{\mathbf{w}}_j^k + \rho |\mathcal{N}_i| \tilde{\mathbf{w}}_i^k + \eta_i^{k,r} \tilde{\mathbf{w}}_i^{k,r-1} \right). \end{aligned} \quad (14)$$

Then, we have:

$$\begin{aligned} \|\mathbf{w}_{i,\mathcal{D}_i}^{k,r} - \mathbf{w}_{i,\mathcal{D}'_i}^{k,r}\| &= \frac{\|\nabla \ell(\mathbf{a}_{i,j}, b_{i,j}, \tilde{\mathbf{w}}_i^{k,r-1}) - \nabla \ell(\mathbf{a}'_{i,j}, b'_{i,j}, \tilde{\mathbf{w}}_i^{k,r-1})\|}{(2\rho|\mathcal{N}_i| + \eta_i^{k,r})|\mathcal{D}_i|} \\ &\leq \frac{2\|\nabla \ell(\cdot)\|}{(2\rho|\mathcal{N}_i| + \eta_i^{k,r})|\mathcal{D}_i|}. \end{aligned} \quad (15)$$

Since function $\ell(\cdot)$ is c_1 -Lipschitz, we obtain the result: $s_i^{k,r} = \frac{2c_1}{(2\rho|\mathcal{N}_i| + \eta_i^{k,r})|\mathcal{D}_i|}$. ■

Theorem 4 (Privacy Guarantee) Let $\epsilon, \delta \geq 0$ be arbitrary. There exists constant c_0 so that Algorithm 1 achieves (ϵ, δ) -differential privacy if we set the noise magnitude σ in Gaussian distribution $\mathcal{N}(0, s_i^{k,r} \sigma^2 \mathbf{I}_d)$ by:

$$\sigma = \frac{c_0 \sqrt{t \cdot l \cdot 2 \ln(1.25/\delta)}}{\epsilon}. \quad (16)$$

Proof Due to the limited space, we only provide the proof sketch here. We first follow Theorem A.1. in (Dwork et al., 2014) to demonstrate that by setting $\sigma = \frac{c_0 \sqrt{t \cdot l \cdot 2 \ln(1.25/\delta)}}{\epsilon}$, the gradient function with Gaussian noise satisfies $(\frac{1}{c_0 \sqrt{t \cdot l}} \epsilon, \delta)$ -differential privacy. Since our approach includes a class of $t \times l$ differentially private gradient functions. By adopting the advanced composition (Theorem 2), we prove that Algorithm 1 achieves (ϵ, δ) -differential privacy. ■

4. Utility Analysis

In this section, we theoretically analyze the utility of Algorithm 1, which can be measured by the expected excess empirical risk with feasibility violation, namely:

$$\mathbb{E}[L_{\mathcal{D}}(\{\hat{\mathbf{w}}_i\}_{i \in \mathcal{V}}) - L_{\mathcal{D}}(\{\mathbf{w}^*\})] + \beta \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j\|, \quad (17)$$

where

$$\begin{aligned} L_{\mathcal{D}}(\{\mathbf{w}_i\}_{i \in \mathcal{V}}) &= \sum_{i \in \mathcal{V}} L_{\mathcal{D}_i}(\mathbf{w}_i) \\ &= \sum_{i \in \mathcal{V}} \left(\sum_{j \in \mathcal{D}_i} \frac{1}{|\mathcal{D}_i|} \ell(\mathbf{a}_{i,j}, b_{i,j}, \mathbf{w}_i) + \frac{\lambda}{n} \phi(\mathbf{w}_i) \right), \end{aligned} \quad (18)$$

$\{\hat{\mathbf{w}}_i\}_{i \in \mathcal{V}}$ and $\{\hat{\mathbf{w}}_i\}_{i \in \mathcal{V}}$ are the outputs of our proposed algorithm, and \mathbf{w}^* is the true minimizer of problem (1). The excess empirical risk measures the accuracy of the trained model by our approach while the feasibility violation measures the differences of local models.

4.1 Main Results

We analyze the excess empirical risk of our approach under the assumption: the objective function is convex and Lipschitz.

Theorem 5 (Utility Analysis) *Assume the objective function $L(\cdot)$ is c_2 -Lipschitz. and the diameter of the \mathcal{W} is bounded by D , namely $\sup_{\mathbf{w}, \mathbf{w}' \in \mathcal{W}} \|\mathbf{w} - \mathbf{w}'\| \leq D$. Let*

$$\hat{\mathbf{w}}_i = \frac{1}{t} \frac{1}{l} \sum_{k=1}^t \sum_{r=0}^{l-1} \tilde{\mathbf{w}}_i^{k,r}. \quad (19)$$

If we set the learning rate:

$$\eta_i^{k,r} = \frac{\sqrt{2kr}}{D} \sqrt{\frac{c_2^2}{n^2} + \frac{dc_0^2 c_1^2 t l 8 \ln(1.25/\delta)}{\epsilon^2 |\mathcal{D}_i|^2}}, \quad (20)$$

we have the following expected error bound:

$$\begin{aligned} &\mathbb{E} \left[L_{\mathcal{D}}(\{\hat{\mathbf{w}}_i\}) - L_{\mathcal{D}}(\{\mathbf{w}^*\}) \right] + \beta \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j\| \\ &\leq \sum_{i \in \mathcal{V}} \left(\frac{\sqrt{2}D}{\sqrt{t \cdot l}} \left(\frac{c_2^2}{n^2} + \frac{dc_0^2 c_1^2 t l 8 \ln(1.25/\delta)}{\epsilon^2 |\mathcal{D}_i|^2} \right)^{\frac{1}{2}} + \frac{\rho |\mathcal{N}_i| D^2 + |\mathcal{N}_i| \beta^2 / \rho}{t} \right). \end{aligned} \quad (21)$$

Proof See the proof in Appendix 8.1. ■

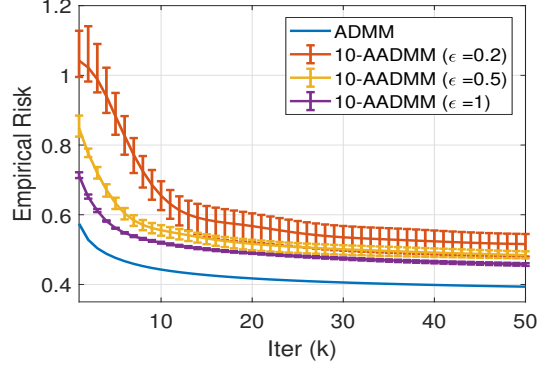


Figure 1: Utility and Privacy Trade-off.

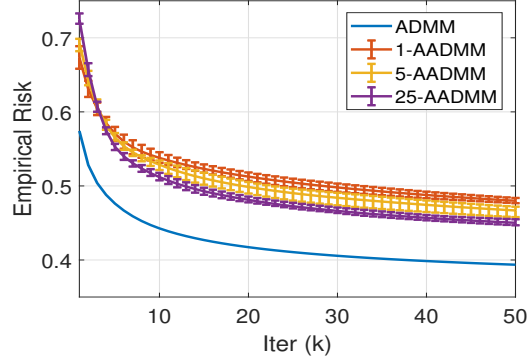


Figure 2: Improve Utility by More Computation.

4.2 Discussion

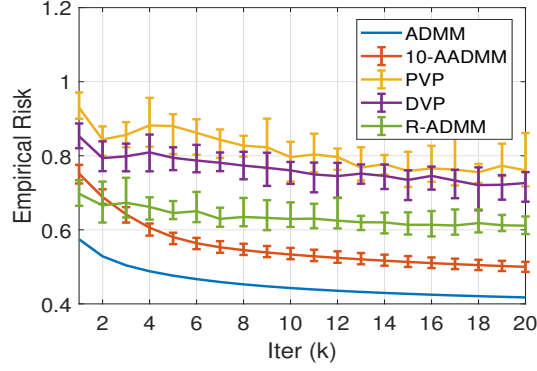
Our main result (Theorem 5) demonstrates the privacy-utility trade-off of our approach. When we want to ensure smaller privacy leakage from our output by setting a smaller ϵ , the utility of our approach will decrease.

Theorem 5 demonstrates that our approach can be improved by setting a larger l , which allows each node to perform more computations in each iteration, but it brings higher computation cost.

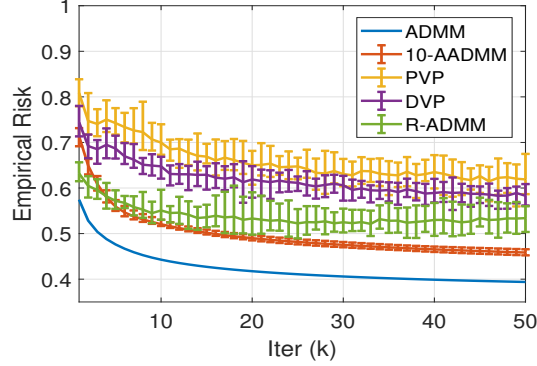
From Theorem 5, if the iteration number t is sufficiently large enough our approach achieve the error bound $O\left(\sum_{i \in \mathcal{V}} \frac{\sqrt{d \ln(1/\delta)}}{|\mathcal{D}_i| \epsilon}\right)$, which is comparable to the state-of-art error bound $O\left(\frac{\sqrt{d \ln(1/\delta)}}{N \epsilon}\right)$ for differentially private empirical risk minimization under the assumption that the objective is Lischipz and convex, here N is the total number of training data.

5. Numerical Results

In this section, we give the numerical results of our approach and compare our proposed algorithm (Algorithm 1) with algorithms proposed by prior works by the simulation in MATLAB.



(a) $\epsilon = 0.5, \delta = 10^{-5}$



(b) $\epsilon = 1, \delta = 10^{-5}$

Figure 3: Accuracy comparison in empirical risk.

5.1 L_2 Regularized Logistic Regression

We evaluate our approach by l_2 regularized logistic regression. Logistic regression is a widely used statistical model for classification, and its loss function is described as:

$$\ell(\mathbf{a}, b, \mathbf{w}) = \log(1 + \exp(-b\mathbf{w}^\top \mathbf{a})). \quad (22)$$

Thus by l_2 regularized logistic regression, the objective of our regularized empirical risk minimization problem can be formulated as:

$$L_{\mathcal{D}}(\{\mathbf{w}_i\}) = \sum_{i \in \mathcal{V}} \left(\sum_{j \in \mathcal{D}_i} \frac{1}{|\mathcal{D}_i|} \log(1 + \exp(-b_{i,j} \mathbf{w}_i^\top \mathbf{a}_{i,j})) + \frac{\lambda}{2n} \|\mathbf{w}_i\|^2 \right), \quad (23)$$

5.2 Dataset

The dataset used in our simulation is Adult dataset (Asuncion and Newman, 2007) from UCI Machine Learning Repository. Adult dataset includes 48,842 instances, each of which has 14 personal attributes with a label representing whether the income is above \$50,000 or not. We follow the previous works (Zhang et al., 2018a; Huang et al., 2019) to preprocess the data by removing all the instances with missing values, converting the categorical attributes

into binary vectors, normalizing columns to guarantee the maximum value of each column is 1, normalizing rows to enforce their l_2 norm to be less than 1, and converting the labels $\{> 50k, < 50k\}$ into $\{+1, -1\}$. After the data preprocessing, we obtain 45,222 data instances each with a 104-dimensional feature vector and a label belonging to $\{+1, -1\}$.

5.3 Baseline Algorithms

We compare our approach: l -Approx ADMM (l -AADMM) with four baseline algorithms: (1) non-private distributed ADMM algorithm, (2) ADMM algorithm with PVP in (Zhang and Zhu, 2017), (3) ADMM with dual variable perturbation (DVP) in (Zhang and Zhu, 2017), and (4) Recycled ADMM (R-ADMM) proposed in (Zhang et al., 2018b).

5.4 Simulation Results

We mainly evaluate our approach on the utility-privacy trade-off and the effect of the choice of l on utility, and compare our approach with the baseline algorithms. In our simulation, we set ρ to be 0.001 and μ to be 0.0001, and by data preprocessing, we can enforce the objective function to be $(n + \frac{\mu D}{n})$ -Lipschitz, where n is the number of agents and D is the diameter of \mathcal{W} . Here, we consider a network consisting of 100 agents fully connected and with evenly divided data. Our numerical results are the averaged results from 10 simulations.

Figure 1 shows the utility-privacy trade-off of our approach. Here we fix δ to be 10^{-5} and l to be 10, and only change ϵ . With increasing ϵ indicating weaker privacy guarantee, our approach has less empirical risk achieving better utility, which is consistent with Theorem 5.

Figure 2 demonstrates the impact of the choice of l on the utility of our approach. In this simulation, we fix ϵ to be 1 and change l from 1 to 25. When we set a larger l , the accuracy of our algorithm can be improved, but it also brings higher computation cost.

Figure 3 compares our approach with the four baseline algorithms on empirical risk when we set the privacy parameter ϵ to be 0.5 and 1, and δ to be 0.00001. The results show that our approach has more stable update processing and has better utility than the other differentially private ADMM algorithms.

6. Related Work

6.1 Distributed ADMM

ADMM demonstrates fast convergence in many applications and it is widely used to solve distributed optimizations. Shi et al. (Shi et al., 2014) focus the theoretical aspects on the convergence rate of distributed ADMM, and demonstrate that the iterates from distributed ADMM can converge to the optimal solution linearly under the assumptions that the objective function is strongly convex and Lipschitz smooth. Zhang and Kwok (Zhang and Kwok, 2014) propose an asynchronous distributed ADMM by using partial barrier and bounded delay. Ling et al. (Ling et al., 2015) design a linearized distributed ADMM where the augmented Lagrangian function is replaced by its first-order approximation to reduce the local computation cost. Song et al. (Song et al., 2016) show that distributed ADMM can converge faster by adaptively choosing the penalty parameter. Makhdomi and Ozdaglar

(Makhdoumi and Ozdaglar, 2017) demonstrate that the objective value with feasibility violation converges to the optimum at a rate of $O(1/t)$ by distributed ADMM, where t is the number of iterations.

6.2 Differentially Private Empirical Risk Minimization

There have been tremendous research efforts on differentially private empirical risk minimization (Chaudhuri et al., 2011; Bassily et al., 2014; Wang et al., 2017; Thakurta and Smith, 2013). Chaudhuri et al. (Chaudhuri et al., 2011) propose two perturbation methods: output perturbation and objective perturbation to guarantee ϵ -differential privacy. Bassily et al. (Bassily et al., 2014) provide a systematic investigation of differentially private algorithms for convex empirical risk minimization and propose efficient algorithms with tighter error bound. Wang et al. (Wang et al., 2017) focus on a more general problem: non-convex problem, and propose a faster algorithm based on a proximal stochastic gradient method.

6.3 Differentially Private ADMM-based Distributed Learning

Recently, there are some works focusing on differentially private ADMM-based distributed learning algorithms. Zhang and Zhu (Zhang and Zhu, 2017) propose two perturbation methods: primal perturbation and dual perturbation to guarantee dynamic differential privacy in ADMM-based distributed learning. Zhang et al. (Zhang et al., 2018a) propose to perturb the penalty parameter of ADMM to guarantee differential privacy. Huang et al. (Huang et al., 2019) propose an algorithm named DP-ADMM, where an approximate augmented Lagrangian function with time-varying Gaussian noise addition is adopted to update iterates while guaranteeing differential privacy, and theoretically analyze the convergence rate of their approach. Zhang et al. (Zhang et al., 2018b) propose recycled ADMM with differential privacy guarantee where the results from odd iterations could be re-utilized by the even iterations, and thus half of updates incur no privacy leakage. Hu et al. (Hu et al., 2019) consider a setting where data features are distributed, and use ADMM with primal variable perturbation for distributed learning while guaranteeing differential privacy.

7. Conclusion

In this paper, we have proposed a new differentially private distributed ADMM algorithm for a class of convex learning problems. In our approach, we have adopted the approximation when updating the primal variables with differential privacy and have allowed each node to perform such primal variable updates for l times in each iteration. We have analyzed the privacy guarantee of our proposed algorithm by properly setting the noise magnitude in Gaussian distribution and using the moments accountant method. We have theoretically analyzed the utility of our approach by the excess empirical risk with feasibility violation under the setting that the objective is Lipschitz and convex. Our theoretical results have shown that our approach can obtain higher utility if we set a larger l and can achieve error bounds, which are comparable to the state-of-art error bounds for differentially private empirical risk minimization.

8. Appendix

8.1 Proof of Theorem 5

Firstly, by assuming that $\|\gamma_{i,j}\| \leq \beta$, we have:

$$\begin{aligned} & \mathbb{E} \left[L_{\mathcal{D}}(\{\hat{\mathbf{w}}_i\}) - L_{\mathcal{D}}(\{\mathbf{w}^*\}) \right] + \beta \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j\| \\ &= \max_{\gamma_{i,j}: \|\gamma_{i,j}\| \leq \beta} \mathbb{E} \left[L_{\mathcal{D}}(\{\hat{\mathbf{w}}_i\}) - L_{\mathcal{D}}(\{\mathbf{w}^*\}) - \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \langle \gamma_{i,j}, \hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j \rangle \right]. \end{aligned} \quad (24)$$

Due to the convexity of $L_{\mathcal{D}}(\cdot)$ and the definition of $\hat{\mathbf{w}}_i$, we have:

$$\begin{aligned} & L_{\mathcal{D}}(\{\hat{\mathbf{w}}_i\}) - L_{\mathcal{D}}(\{\mathbf{w}^*\}) - \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \langle \gamma_{i,j}, \hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j \rangle \\ & \leq \frac{1}{t} \sum_{k=1}^t \frac{1}{l} \sum_{r=0}^{l-1} (L_{\mathcal{D}}(\{\tilde{\mathbf{w}}_i^{k,r}\}) - L_{\mathcal{D}}(\{\mathbf{w}^*\}) - \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \langle \gamma_{i,j}, \tilde{\mathbf{w}}_i^{k,r} - \tilde{\mathbf{w}}_j^{k,r} \rangle) \\ & \leq \frac{1}{t} \sum_{k=1}^t \frac{1}{l} \sum_{r=0}^{l-1} \sum_{i \in \mathcal{V}} (\langle \nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}), \tilde{\mathbf{w}}_i^{k,r} - \mathbf{w}^* \rangle - \sum_{j \in \mathcal{N}_i} \langle \gamma_{i,j}, \tilde{\mathbf{w}}_i^{k,r} - \tilde{\mathbf{w}}_j^{k,r} \rangle) \end{aligned} \quad (25)$$

Next, we analyze $\langle \nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}), \tilde{\mathbf{w}}_i^{k,r} - \mathbf{w}^* \rangle$:

$$\begin{aligned} \sum_{i \in \mathcal{V}} \langle \nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}), \tilde{\mathbf{w}}_i^{k,r} - \mathbf{w}^* \rangle &= \sum_{i \in \mathcal{V}} (\langle \nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}) + \boldsymbol{\xi}_i, \tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^* \rangle \\ & \quad + \langle \boldsymbol{\xi}_i, \mathbf{w}^* - \tilde{\mathbf{w}}_i^{k,r+1} \rangle + \langle \nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}) + \boldsymbol{\xi}_i, \tilde{\mathbf{w}}_i^{k,r} - \tilde{\mathbf{w}}_i^{k,r+1} \rangle). \end{aligned} \quad (26)$$

If we define: $\boldsymbol{\xi}_i = \boldsymbol{\eta}_i^{k,r} / (2\rho|\mathcal{N}_i| + \eta_i^{k,r})$, according to the primal variable update (9a) and (9b), we have:

$$\begin{aligned} & \frac{1}{l} \sum_{r=0}^{l-1} \langle \nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}) + \boldsymbol{\xi}_i, \tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^* \rangle \\ &= \frac{1}{l} \sum_{r=0}^{l-1} (\langle \nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}) - 2\gamma_i^{k-1} - 2\rho \sum_{j \in \mathcal{N}_i} (\tilde{\mathbf{w}}_i^{k,r} - \frac{1}{2}(\tilde{\mathbf{w}}_i^{k-1} + \tilde{\mathbf{w}}_j^{k-1})) + \boldsymbol{\xi}_i, \tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^* \rangle \\ & \quad + 2\langle \gamma_i^{k-1} - \rho \sum_{j \in \mathcal{N}_i} (\tilde{\mathbf{w}}_i^{k,r} - \frac{1}{2}(\tilde{\mathbf{w}}_i^{k-1} + \tilde{\mathbf{w}}_j^{k-1})), \tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^* \rangle) \\ &= \frac{1}{l} \sum_{r=0}^{l-1} ((\eta_i^{k,r+1} + 2\rho|\mathcal{N}_i|) \langle \tilde{\mathbf{w}}_i^{k,r} - \tilde{\mathbf{w}}_i^{k,r+1}, \tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^* \rangle \\ & \quad + 2\langle \gamma_i^{k-1} - \rho \sum_{j \in \mathcal{N}_i} (\tilde{\mathbf{w}}_i^{k,r} - \frac{1}{2}(\tilde{\mathbf{w}}_i^{k-1} + \tilde{\mathbf{w}}_j^{k-1})), \tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^* \rangle). \end{aligned} \quad (27)$$

Based on the dual update (11) and the definition of $\tilde{\mathbf{w}}_i^k$, we have:

$$\begin{aligned}
& \frac{1}{l} \sum_{r=0}^{l-1} 2 \langle \gamma_i^{k-1} - \rho \sum_{j \in \mathcal{N}_i} (\tilde{\mathbf{w}}_i^{k,r} - \frac{1}{2}(\tilde{\mathbf{w}}_i^{k-1} + \tilde{\mathbf{w}}_j^{k-1})), \tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^* \rangle \\
&= \frac{1}{l} \sum_{r=0}^{l-1} 2\rho \sum_{j \in \mathcal{N}_i} \langle \tilde{\mathbf{w}}_i^{k,r+1} - \tilde{\mathbf{w}}_i^k, \mathbf{w}^* - \tilde{\mathbf{w}}_i^{k,r+1} \rangle + \frac{1}{l} \sum_{r=0}^{l-1} 2\rho \sum_{j \in \mathcal{N}_i} \langle \tilde{\mathbf{w}}_i^{k,r} - \tilde{\mathbf{w}}_i^{k,r+1}, \mathbf{w}^* - \tilde{\mathbf{w}}_i^{k,r+1} \rangle \\
& \quad + 2\rho \sum_{j \in \mathcal{N}_i} \langle \frac{1}{2}(\tilde{\mathbf{w}}_i^k + \tilde{\mathbf{w}}_j^k) - \frac{1}{2}(\tilde{\mathbf{w}}_i^{k-1} + \tilde{\mathbf{w}}_j^{k-1}), \mathbf{w}^* - \tilde{\mathbf{w}}_i^k \rangle + 2 \sum_{j \in \mathcal{N}_i} \langle \gamma_{i,j}^k, \tilde{\mathbf{w}}_i^k - \mathbf{w}^* \rangle.
\end{aligned} \tag{28}$$

Since we have:

$$\begin{aligned}
& \langle \tilde{\mathbf{w}}_i^{k,r} - \tilde{\mathbf{w}}_i^{k,r+1}, \tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^* \rangle \\
&= \frac{1}{2} (\|\tilde{\mathbf{w}}_i^{k,r} - \mathbf{w}^*\|^2 - \|\tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^*\|^2 - \|\tilde{\mathbf{w}}_i^{k,r} - \tilde{\mathbf{w}}_i^{k,r+1}\|^2),
\end{aligned} \tag{29}$$

and

$$\sum_{r=0}^{l-1} \langle \tilde{\mathbf{w}}_i^{k,r+1} - \tilde{\mathbf{w}}_i^k, \mathbf{w}^* - \tilde{\mathbf{w}}_i^{k,r+1} \rangle = \frac{1}{l} \sum_{r=0}^{l-1} \sum_{a=0}^{l-1} \langle \tilde{\mathbf{w}}_i^{k,r+1} - \tilde{\mathbf{w}}_i^{k,a+1}, \mathbf{w}^* - \tilde{\mathbf{w}}_i^{k,r+1} \rangle < 0, \tag{30}$$

and

$$\begin{aligned}
& \langle \frac{1}{2}(\tilde{\mathbf{w}}_i^k + \tilde{\mathbf{w}}_j^k) - \frac{1}{2}(\tilde{\mathbf{w}}_i^{k-1} + \tilde{\mathbf{w}}_j^{k-1}), \mathbf{w}^* - \tilde{\mathbf{w}}_i^k \rangle \\
& \leq \frac{1}{2} (\|\frac{1}{2}(\tilde{\mathbf{w}}_i^{k-1} + \tilde{\mathbf{w}}_j^{k-1}) - \mathbf{w}^*\|^2 - \|\frac{1}{2}(\tilde{\mathbf{w}}_i^k + \tilde{\mathbf{w}}_j^k) - \mathbf{w}^*\|^2 + \|\frac{1}{2}(\tilde{\mathbf{w}}_i^k - \tilde{\mathbf{w}}_j^k)\|^2) \\
& = \frac{1}{2} (\|\frac{1}{2}(\tilde{\mathbf{w}}_i^{k-1} + \tilde{\mathbf{w}}_j^{k-1}) - \mathbf{w}^*\|^2 - \|\frac{1}{2}(\tilde{\mathbf{w}}_i^k + \tilde{\mathbf{w}}_j^k) - \mathbf{w}^*\|^2) - \frac{1}{2\rho^2} \|\gamma_{i,j}^{k-1} - \gamma_{i,j}^k\|^2,
\end{aligned} \tag{31}$$

and by Young's inequality:

$$\begin{aligned}
& \langle \nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}) + \boldsymbol{\xi}_i, \tilde{\mathbf{w}}_i^{k,r} - \tilde{\mathbf{w}}_i^{k,r+1} \rangle \\
& \leq \frac{1}{2\eta_i^{k,r+1}} \|\nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}) + \boldsymbol{\xi}_i\|^2 + \frac{\eta_i^{k,r+1}}{2} \|\tilde{\mathbf{w}}_i^{k,r} - \tilde{\mathbf{w}}_i^{k,r+1}\|^2,
\end{aligned} \tag{32}$$

we can obtain:

$$\begin{aligned}
& \frac{1}{l} \sum_{r=0}^{l-1} \langle \nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}), \tilde{\mathbf{w}}_i^{k,r} - \mathbf{w}^* \rangle \\
& \leq \frac{1}{l} \sum_{r=0}^{l-1} (\frac{\eta_i^{k,r+1}}{2} (\|\tilde{\mathbf{w}}_i^{k,r} - \mathbf{w}^*\|^2 + \|\tilde{\mathbf{w}}_i^{k,r+1} - \mathbf{w}^*\|^2) + \frac{1}{2\eta_i^{k,r+1}} \|\nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}) + \boldsymbol{\xi}_i\|^2) \\
& \quad + \rho \sum_{j \in \mathcal{N}_i} (\|\frac{1}{2}(\tilde{\mathbf{w}}_i^{k-1} + \tilde{\mathbf{w}}_j^{k-1}) - \mathbf{w}^*\|^2 - \|\frac{1}{2}(\tilde{\mathbf{w}}_i^k + \tilde{\mathbf{w}}_j^k) - \mathbf{w}^*\|^2) \\
& \quad + \sum_{j \in \mathcal{N}_i} \frac{1}{\rho} \|\gamma_{i,j}^{k-1} - \gamma_{i,j}^k\|^2 + \frac{1}{l} \sum_{r=0}^{l-1} \langle \boldsymbol{\xi}_i, \mathbf{w}^* - \tilde{\mathbf{w}}_i^{k,r+1} \rangle + \frac{1}{l} \sum_{r=0}^{l-1} \sum_{j \in \mathcal{N}_i} 2 \langle \gamma_{i,j}^k, \tilde{\mathbf{w}}_i^{k+1} - \mathbf{w}^* \rangle.
\end{aligned} \tag{33}$$

Next, we analyze $\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \langle -\gamma_{i,j}, \tilde{\mathbf{w}}_i^k - \tilde{\mathbf{w}}_j^k \rangle$:

$$\begin{aligned}
& \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \langle -\gamma_{i,j}, \tilde{\mathbf{w}}_i^k - \tilde{\mathbf{w}}_j^k \rangle \\
&= \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} (\langle -\gamma_{i,j}^k, \tilde{\mathbf{w}}_i^k - \mathbf{w}^* \rangle + \langle \gamma_{i,j}^k, \tilde{\mathbf{w}}_j^k - \mathbf{w}^* \rangle + 2\langle \gamma_{i,j}^k - \gamma_{i,j}, \frac{1}{2}(\tilde{\mathbf{w}}_i^k - \tilde{\mathbf{w}}_j^k) \rangle) \\
&= \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} (\langle -2\gamma_{i,j}^k, \tilde{\mathbf{w}}_i^k - \mathbf{w}^* \rangle + 2\langle \gamma_{i,j}^k - \gamma_{i,j}, \frac{1}{2}(\tilde{\mathbf{w}}_i^k - \tilde{\mathbf{w}}_j^k) \rangle).
\end{aligned} \tag{34}$$

Furthermore, we have:

$$\begin{aligned}
\langle \gamma_{i,j}^k - \gamma_{i,j}, \frac{1}{2}(\tilde{\mathbf{w}}_i^k - \tilde{\mathbf{w}}_j^k) \rangle &= \frac{1}{\rho} \langle \gamma_{i,j}^k - \gamma_{i,j}, \gamma_{i,j}^{k-1} - \gamma_{i,j}^k \rangle \\
&= \frac{1}{2\rho} (\|\gamma_{i,j}^{k-1} - \gamma_{i,j}\|^2 - \|\gamma_{i,j}^k - \gamma_{i,j}\|^2 - \|\gamma_{i,j}^{k-1} - \gamma_{i,j}^k\|^2).
\end{aligned} \tag{35}$$

Since we assume $L_{\mathcal{D}}(\cdot)$ is c_2 -Lipschitz, we have:

$$\mathbb{E}[\|\nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}) + \boldsymbol{\xi}_i\|^2] = \frac{c_2^2}{n^2} + \frac{dc_0^2c_1^2t\ln(1.25/\delta)}{\epsilon^2|\mathcal{D}_i|^2}. \tag{36}$$

Since we have $\mathbb{E}[\langle \boldsymbol{\xi}_i, \mathbf{w}^* - \mathbf{w}_i^{k,r} \rangle] = 0$, by assuming that the diameter of the \mathcal{W} is bounded by D , and let $\eta_i^{k,r} = \frac{\sqrt{2kr}}{D} \sqrt{\frac{c_2^2}{n^2} + \frac{dc_0^2c_1^2t\ln(1.25/\delta)}{\epsilon^2|\mathcal{D}_i|^2}}$, we can obtain:

$$\begin{aligned}
& \mathbb{E} \left[L_{\mathcal{D}}(\{\hat{\mathbf{w}}_i\}) - L_{\mathcal{D}}(\{\mathbf{w}^*\}) \right] + \beta \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j\| \\
&\leq \sum_{i \in \mathcal{V}} \left(\frac{1}{t} \frac{1}{l} \sum_{k=1}^t \sum_{r=0}^{l-1} \mathbb{E} \left[\frac{\|\nabla L_{\mathcal{D}_i}(\tilde{\mathbf{w}}_i^{k,r}) + \boldsymbol{\xi}_i\|^2}{2\eta_i^{k,r+1}} \right] + \frac{1}{t} \frac{1}{l} \frac{\eta_i^{t,l}}{2} D^2 + \frac{1}{t} \rho |\mathcal{N}_i| D^2 \right. \\
&\quad \left. + \frac{1}{t} \frac{|\mathcal{N}_i|}{\rho} \max_{\gamma_{i,j}: \|\gamma_{i,j}\| \leq \beta} \|\gamma_{i,j}^0 - \gamma_{i,j}\|^2 \right) \\
&\leq \sum_{i \in \mathcal{V}} \left(\frac{\sqrt{2}D}{\sqrt{t \cdot l}} \left(\frac{c_2^2}{n^2} + \frac{dc_0^2c_1^2t\ln(1.25/\delta)}{\epsilon^2|\mathcal{D}_i|^2} \right)^{\frac{1}{2}} + \frac{\rho |\mathcal{N}_i| D^2 + |\mathcal{N}_i| \beta^2 / \rho}{t} \right).
\end{aligned} \tag{37}$$

References

- Martín Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 308–318. ACM, 2016.
- Arthur Asuncion and David Newman. Uci machine learning repository, 2007.
- Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pages 464–473. IEEE, 2014.

- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.
- Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(Mar):1069–1109, 2011.
- John C Duchi, Alekh Agarwal, and Martin J Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control*, 57(3):592–606, 2011.
- Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni 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, 2006a.
- Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography Conference*, pages 265–284. Springer, 2006b.
- Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- Pedro A Forero, Alfonso Cano, and Georgios B Giannakis. Consensus-based distributed support vector machines. *Journal of Machine Learning Research*, 11(May):1663–1707, 2010.
- Yanmin Gong, Yuguang Fang, and Yuanxiong Guo. Private data analytics on biomedical sensing data via distributed computation. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 13(3):431–444, 2016.
- Shuo Han, Ufuk Topcu, and George J Pappas. Differentially private distributed constrained optimization. *IEEE Transactions on Automatic Control*, 62(1):50–64, 2017.
- Yaochen Hu, Peng Liu, Linglong Kong, and Di Niu. Learning privately over distributed features: An admm sharing approach. *arXiv preprint arXiv:1907.07735*, 2019.
- Zonghao Huang, Rui Hu, Yuanxiong Guo, Eric Chan-Tin, and Yanmin Gong. Dp-admm: Admm-based distributed learning with differential privacy. *IEEE Transactions on Information Forensics and Security*, 2019.
- Qing Ling and Alejandro Ribeiro. Decentralized linearized alternating direction method of multipliers. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5447–5451. IEEE, 2014.
- Qing Ling, Wei Shi, Gang Wu, and Alejandro Ribeiro. Dlm: Decentralized linearized alternating direction method of multipliers. *IEEE Transactions on Signal Processing*, 63(15):4051–4064, 2015.

- Ali Makhdoumi and Asuman Ozdaglar. Convergence rate of distributed admm over networks. *IEEE Transactions on Automatic Control*, 62(10):5082–5095, 2017.
- H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, et al. Communication-efficient learning of deep networks from decentralized data. *arXiv preprint arXiv:1602.05629*, 2016.
- Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48, 2009.
- Angelia Nedic, Alex Olshevsky, Asuman Ozdaglar, and John N Tsitsiklis. Distributed subgradient methods and quantization effects. In *2008 47th IEEE Conference on Decision and Control*, pages 4177–4184. IEEE, 2008.
- Wei Shi, Qing Ling, Kun Yuan, Gang Wu, and Wotao Yin. On the linear convergence of the admm in decentralized consensus optimization. *IEEE Transactions on Signal Processing*, 62(7):1750–1761, 2014.
- Changkyu Song, Sejong Yoon, and Vladimir Pavlovic. Fast admm algorithm for distributed optimization with adaptive penalty. In *Thirtieth AAAI conference on artificial intelligence*, 2016.
- Abhradeep Guha Thakurta and Adam Smith. Differentially private feature selection via stability arguments, and the robustness of the lasso. In *Conference on Learning Theory*, pages 819–850, 2013.
- Konstantinos I Tsianos, Sean Lawlor, and Michael G Rabbat. Push-sum distributed dual averaging for convex optimization. In *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, pages 5453–5458. IEEE, 2012.
- Di Wang, Minwei Ye, and Jinhui Xu. Differentially private empirical risk minimization revisited: Faster and more general. In *Advances in Neural Information Processing Systems*, pages 2722–2731, 2017.
- Ruiliang Zhang and James Kwok. Asynchronous distributed admm for consensus optimization. In *International Conference on Machine Learning*, pages 1701–1709, 2014.
- Tao Zhang and Quanyan Zhu. Dynamic differential privacy for admm-based distributed classification learning. *IEEE Transactions on Information Forensics and Security*, 12(1):172–187, 2017.
- Xueru Zhang, Mohammad Mahdi Khalili, and Mingyan Liu. Improving the privacy and accuracy of admm-based distributed algorithms. *arXiv preprint arXiv:1806.02246*, 2018a.
- Xueru Zhang, Mohammad Mahdi Khalili, and Mingyan Liu. Recycled admm: Improve privacy and accuracy with less computation in distributed algorithms. In *2018 56th Annual Allerton Conference on Communication, Control, and Computing*, pages 959–965. IEEE, 2018b.