

# Towards Plausible Differentially Private ADMM Based Distributed Machine Learning

Jiahao Ding  
University of Houston  
jding7@uh.edu

Jingyi Wang  
San Francisco State University  
jingyiwang@sfsu.edu

Guannan Liang  
University of Connecticut  
guannan.liang@uconn.edu

Jinbo Bi  
University of Connecticut  
jinbo.bi@uconn.edu

Miao Pan  
University of Houston  
mpan2@uh.edu

## ABSTRACT

The Alternating Direction Method of Multipliers (ADMM) and its distributed version have been widely used in machine learning. In the iterations of ADMM, model updates using local private data and model exchanges among agents impose critical privacy concerns. Despite some pioneering works to relieve such concerns, differentially private ADMM still confronts many research challenges. For example, the guarantee of differential privacy (DP) relies on the premise that the optimality of each local problem can be perfectly attained in each ADMM iteration, which may never happen in practice. The model trained by DP ADMM may have low prediction accuracy. In this paper, we address these concerns by proposing a novel (Improved) Plausible differentially Private ADMM algorithm, called PP-ADMM and IPP-ADMM. In PP-ADMM, each agent approximately solves a perturbed optimization problem that is formulated from its local private data in an iteration, and then perturbs the approximate solution with Gaussian noise to provide the DP guarantee. To further improve the model accuracy and convergence, an improved version IPP-ADMM adopts sparse vector technique (SVT) to determine if an agent should update its neighbors with the current perturbed solution. The agent calculates the difference of the current solution from that in the last iteration, and if the difference is larger than a threshold, it passes the solution to neighbors; or otherwise the solution will be discarded. Moreover, we propose to track the total privacy loss under the zero-concentrated DP (zCDP) and provide a generalization performance analysis. Experiments on real-world datasets demonstrate that under the same privacy guarantee, the proposed algorithms are superior to the state of the art in terms of model accuracy and convergence rate.

## KEYWORDS

differential privacy; distributed machine learning; ADMM; decentralized optimization

## ACM Reference Format:

Jiahao Ding, Jingyi Wang, Guannan Liang, Jinbo Bi, and Miao Pan. 2020. Towards Plausible Differentially Private ADMM Based Distributed Machine Learning. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Nowadays, the development of machine learning creates many emerging applications that can improve the quality of our life, such as medical diagnosis, autonomous driving, face recognition, etc. With the proliferation of mobile phones and Internet-of-things devices, a vast amount of data have been generated at an ever-increasing rate, which leads to significant computational complexity for data collection and processes via a centralized machine learning approach. Therefore, distributed machine learning plays an increasingly important role in dealing with large scale machine learning tasks. There are many research efforts on distributed training a large scale optimization problem, which mainly consist of two types: (sub)gradient based methods, and Alternating Direction Method of Multipliers (ADMM) based methods. As shown in [2], the convergence of (sub)gradient based methods are usually slow, which is  $O(1/\sqrt{T})$  under general convex objectives, while ADMM based algorithms can achieve  $O(1/T)$  convergence rate, where  $T$  is the number of iterations. Therefore, ADMM has been widely used in distributed machine learning [17, 22], and in this paper, we focus on the distributed ADMM.

In the framework of distributed ADMM, data providers (agents) collaboratively solve a learning problem, which can be decomposed into several subproblems, via an interactive procedure of local computation and message passing. However, the information exchanges during this process raise serious privacy concerns, and the adversary can extract private information from the shared learning models via various inference attacks [16, 19]. To prevent such privacy leakage, differential privacy (DP) [9] provides a de facto standard of privacy definition for protecting data privacy, which guarantees that the adversary with arbitrary background knowledge cannot extract any sensitive information about the training dataset.

Many pioneering works have studied how to effectively integrate ADMM with DP, e.g., [5–7, 12, 23, 25, 26], which can be classified into two categories in general. The first type of works is to add a noisy term to perturb the objective function in each ADMM iteration using an objective perturbation approach [3]. The second type of works is to perturb the updates of original distributed ADMM

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference'17, July 2017, Washington, DC, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

algorithm via an output perturbation approach [3]. Specifically, as an objective perturbation method, [23] proposed to inject noise to the update of the dual variable to provide DP guarantee, while the total privacy loss over the whole iterative process is not quantified. Further, [25, 26] proposed to perturb the penalty parameter of ADMM and re-utilize the previous iteration’s results to save the privacy loss. These methods also quantify the total privacy loss over the entire process. Moreover, [12] perturbed the augmented Lagrangian with time-varying Gaussian noise and considered a centralized network structure to perform ADMM updates. As output perturbation methods, [6, 7] proposed to perturb the primal variables by Gaussian noise with linearly decaying Gaussian noise to preserve DP and maintain the utility.

However, the guarantee of DP in the above works relies on the premise that the optimality of each local problem can be perfectly attained in each iteration during the whole training procedure, which is seldom seen in practice. Further, the trained models from the above works exhibit severe degradation in terms of the convergence performance and model accuracy, compared to their non-private versions.

In this paper, we propose (Improved) Plausible differentially Private **ADMM** based distributed machine learning algorithm called PP-ADMM and IPP-ADMM, respectively. Instead of requiring each local problem to reach the optimality, PP-ADMM is able to release a noisy approximate solution of the local optimization with Gaussian noise related to the optimization accuracy, while preserving DP.

To further improve the utility, we propose an improved version of PP-ADMM, i.e., IPP-ADMM, by exploiting the sparse vector technique (SVT) to check whether the current approximate solution has enough difference from that of the previous iteration. Moreover, the privacy analysis of our algorithms based on the zero-concentrated DP (zCDP) yields a tight privacy loss bound. We analyze the generalization performance of PP-ADMM. Our salient contributions are summarized as follows.

- To release the shackles of “exact local optimal results” and make ADMM based distributed machine learning achieve DP and plausible, we propose a novel PP-ADMM method by exploiting the inexact solution of the perturbed local optimization over local agent’s private data during each ADMM iteration, while preserving the data privacy.
- We further propose an improved version of PP-ADMM (IPP-ADMM) by employing SVT to evaluate whether the current approximate solution has a big enough difference from that of the previous iteration. If the difference surpasses a pre-defined threshold, the approximate solution with Gaussian noise will be shared with neighbors; otherwise, the current approximate solution will be discarded. By this way, the redundant privacy loss accumulation and the transmissions of “low quality” can be avoided during the ADMM iterative process.
- To best track the privacy loss accumulation, we leverage the serial and parallel composition theorems of the zCDP to theoretically quantify and analyze the overall privacy guarantees of the PP-ADMM and IPP-ADMM algorithms. Moreover, we provide a generalization performance analysis

of PP-ADMM by measuring the number of data samples required to achieve a certain criteria.

- Through extensive experiments on the real-world datasets, we show the superior performance of the proposed algorithms over the state-of-the-art differentially private ADMM algorithms.

The rest of the paper is organized as follows. Section 2 formulates the optimization problem, and describes preliminaries of ADMM and differential privacy. Then, the plausible private robust ADMM algorithm and its corresponding privacy and sample complexity analysis are presented in Section 3. In Section 4, we provide the IPP-ADMM, an improved version of PP-ADMM, and the corresponding privacy analysis. The experimental results on real-world datasets are shown in Section 5. Finally, we draw conclusive remarks in Section 6.

## 2 PROBLEM FORMULATION AND PRELIMINARIES

In this paper, we consider a connected network contains  $N$  agents with node set  $\mathcal{N} = \{1, \dots, N\}$ , and each agent  $i$  has a dataset  $D_i$  with  $D_i = \{(x_i^n, y_i^n)\}_{n=1}^{|D_i|}$ , where  $x_i^n \in \mathcal{X}$  is a feature vector and  $y_i^n \in \mathcal{Y}$  is a label. The communication among agents can be represented by an undirected graph  $G = \{\mathcal{N}, \mathcal{E}\}$ , where  $\mathcal{E}$  denotes the set of communication links between agents. Note that two agents  $i$  and  $j$  can communicate with each other only when they are neighbors, i.e.,  $(i, j) \in \mathcal{E}$ . We also denote the set of neighbors of agent  $i$  as  $\mathcal{B}_i$ . The goal is to cooperatively train a classifier  $\theta \in \mathbb{R}^d$  over the union of all local datasets in a decentralized fashion (i.e., no centralized controller) while keeping the privacy for each data sample, which can be formulated as an Empirical Risk Minimization (ERM) problem.

$$\min_{\theta \in \mathbb{R}^d} \sum_{i=1}^N \frac{1}{|D_i|} \sum_{n=1}^{|D_i|} \mathcal{L}(y_i^n \theta^T x_i^n) + \hat{\lambda} \mathcal{R}(\theta), \quad (1)$$

where  $\mathcal{L}(\cdot) : \mathcal{X} \times \mathcal{Y} \times \mathbb{R}^d \rightarrow \mathbb{R}$  stands for a convex loss function with  $|\mathcal{L}'(\cdot)| \leq 1$  and  $0 < \mathcal{L}''(\cdot) \leq c_1$ ,  $\mathcal{R}(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}$  is a differentiable and 1-strongly convex regularizer to prevent overfitting, and  $\hat{\lambda} \geq 0$  refers to a regularizer parameter that controls the impact of regularizer. We assume that each feature vector  $x_i^n$  is normalized to  $\|x_i^n\|_2 \leq 1$ . Note that the formulations of classification in machine learning like logistic regression, or support vector machines, can also be fallen into the framework of ERM. In order to solve the ERM problem (1) in a decentralized manner, we adopt the simple but efficient optimization method, ADMM. We then in the following subsection review some preliminaries about ADMM algorithm for solving Problem (1).

### 2.1 ADMM

It is easy to see that the ERM problem (1) can be equivalently reformulated as the following consensus form by introducing  $\theta_i$ , that is, the local copy of common classifier  $\theta$  at agent  $i$ .

$$\begin{aligned} \min_{\{\theta_i\}, \{\rho_{ij}\}} \quad & \sum_{i=1}^N f_i(\theta_i) \\ \text{s.t.} \quad & \theta_i = \rho_{ij}, \theta_j = \rho_{ij}, i \in \mathcal{N}, j \in \mathcal{B}_i, \end{aligned} \quad (2)$$

where  $\{\rho_{ij} | i \in \mathcal{N}, j \in \mathcal{B}_i\}$  is a set of slack variables to enforce all local copies are equal to each other, i.e.,  $\theta_1 = \theta_2 = \dots = \theta_N$ , and  $f_i(\theta_i) = \frac{1}{|D_i|} \sum_{n=1}^{|D_i|} \mathcal{L}(y_i^n \theta_i^T x_i^n) + \frac{\lambda}{N} \mathcal{R}(\theta_i)$ . According to Problem (2), each agent  $i$  can minimize local function  $f_i(\theta_i)$  over its own private dataset with respect to  $\theta_i$ , under the consensus constraints. In [25], ADMM is employed to optimize Problem (2) in a decentralized fashion. By defining a dual variable  $\lambda_i$  for agent  $i$ , and introducing the following notion,  $\mathcal{L}_{non}(\theta_i, D_i) = f_i(\theta_i) + (2\lambda_i^t)^T \theta_i + \eta \sum_{j \in \mathcal{B}_i} \|\frac{1}{2}(\theta_i^t + \theta_j^t) - \theta_i\|_2^2$ , ADMM then has the following iterative updates in the  $(t+1)$ -th iteration:

$$\theta_i^{t+1} = \underset{\theta_i}{\operatorname{argmin}} \mathcal{L}_{non}(\theta_i, D_i); \quad (3)$$

$$\lambda_i^{t+1} = \lambda_i^t + \frac{\eta}{2} \sum_{j \in \mathcal{B}_i} (\theta_i^{t+1} - \theta_j^{t+1}), \quad (4)$$

where  $\eta > 0$  is a penalty parameter. Note that the reason why the variable  $\rho_{ij}$  is not appeared in (3) and (4) is that it can be expressed by using the primal variable  $\theta_i$ , as shown in [15]. In the iteration  $t+1$ , each agent  $i \in \mathcal{N}$  updates its local  $\theta_i^{t+1}$  via (3) by using its previous results  $\theta_i^t$  and  $\lambda_i^t$ , and the shared local classifiers  $\theta_j^t$  from its neighbors  $j \in \mathcal{B}_i$ . Next, agent  $i$  broadcasts  $\theta_i^{t+1}$  to all its neighboring agents. After obtaining all of its neighboring computation results, each agent updates the dual variable  $\lambda_i^{t+1}$  through (4).

## 2.2 Differential Privacy

For the privacy-preserving data analysis, the standard privacy metric, Differential privacy (DP) [9, 11], is proposed to measure the privacy risk of each data sample in the dataset, and has already been adopted in many machine learning domains [4, 8, 18, 20, 24]. Basically, under DP framework, privacy protection is guaranteed by limiting the difference of the distribution of the output regardless of the value change of any one sample in the dataset.

**Definition 2.1 (( $\epsilon, \delta$ )-DP [9]).** A randomized mechanism  $\mathcal{M}$  satisfies ( $\epsilon, \delta$ )-DP if for any two neighboring datasets  $D$  and  $\hat{D}$  differing in at most one single data sample, and for any possible output  $o \in \operatorname{Range}(\mathcal{M})$ , we have  $\Pr[\mathcal{M}(D) = o] \leq e^\epsilon \Pr[\mathcal{M}(\hat{D}) = o] + \delta$ .

Here  $\epsilon, \delta$  are privacy loss parameters that indicate the strength of the privacy protection from the mechanism  $\mathcal{M}$ . The privacy protection is stronger if they are smaller. The above privacy definition reduces to pure DP when  $\delta = 0$  and when  $\delta > 0$ , it is referred to as approximate DP. We can achieve pure and approximate DP by utilizing two popular approaches called Laplace and Gaussian Mechanism, both of which share the same form  $\mathcal{M}(D) = \mathcal{M}_q(D) + \mathbf{u}$ , where  $\mathcal{M}_q(D)$  is a query function over dataset  $D$ , and  $\mathbf{u}$  is random noise. We also denote two neighboring datasets  $D$  and  $\hat{D}$  as  $D \sim \hat{D}$ , and denote  $\operatorname{Lap}(\lambda)$  as a zero-mean Laplace distribution with scale parameter  $\lambda$ .

**Definition 2.2 (Laplace Mechanism [9]).** Given any function  $\mathcal{M}_q : \mathcal{D} \rightarrow \mathbb{R}^d$ , the Laplace Mechanism is defined as:  $\mathcal{M}_L(D, q, \epsilon) = \mathcal{M}_q(D) + \mathbf{u}$ , where  $\mathbf{u}$  is drawn from a Laplace distribution  $\operatorname{Lap}(\frac{\Delta_1}{\epsilon})$  with scale parameter proportional to the  $L_1$ -sensitivity of  $\mathcal{M}_q$  given

as  $\Delta_1 = \sup_{D \sim \hat{D}} \|\mathcal{M}_q(D) - \mathcal{M}_q(\hat{D})\|_1$ . Laplace Mechanism preserves  $\epsilon$ -differential privacy.

**Definition 2.3 (Gaussian Mechanism [9]).** Given any function  $\mathcal{M}_q : \mathcal{D} \rightarrow \mathbb{R}^d$ , the Gaussian Mechanism is defined as:  $\mathcal{M}_G(D, q, \epsilon, \delta) = \mathcal{M}_q(D) + \mathbf{u}$ , where  $\mathbf{u}$  is drawn from a Gaussian distribution  $\mathcal{N}(0, \sigma^2 I_d)$  with  $\sigma \geq \frac{\sqrt{2 \ln(1.25/\delta)} \Delta_2}{\epsilon}$ , and  $\Delta_2$  is the  $L_2$ -sensitivity of function  $\mathcal{M}_q$ , i.e.,  $\Delta_2 = \sup_{D \sim \hat{D}} \|\mathcal{M}_q(D) - \mathcal{M}_q(\hat{D})\|_2$ . Gaussian Mechanism provides  $(\epsilon, \delta)$ -differential privacy.

Next, we introduce a generalization of DP, which is called the zero-concentrated DP (zCDP) [1] that uses the Rényi divergence between  $\mathcal{M}(D)$  and  $\mathcal{M}(\hat{D})$ , which can achieve a much tighter privacy loss bound under multiple privacy mechanisms composition.

**Definition 2.4 ( $\rho$ -zCDP [1]).** We say that a randomized algorithm  $\mathcal{M}$  provides  $\rho$ -zCDP, if for all neighboring datasets  $D$  and  $\hat{D}$  and for all  $\tau \in (1, \infty)$ , we have  $D_\tau(\mathcal{M}(D) \| \mathcal{M}(\hat{D})) \leq \rho \tau$ , where  $D_\tau(\mathcal{M}(D) \| \mathcal{M}(\hat{D}))$  is the  $\tau$ -Rényi divergence<sup>1</sup> between the distribution  $\mathcal{M}(D)$  and the distribution  $\mathcal{M}(\hat{D})$ .

The following lemmas show that the Gaussian mechanism satisfies zCDP, the composition theorem of  $\rho$ -zCDP, and the relationship among  $\rho$ -zCDP,  $\epsilon$ -DP, and  $(\epsilon, \delta)$ -DP.

**LEMMA 2.5 ([1]).** *The Gaussian mechanism with noise  $\mathcal{N}(0, \sigma^2)$  satisfies  $\Delta_2^2 / (2\sigma^2)$ -zCDP.*

**LEMMA 2.6 (Serial Composition [1]).** *If randomized mechanisms  $\mathcal{M}_1 : \mathcal{D} \rightarrow \mathcal{R}_1$  and  $\mathcal{M}_2 : \mathcal{D} \rightarrow \mathcal{R}_2$  obey  $\rho_1$ -zCDP and  $\rho_2$ -zCDP, respectively. Then their composition defined as  $\mathcal{M}'' : \mathcal{D} \rightarrow \mathcal{R}_1 \times \mathcal{R}_2$  by  $\mathcal{M}'' = (\mathcal{M}_1, \mathcal{M}_2)$  obeys  $(\rho_1 + \rho_2)$ -zCDP.*

**LEMMA 2.7 (DP to zCDP conversion [1]).** *If a randomized mechanism  $\mathcal{M}$  provides  $\epsilon$ -DP, then  $\mathcal{M}$  is  $\frac{1}{2}\epsilon^2$ -zCDP. Moreover, for  $\mathcal{M}$  to satisfy  $(\epsilon, \delta)$ -DP, it suffices to satisfy  $\rho$ -zCDP with  $\rho = \frac{\epsilon^2}{4 \ln(1/\delta)}$ .*

**LEMMA 2.8 (zCDP to DP conversion [1]).** *If a randomized mechanism  $\mathcal{M}$  obeys  $\rho$ -zCDP, then  $\mathcal{M}$  obeys  $(\rho + 2\sqrt{\rho \ln(1/\delta)}, \delta)$ -DP for all  $0 < \delta < 1$ .*

## 2.3 Sparse Vector Technique

A powerful approach for achieving DP employs the sparse vector technique (SVT) [10], which takes a sequence of queries with bounded sensitivity  $\Delta$  against a fixed sensitive dataset and outputs a vector representing whether each answer to the query exceeds a threshold or not. A unique advantage of SVT is that it can output some query answer without paying additional privacy cost. Specifically, as shown in [14], SVT has the following four steps. (i), We first compute a noisy threshold  $\hat{y}$  by adding a threshold noise  $\operatorname{Lap}(\frac{\Delta}{\epsilon_1})$  to the predefined threshold  $y$ . (ii), We then utilize a noise  $v_i \sim \operatorname{Lap}(\frac{2c\Delta}{\epsilon_2})$  to perturb each query  $q_i$ . (iii), We compare each noisy query answer  $q_i(D) + v_i$  with the noisy threshold  $\hat{y}$  and respond whether it is higher ( $\top$ ) or lower ( $\perp$ ) than the threshold. (iv), This procedure continues until the number of  $\top$ 's in the output meets the predefined bound  $c$ . According to [14], the SVT algorithm satisfies the  $\epsilon$ -DP with  $\epsilon = \epsilon_1 + \epsilon_2$ . In order to analyze the privacy guarantee of SVT under the zCDP framework, we utilize

<sup>1</sup>Definition can be found in [1]

the conversion result in Lemma 2.7. We can see that SVT satisfies  $\frac{1}{2}\epsilon^2$ -zCDP.

### 3 PLAUSIBLE PRIVATE ADMM

In this section, we will present our plausible differentially private (PP-ADMM) by adding Gaussian noise related to the maximum tolerable gradient norm of perturbed objective in each ADMM iteration, which relaxes the requirement of exact optimal solution as shown in [23, 25, 26], to provide differential privacy guarantee of each training data sample during the iterative procedure. We also adopt the privacy framework of zCDP to compute much tighter privacy loss estimation of PP-ADMM. In addition, the generalization performance guarantees of PP-ADMM is provided by measuring the number of data samples at each agent to achieve a specific criteria.

Specifically, in each iteration, we perturb the subproblem (3) with the objective perturbation method the same as used in previous studies [23, 25, 26], where a random linear vector  $(b_{i1})^T \theta_i$  is injected to the objective function, and  $b_{i1}$  is a random vector drawn from a zero mean Gaussian distribution  $\mathcal{N}(0, \sigma_{i1}^2 I_d)$ . Consequently the objective function (3) used to update the primal variable  $\theta_i^{t+1}$  becomes the following modified function:

$$\mathcal{L}_{per}(\theta_i, D_i) = f_i(\theta_i) + (2\lambda_i^t + b_{i1})^T \theta_i + \eta \sum_{j \in \mathcal{B}_i} \|\frac{1}{2}(\theta_i^t + \theta_j^t) - \theta_i\|_2^2 \quad (5)$$

where  $f_i(\theta_i) = \frac{1}{|D_i|} \sum_{n=1}^{|D_i|} \mathcal{L}(y_i^n \theta_i^T x_i^n) + \frac{\lambda}{N} \mathcal{R}(\theta_i)$ . In order to ensure DP guarantee, as pointed out in [23, 25, 26], each agent  $i \in \mathcal{N}$  needs to find the optimal solution  $\hat{\theta}_i^{t+1}$  of the perturbed objective function  $\mathcal{L}_{per}(\theta_i, D_i)$ , i.e.,

$$\hat{\theta}_i^{t+1} = \underset{\theta_i}{\operatorname{argmin}} \mathcal{L}_{per}(\theta_i, D_i). \quad (6)$$

However, the subproblem (6) may not be easy to solve and obtain an optimal solution in a finite time. For instance, if we choose logistic regression as loss function, the subproblem (6) cannot yield an analytical solution due to the complicated form of logistic regression. Especially when the problem dimension or the number of training samples is large, obtaining optimal solution might not be feasible in every iteration.

Thus, we consider obtaining the approximate solution of perturbed objective function  $\mathcal{L}_{per}(\theta_i, D_i)$  to provide privacy guarantees when the optimal solution is not attainable.

Specifically, we approximately solve the perturbed problem until the norm of gradient of  $\mathcal{L}_{per}$  is within a pre-defined threshold  $\beta$ . However, due to the limitations of objective perturbation method [3], releasing this inexact solution leads to the failure of providing DP guarantee. We thus perturb the approximated solution  $\hat{\theta}_i^{t+1}$  with another random noise  $b_{i2}$  from Gaussian distribution  $\mathcal{N}(0, \sigma_{i2}^2 I_d)$ , to "fuzz" the difference between  $\hat{\theta}_i^{t+1}$  and the optimal solution  $\hat{\theta}_i^{t+1}$ . Note that the noise variance  $\sigma_{i2}^2$  has the parameter  $\beta$  about the maximum tolerable gradient norm, which leads to a trade-off between the gradient norm bound and the difficulty of obtaining an approximate solution within the norm bound.

---

#### Algorithm 1 Plausible Private ADMM

---

- 1: **Input:** datasets  $\{D_i\}_{i=1}^N$ ; initial variables  $\theta_i^0 \in \mathbb{R}^d$  and  $\lambda_i^0 = 0_d$ ; step size  $\eta$ ; privacy parameters,  $\epsilon_{i1}, \delta_{i1}, \epsilon_{i3}, \rho_{i2}$ ; Optimizer  $\mathcal{O}(\cdot, \cdot) : \mathcal{F} \times \boldsymbol{\beta} \rightarrow \mathbb{R}^d$  ( $\mathcal{F}$  is the class of objectives, and  $\boldsymbol{\beta}$  is the optimization accuracy, i.e., the gradient norm of objectives); gradient norm threshold  $\beta \in \boldsymbol{\beta}$ .
  - 2: Set  $\epsilon_{i1}, \delta_{i1}, \epsilon_{i3}, \rho_{i2} > 0$  such that  $\epsilon_{i1} > \epsilon_{i3}$ .
  - 3: Set regularizer parameter  $\hat{\lambda} \geq \max_i \frac{2.8N\epsilon_{i1}}{(\epsilon_{i1} - \epsilon_{i3})|D_i|}$ .
  - 4: **for**  $t = 0, \dots, T - 1$  **do**
  - 5:     **for**  $i = 1, \dots, N$  **do**
  - 6:         Generate noise  $b_{i1} \sim \mathcal{N}(0, \sigma_{i1}^2 I_d)$  with  $\sigma_{i1} = 2\sqrt{2 \ln(1.25/\delta_{i1})}/(|D_i|\epsilon_{i3})$ .
  - 7:         Construct the perturbed objective function  $\mathcal{L}_{per}(\theta_i, D_i)$  according to (5).
  - 8:         Compute an approximate solution  $\hat{\theta}_i^{t+1}$ :  $\hat{\theta}_i^{t+1} = \mathcal{O}(\mathcal{L}_{per}(\theta_i, D_i), \beta)$ .
  - 9:         Generate noise  $b_{i2} \sim \mathcal{N}(0, \sigma_{i2}^2 I_d)$  with  $\sigma_{i2} = \beta / [\sqrt{2\rho_{i2}}(\frac{\hat{\lambda}}{N} + 2\eta|\mathcal{B}_i|)]$ .
  - 10:         Perturb  $\hat{\theta}_i^{t+1}$ :  $\theta_i^{t+1} = \hat{\theta}_i^{t+1} + b_{i2}$ .
  - 11:     **end for**
  - 12:     **for**  $i = 1, \dots, N$  **do**
  - 13:         Broadcast  $\theta_i^{t+1}$  to all neighbors  $j \in \mathcal{B}_i$ .
  - 14:     **end for**
  - 15:     **for**  $i = 1, \dots, N$  **do**
  - 16:         Update local dual variables  $\lambda_i^{t+1}$  from  $\lambda_i^{t+1} = \lambda_i^t + \frac{\eta}{2} \sum_{j \in \mathcal{B}_i} (\theta_i^{t+1} - \theta_j^{t+1})$ .
  - 17:     **end for**
  - 18: **end for**
- 

The key steps of PP-ADMM algorithm are summarized in Algorithm 1. The privacy parameters  $(\epsilon_{i1}, \delta_{i1})$  are used to perturb the objective function while the parameter  $\rho_{i2}$  being used to perturb the approximate solution. Moreover, the parameter  $\epsilon_{i3}$ , a portion of  $\epsilon_{i1}$ , is used to scale the noise injected to the objective function, and the remaining privacy budget  $(\epsilon_{i1} - \epsilon_{i3})$  is allocated to setting the regularizer parameter. Notice that we also define an **Optimizer**  $\mathcal{O}(\cdot, \cdot) : \mathcal{F} \times \boldsymbol{\beta} \rightarrow \mathbb{R}^d$ , where  $\mathcal{F}$  is the class of objectives, and  $\boldsymbol{\beta}$  is the optimization accuracy, i.e., the gradient norm of objectives. Each agent  $i$  then constructs the perturbed function  $\mathcal{L}_{per}(\theta_i, D_i)$  with a Gaussian random vector  $b_{i1}$  and finds an inexact solution  $\hat{\theta}_i^{t+1}$  where the norm of gradient is lower than  $\beta$ , i.e.,  $\hat{\theta}_i^{t+1} = \mathcal{O}(\mathcal{L}_{per}(\theta_i, D_i), \beta)$ . After that each agent  $i$  generates a random Gaussian noise  $b_{i2}$  and transmits  $\theta_i^{t+1} = \hat{\theta}_i^{t+1} + b_{i2}$  to its neighbors  $j \in \mathcal{B}_i$ . Finally, each agent updates the local dual variables  $\lambda_i^{t+1}$  via  $\lambda_i^{t+1} = \lambda_i^t + \frac{\eta}{2} \sum_{j \in \mathcal{B}_i} (\theta_i^{t+1} - \theta_j^{t+1})$ .

### 3.1 Privacy Analysis

Here, we provide the privacy guarantee of PP-ADMM (Algorithm 1) in the following two theorems. Due to the limited space, we only provide a proof idea of Theorem 3.1, and the detailed proof can be found in Appendix.

**THEOREM 3.1.** *The PP-ADMM in Algorithm 1 satisfies  $\rho_i$ -zCDP for each agent  $i$  with  $\rho_i = T(\rho_{i1} + \rho_{i2})$ , where  $\rho_{i1} = \epsilon_{i1}^2 / (4 \ln(1/\delta_{i1}))$ , and  $\rho_{i2} > 0$  is the privacy budget for perturbing the approximate solution.*

**PROOF SKETCH.** For achieving  $\rho_i$ -zCDP for each agent  $i$  at  $t + 1$  iteration in Algorithm 1, we first divide the output of  $t + 1$  iteration into two parts. The first part is to obtain the optimal solution  $\tilde{\theta}_i^{t+1}$  of the perturbed objective function  $\mathcal{L}_{per}(\theta_i, D_i)$ , and the second part is to obtain the approximate solution with Gaussian noise  $\theta_i^{t+1}$ . We then show obtaining the optimal solution  $\tilde{\theta}_i^{t+1}$  provides  $\rho_{i1}$ -zCDP with  $\rho_{i1} = \epsilon_{i1}^2 / (4 \ln(1/\delta_{i1}))$  for the first part, and releasing an approximate solution in the second part is  $\rho_{i2}$ -zCDP. By using the composition of zCDP in Lemma 2.6, we can get releasing the perturbed primal variable  $\theta_i^{t+1}$  at  $t + 1$  iteration provides  $(\rho_{i1} + \rho_{i2})$ -zCDP. Considering  $T$  iterations, the total privacy loss for each agent  $i$  is bounded by  $\rho_i = T(\rho_{i1} + \rho_{i2})$ .  $\square$

We then give the following parallel composition theorem of  $\rho$ -zCDP to provides a together characterization of total privacy loss for distributed algorithms.

**LEMMA 3.2 (PARALLEL COMPOSITION [21]).** *Suppose that a mechanism  $\mathcal{M}$  consists of a sequence of  $k$  adaptive mechanism  $\mathcal{M}_1, \dots, \mathcal{M}_k$  where each  $\mathcal{M}_i : \prod_{j=1}^{i-1} \mathcal{R}_j \times \mathcal{D} \rightarrow \mathcal{R}_i$  and  $\mathcal{M}_i$  satisfies  $\rho_i$ -zCDP. Let  $\mathbb{D}_1, \mathbb{D}_2, \dots, \mathbb{D}_k$  be the result of a randomized partition of the input domain  $\mathbb{D}$ . The mechanism  $\mathcal{M}(D) = (\mathcal{M}_1(D \cap \mathbb{D}_1), \dots, \mathcal{M}_k(D \cap \mathbb{D}_k))$  satisfies  $(\max_i \rho_i)$ -zCDP.*

Based on Lemma 3.2, we can directly obtain the total privacy loss of PP-ADMM given as follows.

**THEOREM 3.3.** *The PP-ADMM in Algorithm 1 satisfies  $\rho$ -zCDP with  $\rho = \max_i \rho_i$  and satisfies  $(\epsilon, \delta)$ -DP with  $\epsilon = \rho + 2\sqrt{\rho \ln(1/\delta)}$ .*

## 3.2 Sample Complexity Analysis

We next measure the generalization performance of PP-ADMM by focusing on the ERM problem given in Section 2. We also assume that data samples of each agent  $i$  are drawn from a data distribution  $\mathcal{P}$ . The expected loss of classifier  $\theta_i^t$  at iteration  $t$  is defined as

$$\mathbb{L}(\theta_i^t) = \mathbb{E}_{(x,y) \sim \mathcal{P}} \left( \mathcal{L}(y(\theta_i^t)^T x) \right).$$

Following the similar analysis in [3, 23], we first introduce a reference classifier  $\theta_{ref}$  with expected loss  $\mathbb{L}(\theta_{ref})$ , and we then measure the generalization performance using the number of samples  $D_i$  required at each agent to achieve  $\mathbb{L}(\theta_i^t) \leq \mathbb{L}(\theta_{ref}) + a_{acc}$ , where  $a_{acc}$  is the generalization error.

**3.2.1 PP-ADMM without Noise.** Here, we consider the learning performance at all iterations rather than only the final output. Let the intermediate updated classifier  $\hat{\theta}_{i,non}^{t+1}$  at iteration  $t + 1$  be  $\hat{\theta}_{i,non}^{t+1} = \mathcal{O}(\mathcal{L}_{non}(\theta_i, D_i), \beta)$ . Note that  $\{\hat{\theta}_{i,non}^{t+1}\}$  is a sequence of non-private classifier without adding perturbations. Let  $\theta^* = \operatorname{argmin}_{\theta_i} f_i(\theta_i, D_i)$  be the optimal output of PP-ADMM without Noise. The sequence  $\{\hat{\theta}_{i,non}^{t+1}\}$  is bounded and  $\hat{\theta}_{i,non}^{t+1}$  converges to  $\theta^*$  as  $t \rightarrow \infty$ . Therefore, there exists a constant  $\Delta_{i,non}^{t+1}$  at iteration  $t + 1$  such that

$\mathbb{L}(\hat{\theta}_{i,non}^{t+1}) \leq \mathbb{L}(\theta^*) + \Delta_{i,non}^{t+1}$ . We then have the following result, and the detailed proof can be found in supplemental material.

**THEOREM 3.4.** *Consider a regularized ERM problem with  $\mathcal{R}(\theta) = \frac{1}{2} \|\theta\|_2^2$ , and let  $\theta_{ref}$  be the reference classifier for all agents and  $\{\hat{\theta}_{i,non}^{t+1}\}$  be a sequence of outputs of PP-ADMM without adding noise. If the number of samples at agent  $i$  satisfies,*

$$|D_i| \geq V \max_t \left\{ \frac{\log(1/\xi)}{\frac{(a_{acc} - \Delta_{i,non}^{t+1})^2}{2\|\theta_{ref}\|_2^2} - (1+a)\beta} \right\}$$

for some constant  $V$ , then  $\hat{\theta}_{i,non}^{t+1}$  satisfies

$$\Pr \left[ \mathbb{L}(\hat{\theta}_{i,non}^{t+1}) \leq \mathbb{L}(\theta_{ref}) + a_{acc} \right] \geq 1 - \xi$$

with  $a_{acc} \geq \Delta_{i,non}^{t+1}$ .

**REMARK 1.** *As we can see from Theorem 3.4, the number of data samples  $|D_i|$  relies on the  $l_2$ -norm of reference classifier  $\|\theta_{ref}\|_2^2$  and the parameter  $\beta$  that bounds the optimization accuracy of the non-private intermediate classifier. The results demonstrate that if  $|D_i|$  satisfies  $|D_i| \geq V \max_t \left\{ \frac{\log(1/\xi)}{\frac{(a_{acc} - \Delta_{i,non}^{t+1})^2}{2\|\theta_{ref}\|_2^2} - (1+a)\beta} \right\}$ , each agent's non-*

*private intermediate classifier will have an additional error less than  $a_{acc}$  compared to any classifier with  $\|\theta_{ref}\|_2^2$ . Moreover, if  $\beta = 0$ , the result reduces to  $|D_i| \geq V \max_t \left\{ \frac{2\|\theta_{ref}\|_2^2 \log(1/\xi)}{(a_{acc} - \Delta_{i,non}^{t+1})^2} \right\}$ , the same as given in [23], which shows that the lower optimization accuracy of the non-private intermediate classifier, the more samples required to achieve the same accuracy.*

**3.2.2 PP-ADMM.** We then show the sample complexity of the PP-ADMM algorithm. Similar to the analysis in PP-ADMM without noise, we also consider bounding the generalization error of the intermediate classifier  $\theta_i^{t+1}$  of each agent  $i$  at all iterations. In order to compare the private classifier  $\theta_i^{t+1}$  with a reference classifier  $\theta_{ref}$ , we follow the same strategy used in [23]. We define a new optimization function  $f_i^{new}(\theta_i, D_i) = f_i(\theta_i, D_i) + b_{i1}^T \theta_i$  and then solving PP-ADMM algorithm is equivalent to solving a new optimization problem, where each agent  $i$ 's performs local minimization to get  $\theta_i^{t+1} = \mathcal{O}(f_i^{new}(\theta_i, D_i), \beta) + b_{i2}$ . The sequence of outputs  $\{\theta_i^{t+1}\}$  is bounded and  $\theta_i^{t+1}$  converges to a fixed point  $\theta_{new}^*$  as  $t \rightarrow \infty$ . Thus, there exists a constant  $\Delta_{i,new}^{t+1}$  at  $t + 1$  iteration, such that  $\mathbb{L}(\theta_i^{t+1}) \leq \mathbb{L}(\theta_{new}^*) + \Delta_{i,new}^{t+1}$ . We then give the following result, and the detailed proof can be found in supplemental material.

**THEOREM 3.5.** *Consider a regularized ERM problem with  $\mathcal{R}(\theta) = \frac{1}{2} \|\theta\|_2^2$ , and let  $\theta_{ref}$  be the reference classifier for all agents and  $\{\theta_i^{t+1}\}$  be a sequence of outputs of PP-ADMM. If the number of samples at agent  $i$  satisfies, for some constant  $V$ ,*

$$|D_i| \geq V \max_t \left\{ \frac{\log(1/\xi)}{\frac{(a_{acc} - \Delta_{i,new}^{t+1})^2}{2\|\theta_{ref}\|_2^2} - (1+a)(\beta + \mathcal{H})} \right\}$$

with  $\mathcal{H} = \frac{\sigma_{i2}(a_{acc} - \Delta_{i,new}^{t+1})\sqrt{2d \log \frac{1}{\xi}}}{\|\theta_{ref}\|_2^2} + 2\sigma_{i1}^2 d \log \frac{1}{\xi}$ , then  $\hat{\theta}_{i,new}^{t+1}$  satisfies

$$\Pr \left[ \mathbb{L}(\hat{\theta}_{i,new}^{t+1}) \leq \mathbb{L}(\theta_{ref}) + a_{acc} \right] \geq 1 - 3\xi$$

with  $a_{acc} \geq \Delta_{i,new}^{t+1}$ .

REMARK 2. Compared to Theorem 3.4, we can see that in Theorem 3.5, the privacy constraints impose an additional term  $\mathcal{H}$  with  $\mathcal{H} = \sigma_{i2}(a_{acc} - \Delta_{i,new}^{t+1})\sqrt{2d \log \frac{1}{\xi}} / \|\theta_{ref}\|_2^2 + 2\sigma_{i1}^2 d \log \frac{1}{\xi}$ . If both noise variances  $\sigma_{i1}$  and  $\sigma_{i2}$  are equal to zero, the number of required samples  $|D_i|$  will reduce to the same result shown in Theorem 3.4. Moreover, the additional term  $\mathcal{H}$  demonstrates that the higher dimension of features, the more added noise to achieve the same accuracy requires more data samples.

#### 4 IMPROVED PLAUSIBLE PRIVATE ADMM

In this section, we present an improved version of PP-ADMM algorithm called Improved Plausible Private ADMM (IPP-ADMM) by leveraging sparse vector technique (SVT) to improve the performance and reduce the communication cost of PP-ADMM. Compared with current differentially private ADMM algorithms [23, 25, 26], although the proposed PP-ADMM algorithm can ensure DP guarantee without requiring the optimal solution during each ADMM iteration, the primal variable is updated using the local data in every iteration and frequently broadcasted to neighboring agents, which leads to the privacy loss unavoidably accumulating over the iterations, and compromise the accuracy during the whole training procedure.

Hence, we adopt SVT that can output some local computational results without paying any privacy budget, to check whether current approximate solution has a big enough difference from that of previous iteration,

where the difference is quantified by a quality function,  $f_i(\theta_i^t) - f_i(\hat{\theta}_i^{t+1})$ , based on the change of the values of local function over the primal variable from previous iteration and current approximate solution. If a sufficient level of difference  $\alpha$  is achieved, each agent transmits the current approximate solution with Gaussian noise to its neighbors. Intuitively, if the difference between the current approximate solution  $\hat{\theta}_i^{t+1}$  and the previously transmitted  $\theta_i^t$  is small, then using either one does not help the convergence of the iterative process, which leads to reducing the communication cost.

However, one difficulty in using SVT is that there is no known priori bound on query (i.e., the quality function)  $f_i(\theta_i^t) - f_i(\hat{\theta}_i)$ . To bound the sensitivity of  $f_i(\theta_i^t) - f_i(\hat{\theta}_i)$ , we apply the clipping method to clipping the loss function  $\mathcal{L}(\cdot)$ . Given a fixed clipping threshold  $C_{loss}$ , we compute the value of loss function  $\mathcal{L}(\cdot)$  on each local data sample, clip the values at most  $C_{loss}$ , and compute the value of  $f_i(\theta_i^t) - f_i(\hat{\theta}_i)$  based on the clipped values. Note that we denote this loss function clipping procedure as Clip.

The complete procedure of IPP-ADMM algorithm for a single agent is shown in Algorithm 2. The privacy parameters  $\epsilon_1$  and  $\epsilon_2$  are allocated to perturb the quality function and threshold  $\alpha$ , respectively. In each iteration, each agent  $i$  first constructs the perturbed function  $\mathcal{L}_{per}(\theta_i, D_i)$  with a Gaussian random vector  $b_{i1}$  and finds an inexact solution  $\hat{\theta}_i^{t+1}$ , where the norm of gradient is lower than  $\beta$ , i.e.,  $\hat{\theta}_i^{t+1} = \mathcal{O}(\mathcal{L}_{per}(\theta_i, D_i), \beta)$ . Then each agent apply the clipping method Clip to clip the quality function  $f_i(\theta_i^t) - f_i(\hat{\theta}_i^{t+1})$  with a clipping threshold  $C_{loss}$  to limit the sensitivity of quality function. Further, each agent uses SVT to check whether

---

#### Algorithm 2 Improved Plausible Private ADMM Run by Agent $i$

---

- 1: **Input:** dataset  $D_i$ ; initial variables  $\theta_i^0 \in \mathbb{R}^d$  and  $\lambda_i^0 = 0_d$ ; threshold,  $\alpha$ ; Maximum number of primal variables that can be broadcasted,  $c$ ; loss function clipping threshold  $C_{loss}$ ; step size  $\eta$ ; privacy parameters,  $\epsilon_{i1}, \delta_{i1}, \epsilon_{i3}, \rho_{i2}, \epsilon_1, \epsilon_2$ ; Optimizer  $\mathcal{O}(\cdot, \cdot) : \mathcal{F} \times \boldsymbol{\beta} \rightarrow \mathbb{R}^d$  ( $\mathcal{F}$  is the class of objectives, and  $\boldsymbol{\beta}$  is the optimization accuracy, i.e., the gradient norm of objectives); gradient norm threshold  $\beta \in \boldsymbol{\beta}$ .
  - 2: Set  $\epsilon_{i1}, \delta_{i1}, \epsilon_{i3}, \rho_{i2}, \epsilon_1, \epsilon_2 > 0$  such that  $\epsilon_{i1} > \epsilon_{i3}$ .
  - 3: Set regularizer parameter  $\hat{\lambda} \geq \max_i \frac{2.8Nc_1}{(\epsilon_{i1} - \epsilon_{i3})|D_i|}$ .
  - 4:  $count_i = 0$ .
  - 5: **for**  $t = 0, \dots, T - 1$  **do**
  - 6:     Generate noise  $b_{i1} \sim \mathcal{N}(0, \sigma_{i1}^2 I_d)$  with  $\sigma_{i1} = 2\sqrt{2 \ln(1.25/\delta_{i1})} / (|D_i| \epsilon_{i1})$ .
  - 7:     Construct the perturbed objective function  $\mathcal{L}_{per}(\theta_i, D_i)$  according to (5).
  - 8:     Compute an approximate solution  $\hat{\theta}_i^{t+1}$ :  $\hat{\theta}_i^{t+1} = \mathcal{O}(\mathcal{L}_{per}(\theta_i, D_i), \beta)$ .
  - 9:     **if**  $\text{Clip} \left[ f_i(\theta_i^t) - f_i(\hat{\theta}_i^{t+1}) \right] + \text{Lap} \left( \frac{4cC_{loss}}{\epsilon_2} \right) \geq \alpha + \text{Lap} \left( \frac{2cC_{loss}}{\epsilon_1} \right)$  **then**
  - 10:          $count_i = count_i + 1$ , **Abort** if  $count_i > c$ .
  - 11:         Generate noise  $b_{i2} \sim \mathcal{N}(0, \sigma_{i2}^2 I_d)$  with  $\sigma_{i2} = \beta / [\sqrt{2\rho_{i2}}(\frac{\hat{\lambda}}{N} + 2\eta|\mathcal{B}_i|)]$ .
  - 12:         Perturb  $\hat{\theta}_i^{t+1}$ :  $\theta_i^{t+1} = \hat{\theta}_i^{t+1} + b_{i2}$ .
  - 13:         Broadcast  $\theta_i^{t+1}$  to all neighbors  $j \in \mathcal{B}_i$ .
  - 14:     **else**
  - 15:         Let  $\theta_i^{t+1} = \theta_i^t$ .
  - 16:     **end if**
  - 17:     **if**  $\theta_j^{t+1}$  is not received from neighbor  $j \in \mathcal{B}_i$  **then**
  - 18:         Replace  $\theta_j^{t+1}$  with  $\theta_j^t$ .
  - 19:     **else**
  - 20:         Keep  $\theta_j^{t+1}$ .
  - 21:     **end if**
  - 22:     Update local dual variables  $\lambda_i^{t+1}$  from  $\lambda_i^{t+1} = \lambda_i^t + \frac{\eta}{2} \sum_{j \in \mathcal{B}_i} (\theta_i^{t+1} - \theta_j^{t+1})$ .
  - 23: **end for**
- 

the difference between the approximate solution  $\hat{\theta}_i^{t+1}$  and  $\theta_i^t$  is below a noisy threshold  $\hat{\alpha} = \alpha + \text{Lap} \left( \frac{2cC_{loss}}{\epsilon_1} \right)$  via a noisy quality function,  $\text{Clip} \left[ f_i(\theta_i^t) - f_i(\hat{\theta}_i^{t+1}) \right] + \text{Lap} \left( \frac{4cC_{loss}}{\epsilon_2} \right)$ . If yes, then agent  $i$  does not transmit any computational results and let  $\theta_i^{t+1} = \theta_i^t$ ; otherwise, each agent  $i$  generates a random noise  $b_{i2} \sim \mathcal{N}(0, \sigma_{i2}^2 I_d)$  with  $\sigma_{i2} = \beta / \sqrt{2\rho_{i2}}(\frac{\hat{\lambda}}{N} + 2\eta|\mathcal{B}_i|)$ , and transmits  $\theta_i^{t+1} = \hat{\theta}_i^{t+1} + b_{i2}$  to its neighbors. Moreover, each agent maintains a counter  $count_i$  to bound the total number of broadcasts of primal variables during the whole interactive process. If a predefined transmission number  $c(c \leq T)$  is exceeded, agent  $i$  stops transmitting anything even when the condition in Line 7 is satisfied. Hence, if agent  $i$  does not receive  $\theta_j^{t+1}$  from any neighbor  $j \in \mathcal{B}_i$ , then lets  $\theta_j^{t+1} = \theta_j^t$ .

Finally, each agent updates the local dual variables  $\lambda_i^{t+1}$  via  $\lambda_i^{t+1} = \lambda_i^t + \frac{\eta}{2} \sum_{j \in \mathcal{B}_i} (\theta_j^{t+1} - \theta_j^t)$ .

## 4.1 Privacy Analysis

We provide the privacy guarantee of IPP-ADMM (Algorithm 2) in following theorem.

**THEOREM 4.1.** *The IPP-ADMM in Algorithm 2 satisfies  $\rho'_i$ -zCDP for each agent  $i$  with  $\rho'_i = \rho'_1 + \mathbf{c}(\rho_{i1} + \rho_{i2})$ , where  $\rho'_1 = \frac{(\epsilon_1 + \epsilon_2)^2}{2}$ ,  $\rho_{i1} = \epsilon_{i1}^2 / (4 \ln(1/\delta_{i1}))$ ,  $\rho_{i2} > 0$  is the privacy budget for perturbing the approximate solution, and  $\mathbf{c}$  ( $\mathbf{c} < T$ ) is the maximum number of primal variables that can be broadcasted. Moreover, the total privacy guarantee of IPP-ADMM is  $\rho'$ -zCDP with  $\rho' = \max_i \rho'_i$ .*

**PROOF.** For achieving  $\rho'_i$ -zCDP for each agent  $i$  in Algorithm 2, we first divide the procedure of the algorithm into two parts. The first part is using SVT to compare the noisy threshold and the perturbed query answer (i.e., the value of quality function) to check the quality of the approximate solution obtained in Step 7 of the Algorithm 1. The second part is to share the approximate solution with Gaussian noise, whose value is above the threshold. We prove that DP mechanism used in the first part provides  $\rho'_1$ -zCDP (shown in Lemma 4.3). Moreover, at each iteration, the privacy budget spending on releasing an approximate solution in the second part is  $(\rho_{i1} + \rho_{i2})$ -zCDP (shown in Theorem 3.1). Then, using the composition of zCDP in Lemma 2.6, we obtain the privacy guarantee of IPP-ADMM for each agent  $i$  is  $\rho_i = \rho_1 + \mathbf{c}(\rho_{i1} + \rho_{i2})$  by considering  $\mathbf{c}$  times of broadcasting primal variables. Lastly, we get a total privacy guarantee of IPP-ADMM, i.e.,  $\rho'$ -zCDP with  $\rho' = \max_i \rho'_i$  by adopting the parallel composition in Lemma 3.2.  $\square$

Before presenting the privacy guarantee of the first part, i.e., compare the noisy threshold and the perturbed query answer to check the quality of the approximate solution, we first give the sensitivity of the clipped quality function as follows.

**LEMMA 4.2.** *Given a clipping threshold  $C_{Loss}$  of the loss function  $\mathcal{L}(\cdot)$ , the sensitivity of quality function  $f_i(\theta_i^t) - f_i(\hat{\theta}_i^{t+1})$  is at most  $2C_{Loss}$ , where  $f_i(\theta_i) = \frac{1}{|D_i|} \sum_{n=1}^{|D_i|} \mathcal{L}(y_i^n \theta_i^T x_i^n) + \frac{\lambda}{N} \mathcal{R}(\theta_i)$ .*

**PROOF.** Fix a pair of adjacent datasets  $D_i$  and  $\hat{D}_i$  and we also assume that only the first data point in  $D_i$  and  $\hat{D}_i$  are different, i.e.,  $(x_i^1, y_i^1)$  and  $(\hat{x}_i^1, \hat{y}_i^1)$ . According to the definition of  $L_1$ -sensitivity, we have

$$\begin{aligned} \Delta_f &= \|f_i(\theta_i^t, D_i) - f_i(\hat{\theta}_i^{t+1}, D_i) - f_i(\theta_i^t, \hat{D}_i) + f_i(\hat{\theta}_i^{t+1}, \hat{D}_i)\|_1 \\ &= \|\mathcal{L}(y_i^1(\theta_i^t)^T x_i^1) - \mathcal{L}(\hat{y}_i^1(\theta_i^t)^T \hat{x}_i^1) \\ &\quad - (\mathcal{L}(y_i^1(\theta_i^{t+1})^T x_i^1) - \mathcal{L}(\hat{y}_i^1(\theta_i^{t+1})^T \hat{x}_i^1))\|_1 \\ &\leq \|\mathcal{L}(y_i^1(\theta_i^t)^T x_i^1) - \mathcal{L}(\hat{y}_i^1(\theta_i^t)^T \hat{x}_i^1)\|_1 \\ &\quad + \|\mathcal{L}(y_i^1(\theta_i^{t+1})^T x_i^1) - \mathcal{L}(\hat{y}_i^1(\theta_i^{t+1})^T \hat{x}_i^1)\|_1 \\ &\leq 2C_{Loss}. \end{aligned} \quad \square$$

Then we show the privacy guarantee of the first part in the following lemma.

**LEMMA 4.3.** *Given the maximum number of primal variables that we can broadcast,  $\mathbf{c}$ , using SVT to check whether the approximate solution is above the threshold  $\alpha$  provides  $\rho_1$ -zCDP with  $\rho_1 = \frac{(\epsilon_1 + \epsilon_2)^2}{2}$ .*

**PROOF.** During the whole training process, each agent will receive a stream of queries (i.e., a stream of clipped quality functions  $\text{Clip}[f_i(\theta_i^t) - f_i(\hat{\theta}_i^{t+1})]$ ) with sensitivity  $2C_{Loss}$  and compare them with a noisy threshold  $\alpha + \text{Lap}(\frac{2\mathbf{c}C_{Loss}}{\epsilon_1})$ . According to Theorem 1 in [14], this procedure satisfies  $(\epsilon_1 + \epsilon_2)$ -DP and by Lemma 2.7, it also satisfies  $\frac{(\epsilon_1 + \epsilon_2)^2}{2}$ -zCDP.  $\square$

## 5 NUMERICAL EXPERIMENTS

**Datasets.** Experiments are performed on three benchmark datasets<sup>2</sup>: Adult, US, and Brazil. Adult has 48,842 data samples and 41 features, and the label is to predict whether an annual income is more than \$50k or not. US has 40,000 records and 58 features, and the label is to predict whether the annual income of an individual is more than \$25k. BR has 38,000 samples and 53 features, and the goal is to predict whether the monthly income of an individual is more than \$300.

**Data preprocessing.** We consider the same preprocessing procedure as the method used in [25]. We first normalize each attribute so that the maximum attribute value is 1, and normalize each sample so its  $L_2$ -norm at most 1. As for the label column, we also map it to  $\{-1, 1\}$ . In each simulation, we randomly sample 35,000 records for training and divide them into  $N$  parties, and thus each party includes  $35000/N$  data samples (i.e.,  $|D_i| = 35000/N$ ). We denote the rest of the data records as testing data.

**Baselines.** We compare our proposed algorithms against four baseline algorithms: (i) DVP [23], is a dual variable perturbation method, where the dual variable of each agent at each ADMM iteration is perturbed by Gamma noise. (ii) M-ADMM [25], is a penalty perturbation approach, where each agent's penalty variable is perturbed by Gamma noise at each ADMM iteration. (iii) R-ADMM [26], is based on the penalty approach and the re-utilization of previous iteration's results to save the privacy loss. (iv) Non-private (decentralized ADMM without adding noise). Note that the privacy guarantees of DVP, M-ADMM, and R-ADMM hold only when the optimal solution of the perturbed subproblem is obtained in each iteration. In order to have a fair comparison, we adopt the Newton solver to obtain the optimal solution in each iteration. Notice that we also provide sharpened and tight privacy loss of above private ADMM algorithms under the privacy framework of zCDP.

**Setup.** We adopt logistic loss  $\mathcal{L}(y_i^n \theta_i^T x_i^n) = \log(1 + \exp(-y_i^n \theta_i^T x_i^n))$  as loss function, and the derivative  $\mathcal{L}'(\cdot)$  is bounded with  $|\mathcal{L}'(\cdot)| \leq 1$  and  $c_1$ -Lipschitz with  $c_1 = 1/4$ . We also let  $\mathcal{R}(\theta_i) = \frac{1}{2} \|\theta_i\|_2^2$ . We evaluate the accuracy by classification error rate over the testing set, defined as  $\text{Error rate} = \frac{\text{Number of incorrect predictions}}{\text{Total number of predictions made}}$  and the convergence of algorithms by the average loss over the training samples, given by  $\mathcal{L}_t = \frac{1}{N} \sum_{i=1}^N \frac{1}{|D_i|} \sum_{n=1}^{|D_i|} \mathcal{L}(y_i^n(\theta_i^t)^T x_i^n)$ . We also report the mean and standard deviation of the average loss. The smaller the average loss, the higher accuracy.

<sup>2</sup><http://archive.ics.uci.edu/ml/datasets/Adult>, <http://international.ipums.org>

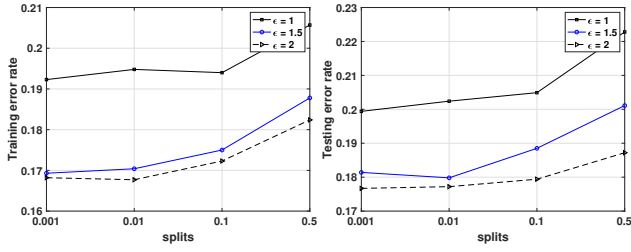


Figure 1: Effects of privacy budget splitting

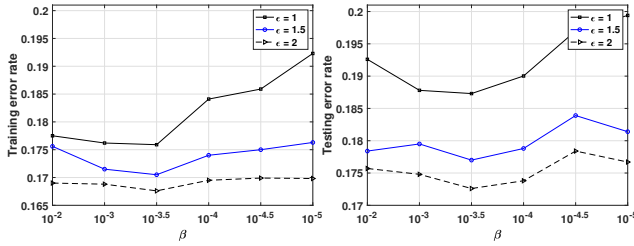


Figure 2: Effects of optimization accuracy  $\beta$

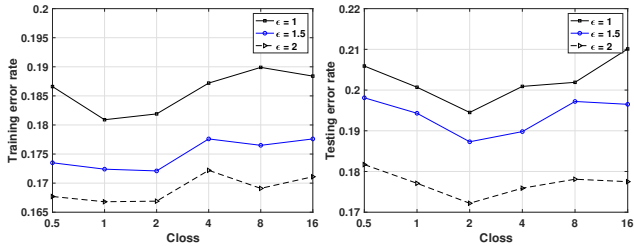


Figure 3: Effects of clipping threshold  $C_{loss}$

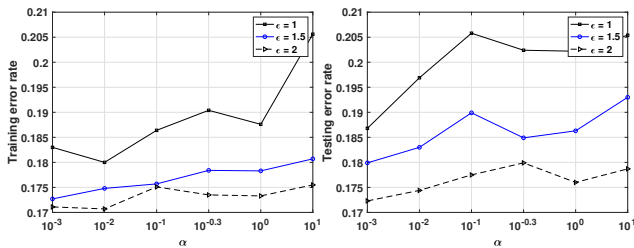


Figure 4: Effects of  $\alpha$

**Parameter settings.** We consider a randomly generated undirected network with  $N = 5$  agents and we fix the step size  $\eta = 0.5$  and the total iteration number  $T = 30$ . We also consider the maximum number of primal variables that can be shared,  $c = 15$ . Moreover, to maximize the utility of SVT, we follow the ratio between  $\epsilon_1$  and  $\epsilon_2$  in [14], i.e.,  $\epsilon_1 : \epsilon_2 = 1 : (2c)^{\frac{2}{3}}$ . In all experiments, we set  $\delta = 10^{-4}$ , and  $\epsilon = \{0.5, 1, 1.5, 2, 10\}$ .

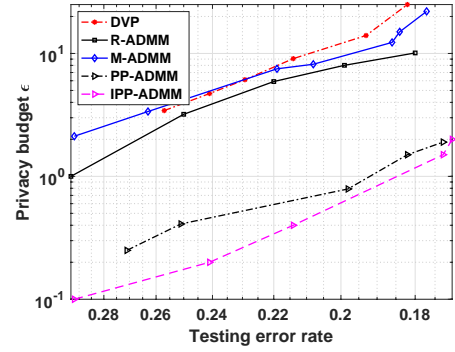


Figure 5: Trade-off between classification error rate and privacy on Adult dataset

**Impacts of parameters.** In this set of experiments on the Adult dataset, we present the effects of privacy budgets splitting and optimization accuracy (i.e., gradient norm threshold)  $\beta$  on the performance of PP-ADMM, and the loss clipping threshold  $C_{loss}$  and the quality function significance threshold  $\alpha$  on the performance of IPP-ADMM. Specifically, we adjust different parameter settings separately, while keeping the rest constant to represent their impacts on training and testing accuracy.

For the privacy budgets splitting of PP-ADMM, we first convert the overall privacy budget parameters  $(\epsilon, \delta)$  to  $\rho_{total} = \frac{\epsilon^2}{4 \ln(1/\delta)}$ . We set  $\rho_{i1} = \frac{\rho_{total}}{T} \cdot (1 - splits)$  and  $\rho_{i2} = \frac{\rho_{total}}{T} \cdot splits$ , where  $splits$  denotes the fraction of  $\rho_{total}$  allocated to  $\rho_{i2}$ . By tuning  $splits$ , we can find the good trade-off between the privacy budget for perturbing the objective and perturbing the approximate solution. In addition, we compute  $\epsilon_{i1} = \rho_{i1} + 2\sqrt{\rho_{i1} \ln(1/\delta_{i1})}$  with  $\delta_{i1} = 10^{-4}$ , and set  $\epsilon_{i3} = 0.99 \cdot \epsilon_{i1}$  to dedicate most of the budget to reduce the amount of noise for perturbing the objective and increase the influence of regularization. Figure 1 shows the effects of privacy budget splitting on the performance of PP-ADMM by setting  $\beta = 10^{-6}$ . As  $splits$  decreases, i.e., allocating less privacy budgets for perturbing the approximate solution, it yields better training and testing accuracy. Thus, we set  $splits = 0.001$  to achieve a good trade-off between amount of noise added to the objective and approximate solution.

Figure 2 shows how classification accuracy changes with varying values of  $\beta$  and fixing  $splits = 0.001$ . The parameter  $\beta$  controls the optimization accuracy of each iteration of PP-ADMM training process and the amount of noise for perturbing the approximate solution. As it can be observed from the figure, due to randomness of objective introduced by the random noise, when  $\beta$  is too small, solving the noisy objective perfectly in each iteration may not help the final performance. Conversely, when  $\beta$  is too large, large amount of noise is added to perturb the approximate solution, which also leads to performance degradation. In our experiments, we thus fix  $\beta = 10^{-3.5}$  that achieves lowest training/testing error rate.

The IPP-ADMM algorithm has two threshold parameters,  $C_{loss}$  and  $\alpha$ . These two parameters are used to bound the sensitivity of the quality function, and the value of quality function, respectively. If the clipping threshold  $C_{loss}$  is set to a small value, it significantly



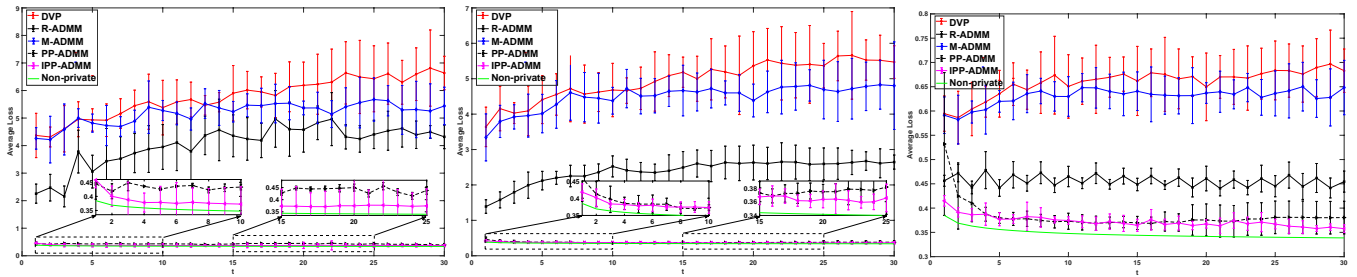


Figure 6: Convergence comparisons on Adult dataset (left:  $\epsilon = 1$ , middle:  $\epsilon = 2$ , right:  $\epsilon = 10$ )

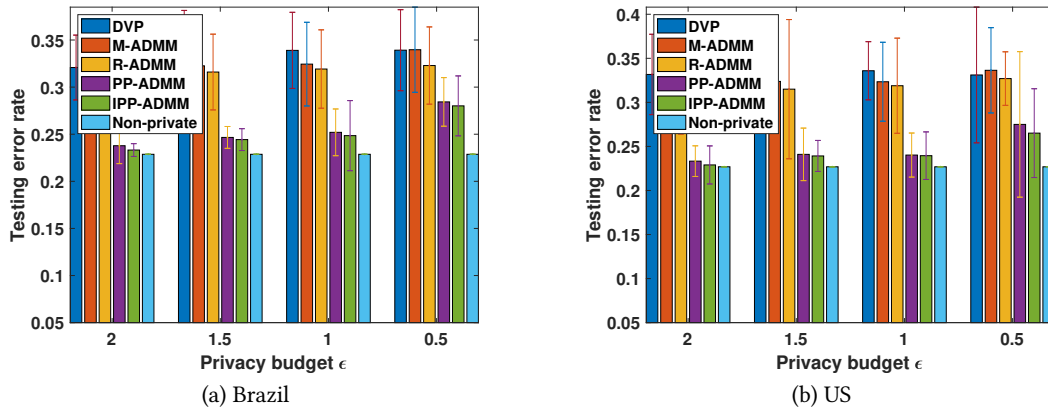


Figure 7: Classification error rate comparisons on Brazil and US datasets.

reduces the sensitivity but at the same time it leads much information loss in the estimation of quality function. On the other hand, if  $C_{loss}$  is large, the sensitivity becomes large that results in adding too much noise to the estimation. Thus, too large or small values of  $C_{loss}$  have a negative effect on employing SVT to check whether the current approximate solution has a big enough difference from that of previous iteration. As we can see from Figure 3,  $C_{loss} = 2$  achieves a good trade-off between high information loss and large sensitivity. In Figure 4, we fix the clipping threshold  $C_{loss} = 2$  and vary  $\alpha$  from  $10^{-3}$  to 10 to see the effect of  $\alpha$  on the performance. Although large value of  $\alpha$  may potentially reduce the releasing of low quality approximate solution and reduce the communication cost, we observe that it also leads the learning performance degradation. We then choose  $\alpha = 10^{-3}$  in our experiments, which achieves the lowest testing/training error rate.

**Performance comparisons.** We also present the trade-off between classification error rate and privacy cost in Figure 5, where we measured the privacy costs of all algorithms to obtain some specified testing error rates. Figure 5 illustrates that both of our methods have consistently lower privacy cost than those baseline algorithms. Compared with PP-ADMM, IPP-ADMM further saves more privacy cost due to limiting the number of releasing low-quality computational results. Additionally, we also inspect the convergence performance (i.e., average loss) of different algorithms under the same budgets, as shown in Fig. 6. We can observe that when budget  $\epsilon$  decreases from 10 to 1, the average loss values of

baseline algorithms increase, which matches the simulation results shown in [7, 25, 26]. Although we also analyze the baseline algorithms using zCDP to provide tight privacy bound, using Gaussian noise instead of Gamma noise might be more beneficial to the performance, which usually has at least  $\sqrt{d}$  times improvement of the empirical risk bound [13], where  $d$  is the dimension of training model. And our proposed algorithms continues to outperform the baseline algorithms significantly.

Figure 7 compares the accuracy (classification error rate) of different algorithms on Brazil and US. The noise parameter of all algorithms are chosen respectively so that they can achieve the same total privacy loss. As expected, the lower privacy budget, the higher classification error rate. As it was observed in the experiments, our proposed algorithms get close to the best achievable classification error rate for a wide range of total privacy loss considered in the experiments.

## 6 CONCLUSIONS

In this paper, we have developed (Improved) plausible differentially private ADMM algorithm, PP-ADMM and IPP-ADMM. In PP-ADMM, in order to release the shackles of the exact optimal solution during each ADMM iteration to ensure differential privacy, we consider outputting a noisy approximate solution for the perturbed objective. To further improve the utility of PP-ADMM, we have adopted SVT in IPP-ADMM to check whether the current

approximate solution has a big enough difference from that of previous iteration. Moreover, we have analyzed privacy loss under the framework of zCDP and generalization performance guarantee. Finally, through the experiments on real-world datasets, we have demonstrated that the proposed algorithms outperform other differentially private ADMM based algorithms while providing the same privacy guarantee. In future work, we plan to extend our privacy analysis to non-convex loss function and apply our methods to deep learning framework. Another research direction is to study the idea of using SVT to other distributed (federated) deep learning framework to save the privacy budget and reduce the communication cost.

## ACKNOWLEDGMENTS

We thank the reviewers for their insightful comments. The work of J. Ding and M. Pan was supported in part by the U.S. National Science Foundation under grants US CNS-1646607, CNS-1801925, and CNS-2029569. The work of J. Bi was partially supported by NSF grants: CCF-1514357 and IIS-1718738, and NIH grant 5K02DA043063-03.

## REFERENCES

- [1] Mark Bun and Thomas Steinke. 2016. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography*. Springer Berlin Heidelberg, Beijing, China, 635–658.
- [2] Tsung-Hui Chang, Mingyi Hong, and Xiangfeng Wang. 2014. Multi-agent distributed optimization via inexact consensus ADMM. *IEEE Transactions on Signal Processing* 63, 2 (2014), 482–497.
- [3] Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. 2011. Differentially private empirical risk minimization. *JMLR* 12, Mar (2011), 1069–1109.
- [4] Jiahao Ding, Sai Mounika Errapotu, Yuanxiong Guo, Haixia Zhang, Dongfeng Yuan, and Miao Pan. 2020. Private Empirical Risk Minimization with Analytic Gaussian Mechanism for Healthcare System. *IEEE Transactions on Big Data* (2020), 1–1.
- [5] Jiahao Ding, Sai Mounika Errapotu, Haijun Zhang, Miao Pan, and Zhu Han. 2019. Stochastic ADMM Based Distributed Machine Learning with Differential Privacy. In *International conference on security and privacy in communication systems*. Springer International Publishing, Orlando, FL, 257–277.
- [6] Jiahao Ding, Yanmin Gong, Chi Zhang, Miao Pan, and Zhu Han. 2019. Optimal Differentially Private ADMM for Distributed Machine Learning. *CoRR* abs/1901.02094 (2019). <https://arxiv.org/abs/1901.02094>
- [7] Jiahao Ding, Xinyue Zhang, Mingsong Chen, Kaiping Xue, Chi Zhang, and Miao Pan. 2019. Differentially Private Robust ADMM for Distributed Machine Learning. In *IEEE International Conference on Big Data*. IEEE, Los Angeles, CA, 1302–1311.
- [8] Jiahao Ding, Xinyue Zhang, Xiaohuan Li, Junyi Wang, Rong Yu, and Miao Pan. 2020. Differentially Private and Fair Classification via Calibrated Functional Mechanism. In *Proceedings of the AAAI Conference on Artificial Intelligence*. AAAI, New York, NY, 622–629.
- [9] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography*. Springer Berlin Heidelberg, Berlin, Heidelberg, 265–284.
- [10] Cynthia Dwork, Moni Naor, Omer Reingold, Guy N Rothblum, and Salil Vadhan. 2009. On the complexity of differentially private data release: efficient algorithms and hardness results. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*. ACM, New York, NY, 381–390.
- [11] Cynthia Dwork and Aaron Roth. 2014. *The algorithmic foundations of differential privacy* (1st ed.). Now Publishers, Inc., Boston, MA, USA.
- [12] Zonghao Huang, Rui Hu, Yuanxiong Guo, Eric Chan-Tin, and Yanmin Gong. 2019. DP-ADMM: ADMM-based distributed learning with differential privacy. *IEEE Transactions on Information Forensics and Security* 15 (July 2019), 1002–1012.
- [13] Daniel Kifer, Adam Smith, and Abhradeep Thakurta. 2012. Private Convex Empirical Risk Minimization and High-dimensional Regression. In *Proceedings of the 25th Annual Conference on Learning Theory*, Vol. 23. PMLR, Edinburgh, Scotland, 25.1–25.40.
- [14] Min Lyu, Dong Su, and Ninghui Li. 2017. Understanding the sparse vector technique for differential privacy. *Proceedings of the VLDB Endowment* 10, 6 (2017), 637–648.
- [15] Gonzalo Mateos, Juan Andrés Bazerque, and Georgios B Giannakis. 2010. Distributed sparse linear regression. *IEEE Transactions on Signal Processing* 58, 10 (2010), 5262–5276.
- [16] Luca Melis, Congzheng Song, Emiliano De Cristofaro, and Vitaly Shmatikov. 2019. Exploiting unintended feature leakage in collaborative learning. In *2019 IEEE S&P*. IEEE, San Francisco, CA, 691–706.
- [17] Joao FC Mota, Joao MF Xavier, Pedro MQ Aguiar, and Markus Püschel. 2013. D-ADMM: A communication-efficient distributed algorithm for separable optimization. *IEEE Transactions on Signal Processing* 61, 10 (2013), 2718–2723.
- [18] Dian Shi, Jiahao Ding, Sai Mounika Errapotu, Hao Yue, Wenjun Xu, Xiangwei Zhou, and Miao Pan. 2019. Deep Q-Network-Based Route Scheduling for TNC Vehicles With Passengers’s Location Differential Privacy. *IEEE Internet of Things Journal* 6, 5 (2019), 7681–7692.
- [19] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks against machine learning models. In *2017 IEEE S&P*. IEEE, San Jose, CA, 3–18.
- [20] Jingyi Wang, Xinyue Zhang, Haijun Zhang, Hai Lin, Hideki Tode, Miao Pan, and Zhu Han. 2018. Data-Driven Optimization for Utility Providers with Differential Privacy of Users’ Energy Profile. In *2018 IEEE Global Communications Conference (GLOBECOM)*. IEEE, Singapore, 1–6.
- [21] L. Yu, L. Liu, C. Pu, M. E. Gursoy, and S. Truex. 2019. Differentially Private Model Publishing for Deep Learning. In *2019 IEEE S&P*. IEEE Computer Society, Los Alamitos, CA, USA, 332–349.
- [22] Ruijiang Zhang and James Kwok. 2014. Asynchronous distributed ADMM for consensus optimization. In *ICML. JMLR*, Beijing, China, 1701–1709.
- [23] Tao Zhang and Quanyan Zhu. 2016. Dynamic differential privacy for ADMM-based distributed classification learning. *IEEE Transactions on Information Forensics and Security* 12, 1 (2016), 172–187.
- [24] Xinyue Zhang, Jiahao Ding, Sai Mounika Errapotu, Xiaoxia Huang, Pan Li, and Miao Pan. 2019. Differentially Private Functional Mechanism for Generative Adversarial Networks. In *2019 IEEE Global Communications Conference (GLOBECOM)*. IEEE, Waikoloa, HI, 1–6.
- [25] Xueru Zhang, Mohammad Mahdi Khalili, and Mingyan Liu. 2018. Improving the Privacy and Accuracy of ADMM-Based Distributed Algorithms. In *ICML. PMLR, Stockholm* Sweden, 5796–5805.
- [26] Xueru Zhang, Mohammad Mahdi Khalili, and Mingyan Liu. 2018. Recycled admm: Improve privacy and accuracy with less computation in distributed algorithms. In *Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, Monticello, IL, 959–965.