

---

# Communication-Efficient Distributed Primal-Dual Algorithm for Saddle Point Problems

---

Yaodong Yu\*

Nanyang Technological University  
ydyu@ntu.edu.sg

Sulin Liu\*

Nanyang Technological University  
liusl@ntu.edu.sg

Sinno Jialin Pan

Nanyang Technological University  
sinnopan@ntu.edu.sg

## Abstract

Primal-dual algorithms, which are proposed to solve reformulated convex-concave saddle point problems, have been proven to be effective for solving a generic class of convex optimization problems, especially when the problems are ill-conditioned. However, the saddle point problem still lacks a distributed optimization framework where primal-dual algorithms can be employed. In this paper, we propose a novel communication-efficient distributed optimization framework to solve the convex-concave saddle point problem based on primal-dual methods. We carefully design local subproblems and a central problem such that our proposed distributed optimization framework is communication-efficient. We provide a convergence analysis of our proposed algorithm, and extend it to address non-smooth and non-strongly convex loss functions. We conduct extensive experiments on several real-world datasets to demonstrate competitive performance of the proposed method, especially on ill-conditioned problems.

## 1 INTRODUCTION

In the era of big data, developing distributed machine learning algorithms has become increasingly important yet challenging. In this work, we focus on developing a new distributed optimization algorithm for regularized *empirical risk minimization* (ERM), which is a generic class of convex optimization problems that arises often from machine learning. Specifically, our goal is to mini-

mize the empirical loss defined over  $n$  data samples:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{P}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \phi_i(a_i^\top \mathbf{x}) + g(\mathbf{x}), \quad (1)$$

where  $a_1, \dots, a_n \in \mathbb{R}^d$  are feature vectors of  $n$  data points,  $\phi_i : \mathbb{R} \rightarrow \mathbb{R}$  is a convex loss function with the linear predictor  $a_i^\top \mathbf{x}$ , for  $i = 1, \dots, n$ , and  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is a convex regularization function for the predictor  $\mathbf{x} \in \mathbb{R}^d$ . Suppose that a distributed system consists of  $K$  machines, and each machine has access only a subset  $\mathcal{P}_k$  of the data  $[n] := \{1, \dots, n\}$ , where  $\{\mathcal{P}_k\}_{k=1}^K$  is a given partition of the dataset  $[n]$ , and we denote by  $n_k = |\mathcal{P}_k|$ .

### 1.1 COMMUNICATION-EFFICIENT DISTRIBUTED OPTIMIZATION

One of the most important issues in distributed optimization is the communication efficiency because communication between machines is much more expensive than reading data from the memory on local machines. Therefore, more and more efforts have been made in proposing communication-efficient methods in distributed optimization (Jaggi et al., 2014, Ma et al., 2015, Reddi et al., 2016, Shamir et al., 2014, Smith et al., 2015, Yang, 2013). The basic idea behind these methods is to carefully design local computation and communication in a distributed system. To achieve this goal, existing methods usually decompose the optimization problem into  $K$  local subproblems, denoted by  $\mathcal{L}_k$ , with respect to local data of each machine. After each local machine performs an arbitrary optimization method on  $\mathcal{L}_k$  and solves it approximately, updated information from local machines is sent to a central node to carry out a central update. This allows one to control the trade-off between communication and local computation, which is more flexible in distributed setting.

In general, existing communication-efficient methods can be classified into two categories. A first category

---

\*Indicates equal contributions.

is referred to as *gradient*-type, which aims to decompose the problem (1) into  $K$  local subproblems and solve each subproblem using gradient descent methods, such as stochastic variance reduced gradient (SVRG) (Johnson and Zhang, 2013), on each local machine independently, and then update the optimization variables by doing an average of local variables of each machine (Reddi et al., 2016, Shamir et al., 2014).

The second one is referred to as *coordinate*-type, which usually focuses on solving the dual problem of (1),

$$\max_{\alpha \in \mathbb{R}^n} \mathcal{D}(\alpha) = \frac{1}{n} \sum_{i=1}^n -\phi_i^*(\alpha_i) - g^*\left(-\frac{1}{n} \sum_{i=1}^n \alpha_i a_i\right), \quad (2)$$

where  $\phi_i^*$  is the conjugate function of  $\phi_i$ ,  $\alpha$  is the dual variable vector with the  $i$ -th element being  $\alpha_i$ . This category of methods decomposes the dual problem (2) into  $K$  local subproblems with regard to each machine, and employs a coordinate-type method, such as stochastic dual coordinate ascent method (SDCA) (Shalev-Shwartz and Zhang, 2013), to solve each local subproblem. After each local subproblem is approximately solved, the central node takes an average/add step according to the local update from each machine (Jaggi et al., 2014, Ma et al., 2015, Yang, 2013). Our proposed method falls into this category.

There are two limitations of the *coordinate*-type approach. One limitation is that most existing coordinate-type methods fail to match the communication complexity lower bounds proved in (Arjevani and Shamir, 2015). For example, when  $\phi_i$  is  $(1/\gamma)$ -smooth,  $g$  is  $\lambda$ -strongly convex, and  $R = \max_i \|a_i\|_2$  in problem (1), then the *condition number* is defined as

$$\kappa = R^2/(\lambda\gamma), \quad (3)$$

and the lower bound of communication complexity obtained in (Arjevani and Shamir, 2015) is  $\tilde{\mathcal{O}}(\sqrt{\kappa} \log(1/\epsilon))$  in order to achieve an  $\epsilon$ -suboptimal solution (Nesterov, 2013). Since these *coordinate*-type methods (Jaggi et al., 2014, Ma et al., 2015, Yang, 2013) could be seen as a distributed version of SDCA, and the iteration complexity of SDCA is  $\tilde{\mathcal{O}}(\kappa \log(1/\epsilon))$ , the communication complexity of these methods (Jaggi et al., 2014, Ma et al., 2015, Yang, 2013) is  $\tilde{\mathcal{O}}(\kappa \log(1/\epsilon))$ , which does not match the optimal communication complexity. The problem could get worse when problem (1) is ill-conditioned, i.e.,  $\kappa \gg 1$ . The other limitation is that it is not straight-forward to extend existing coordinate-type approaches (Jaggi et al., 2014, Ma et al., 2015, Yang, 2013) to deal with different regularization forms of the regularization term  $g(\cdot)$  other than  $\ell_2$  regularization.

## 1.2 OUR APPROACH

To overcome the limitations mentioned above, we propose a novel communication-efficient distributed framework with a primal-dual algorithm for saddle point problems (Chambolle and Pock, 2011, Zhang and Xiao, 2015b). The reasons that we focus on developing distributed framework for saddle point problems are three folds:

1. It has been shown that saddle point algorithms are able to obtain comparable and even more stable performance than other state-of-the-art techniques for convex optimization (Chambolle and Pock, 2011, Yu et al., 2015, Zhang and Xiao, 2015b).
2. Primal-dual coordinate algorithms for saddle point problems are able to reach the optimal iteration complexity as proved by Zhang and Xiao (2015b).
3. Since primal-dual algorithms keep both primal and dual variables in optimization, they are able to deal with different kinds of regularization terms  $g(\cdot)$ , such as  $\ell_1$ -regularization, naturally.

Therefore, we aim to develop a distributed primal-dual algorithm which can inherit the above three properties in a distributed environment. To be specific, we first reformulate problem (1) as a convex-concave saddle point problem through convex conjugation, which keeps both primal and dual variables in optimization. We then decompose the saddle point problem into carefully designed local subproblems, which can be solved independently by each machine with respect to its local data. After the local subproblems are solved approximately, the updated local dual variables of each machine are sent to a central node. Based on all the updates of local dual variables, we construct a central problem and get the central update of the primal variables on the central node. Finally, the central node sends the updated primal variables and aggregated dual variables to each local machine for next local optimization iteration. Details will be described in Section 3.

As will be discussed in Section 4 and Section 5, by carefully designing the local subproblems and the central problem, our algorithm is able to reach the communication complexity lower bounds, and deal with different kinds of regularization terms. As the parameter  $\lambda$  of the regularization function  $g(\cdot)$  is usually on the order of  $1/\sqrt{n}$  or  $1/n$  for many machine learning problems, we are especially interested in solving problem (1) under the distributed setting when it is ill-conditioned. With an extremely large-scale dataset, the *condition number* in (3) could be relatively large ( $\kappa \gg 1$ ). We will

show through experiments in Section 6 that our proposed algorithm can obtain better and more stable performance on ill-conditioned problems compared to other communication-efficient distributed optimization methods. Moreover, we provide a solution to extend our algorithm to deal with non-smooth and non-strongly convex optimization problems in Section 5.

### 1.3 OTHER RELATED WORK

Besides the communication-efficient distributed optimization methods reviewed above, there exist other parallel and distributed optimization techniques. For example, a well-known approach for solving problem (1) is to perform a gradient descent method implemented in a distributed system. Each local machine computes its local gradient and sends it to the central node. The central node aggregates the local gradient to take a gradient step and updates  $\mathbf{x}$ , and then broadcasts it back to each local machine for the next iteration updates. If the accelerated gradient descent method (Nesterov, 2013) is used, one can obtain the iteration complexity  $\tilde{O}(\sqrt{\kappa} \log(1/\epsilon))$ . Another popular technique is distributed alternating direction method of multipliers (ADMM) (Boyd et al., 2011, Shi et al., 2014), whose complexity is  $\tilde{O}(\sqrt{\kappa} \log(1/\epsilon))$  under certain conditions. Recently, Zhang and Xiao (2015a) proposed a distributed algorithm based on the inexact damped Newton method, which matches the communication complexity lower bound  $\tilde{O}(\sqrt{\kappa} \log(1/\epsilon))$ . For more related distributed and parallel optimization methods, we refer the readers to (Arjevani and Shamir, 2015, Chang et al., 2014, Duchi et al., 2013, Zhang et al., 2012) for a more comprehensive literature of the related work.

## 2 PRELIMINARIES

In this section, we introduce first-order primal-dual algorithms for convex-concave saddle point problems in a single machine. To begin with, we reformulate the primal problem (1) as a convex-concave saddle point problem through convex conjugation, and the saddle point reformulation is widely applied in machine learning (Dai et al., 2016, Zhang and Xiao, 2015b). Based on the definition of convex conjugate, we replace each component function  $\phi_i(a_i^\top \mathbf{x})$  by

$$\phi_i(a_i^\top \mathbf{x}) = \sup_{y_i \in \mathbb{R}} \{y_i \langle a_i, \mathbf{x} \rangle - \phi_i^*(y_i)\},$$

where  $\phi_i^*(y_i) = \sup_{\alpha \in \mathbb{R}} \{\alpha y_i - \phi_i(\alpha)\}$  is the convex conjugate of  $\phi_i$ , then we arrive at a convex-concave saddle point problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{x}, \mathbf{y}), \quad (4)$$

where

$$f(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n (y_i \langle a_i, \mathbf{x} \rangle - \phi_i^*(y_i)) + g(\mathbf{x}),$$

where  $\mathbf{y} \in \mathbb{R}^n$  is referred to as a vector of dual variables with its  $i$ -th element denoted by  $y_i$ , and each  $y_i$  is associated with a data point  $a_i$ . Under the assumption that  $\phi_i$  is  $(1/\gamma)$ -smooth,  $g$  is  $\lambda$ -strongly convex, and each  $\phi_i^*$  is  $\gamma$ -strongly convex, the saddle point problem (4) has a unique solution, which is denoted by  $(\mathbf{x}^*, \mathbf{y}^*)$ .

Based on the reformulated problem (4), we could rewrite the optimization problem as

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^n} f(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \langle A\mathbf{y}, \mathbf{x} \rangle - \Phi^*(\mathbf{y}) + g(\mathbf{x}),$$

where  $A = [a_1, \dots, a_n]$  and  $\Phi^*(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \phi_i^*(y_i)$ . This is the generic saddle point problem (Chambolle and Pock, 2011). To solve this optimization problem, the basic idea is to alternately maximize  $f$  with respect to  $\mathbf{y}$ , and minimize  $f$  with respect to  $\mathbf{x}$ , which is

$$\begin{aligned} \mathbf{y}^{(t+1)} &= \arg \max_{\mathbf{y} \in \mathbb{R}^n} \frac{1}{n} \langle A\mathbf{y}, \bar{\mathbf{x}}^{(t)} \rangle - \Phi^*(\mathbf{y}) - \frac{\|\mathbf{y} - \mathbf{y}^{(t)}\|_2^2}{2\sigma n}, \\ \mathbf{x}^{(t+1)} &= \arg \min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{n} \langle A\mathbf{y}^{(t+1)}, \mathbf{x} \rangle + g(\mathbf{x}) + \frac{\|\mathbf{x} - \mathbf{x}^{(t)}\|_2^2}{2\tau}, \\ \bar{\mathbf{x}}^{(t+1)} &= \mathbf{x}^{(t+1)} + \theta (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}), \end{aligned}$$

where the parameters  $\tau$  and  $\sigma$  control the quadratic regularization terms with respect to  $\mathbf{x}$  and  $\mathbf{y}$ , respectively, which is similar to the use of step size in primal methods. And  $\bar{\mathbf{x}}^{(t+1)}$  is the extrapolation from  $\mathbf{x}^{(t)}$  and  $\mathbf{x}^{(t+1)}$  with parameter  $\theta \in [0, 1]$ . Here,  $\theta (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)})$  is similar to the momentum term in Nesterov acceleration (Nesterov, 2013). There have been various primal-dual coordinate algorithms designed based on the above scheme, such as (Chambolle and Pock, 2011, Yu et al., 2015, Zhang and Xiao, 2015b). Among them, Yu et al. (2015), Zhang and Xiao (2015b) developed stochastic versions of the primal-dual algorithm for the saddle point problem. Suppose  $\phi_i$  is  $(1/\gamma)$ -smooth and  $g$  is  $\lambda$ -strongly convex, most existing algorithms (Chambolle and Pock, 2011, Yu et al., 2015, Zhang and Xiao, 2015b) can achieve linear convergence rate to obtain an  $\epsilon$ -accurate solution. The (expected) iteration complexity of these methods is  $\tilde{O}(\sqrt{\kappa} \log(1/\epsilon))$ , which is desired on ill-conditioned problems.

## 3 DISTRIBUTED SADDLE-POINT FRAMEWORK

In this section, we present our proposed communication-efficient Distributed Saddle Point Algorithm (DiSPA) in

detail. Assume that the dataset  $\{a_i\}_{i=1}^n$  is distributed over  $K$  local machines, each machine  $k$  has access to a subset  $\mathcal{P}_k$  of the data with size of  $n_k$ . Here we define  $K$  local vectors of  $\mathbf{y}$  by using the notation  $\mathbf{y}_{[k]} \in \mathbb{R}^n$  for  $k = 1, \dots, K$ :

$$(\mathbf{y}_{[k]})_i = \begin{cases} y_i, & \text{if } i \in \mathcal{P}_k, \\ 0, & \text{otherwise,} \end{cases}$$

and  $K$  local data matrix of  $A$  by using the notation  $A_{[k]} \in \mathbb{R}^{d \times n}$  for  $k = 1, \dots, K$ :

$$(A_{[k]})_i = \begin{cases} a_i, & \text{if } i \in \mathcal{P}_k, \\ \mathbf{0}, & \text{otherwise,} \end{cases}$$

where  $(A_{[k]})_i$  denotes the  $i$ -th column of  $A_{[k]}$ , and  $\mathbf{0}$  is a  $d$ -dimensional vector of all zeros. According to the data partition, we decompose the problem (4) into  $K$  local subproblems, and define the associated central problem. The core idea is to solve the defined **local subproblems** on each local machine independently, and then centralize the updated information of each machine to solve a easy **central problem**.

### 3.1 LOCAL SUBPROBLEMS

We define a local subproblem of the original saddle point problem (4) for machine  $k$ , which only requires accessing data that is available locally, i.e.,  $\mathcal{P}_k$ . Specifically, each local subproblem of machine  $k$  at  $t$ -iteration is defined as

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y}_{[k]} \in \mathbb{R}^n} \mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]}), \quad (5)$$

where

$$\begin{aligned} \mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]}) := & \frac{1}{n} \sum_{j \in \mathcal{P}_k} (y_j \langle a_j, \mathbf{x} \rangle - \phi_j^*(y_j)) + g(\mathbf{x}) \\ & + \frac{1}{n} \sum_{i \notin \mathcal{P}_k} y_i^{(t-1)} \langle a_i, \mathbf{x} \rangle + r_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]}), \end{aligned}$$

and

$$r_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]}) = \frac{\|\mathbf{x} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} - \frac{K \|\mathbf{y}_{[k]} - \mathbf{y}_{[k]}^{(t-1)}\|_2^2}{2\sigma n}.$$

Each local subproblem is related to the solution obtained in the previous iteration  $(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)})$ , and we will define the parameters  $\tau, \sigma$  later. The first two terms can be regarded as a local version of the original saddle point problem (4). The third term  $\frac{1}{n} \sum_{i \notin \mathcal{P}_k} y_i^{(t-1)} \langle a_i, \mathbf{x} \rangle = \frac{1}{n} \langle \sum_{l \neq k} A_{[l]} \mathbf{y}_{[l]}^{(t-1)}, \mathbf{x} \rangle$  in (5) can be interpreted as the interaction with other local subproblems. And the quadratic regularization term  $r_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]})$  is used to enforce  $\mathbf{x}$  and  $\mathbf{y}_{[k]}$  not to move too far away from solution obtained at the previous iteration.

To solve the local subproblem  $\mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]})$ , we could apply primal-dual coordinate type of methods (Yu et al., 2015, Zhang and Xiao, 2015b). The input needed for solving the local subproblem would be  $\mathbf{x}^{(t-1)}$  and  $A\mathbf{y}^{(t-1)}$ , as the term  $\sum_{l \neq k} A_{[l]} \mathbf{y}_{[l]}^{(t-1)}$  can be expressed as  $(A\mathbf{y}^{(t-1)} - A_{[k]} \mathbf{y}_{[k]}^{(t-1)})$  and machine  $k$  has access to  $(A_{[k]}, \mathbf{y}_{[k]}^{(t-1)})$ . By denoting our local optimization algorithm used as ‘Local-DiSPA’, we could denote the procedure of solving the local subproblem as  $(\mathbf{x}_k^{(t)}, \mathbf{y}_{[k]}^{(t)}) \leftarrow \text{Local-DiSPA}(\mathbf{x}^{(t-1)}, A\mathbf{y}^{(t-1)})$ . Note that in our framework, each local subproblem is only required to be solved approximately.

### 3.2 CENTRAL PROBLEM

Let  $(\mathbf{x}_k^{(t)}, \mathbf{y}_{[k]}^{(t)})$  denote the approximate solution obtained by solving  $\mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]})$  on machine  $k$ . Let  $\mathbf{y}^{(t)} = \sum_{k=1}^K \mathbf{y}_{[k]}^{(t)}$ , we define a central problem at  $t$ -iteration as

$$\mathbf{x}^{(t)} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{C}^{(t)}(\mathbf{x}), \quad (6)$$

where

$$\mathcal{C}^{(t)}(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n y_i^{(t)} \langle a_i, \mathbf{x} \rangle + g(\mathbf{x}) + \frac{\|\mathbf{x} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau}.$$

The central update involves two steps. Firstly, updated local  $A_{[k]} \mathbf{y}_{[k]}^{(t)} \in \mathbb{R}^d$  are sent to the central node. After that, we solve (6) on the central node, and then send the updated  $\mathbf{x}^{(t)} \in \mathbb{R}^d$  and  $A\mathbf{y}^{(t)} \in \mathbb{R}^d$  of the  $t$ -iteration back to each local machine. The central update could be viewed as minimizing the original saddle point problem (4) over  $\mathbf{x}$  with updated  $\mathbf{y}^{(t)}$ . Also, we should observe that the parameter  $\tau$  defined in (5) and (6) is the same parameter, which is one of the key points of our framework.

### 3.3 OVERALL ALGORITHM DESCRIPTION

The overall algorithm of DiSPA is presented in Algorithm 1. Firstly, we distribute the initial dual variable  $\mathbf{y}^{(0)}$  to  $K$  local machines, and aggregate  $A_{[k]} \mathbf{y}_{[k]}^{(0)}$  from all machines to compute  $A\mathbf{y}^{(0)} \leftarrow \sum_{k=1}^K A_{[k]} \mathbf{y}_{[k]}^{(0)}$ . After that, the initial primal variables  $\mathbf{x}^{(0)}$  and the transformed dual variables  $A\mathbf{y}^{(0)}$  are sent to each machine (Steps 2-4). Secondly, each machine performs optimization on its own local subproblem  $\mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]})$ , using a primal-dual algorithm, such as SPDC (Zhang and Xiao, 2015b), and sends its local update  $A_{[k]} \mathbf{y}_{[k]}^{(t)} \in \mathbb{R}^d$  to the central

---

**Algorithm 1** DiSPA( $f, \mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \tau, \sigma, T$ )

---

- 1: **Input:** Data points  $\{a_i\}_{i=1}^n$  distributed across  $K$  machines  $\{\mathcal{P}_k\}_{k=1}^K$ , parameters  $\tau, \sigma \in \mathbb{R}$ , initial primal variables  $\mathbf{x}^{(0)} \in \mathbb{R}^d$  and initial dual variables  $\mathbf{y}^{(0)} \in \mathbb{R}^n$ .
  - 2: Distribute  $\mathbf{y}^{(0)}$  to  $K$  machines as  $\mathbf{y}_{[k]}^{(0)}$  for each.
  - 3: Each machine compute  $A_{[k]}\mathbf{y}_{[k]}^{(0)}$ , and send it back to the central node.
  - 4: Perform  $A\mathbf{y}^{(0)} \leftarrow \sum_{k=1}^K A_{[k]}\mathbf{y}_{[k]}^{(0)}$  on the central node, and send  $\mathbf{x}^{(0)}$  and  $A\mathbf{y}^{(0)}$  to each machine
  - 5: **for**  $t = 1, 2, \dots, T$  **do**
  - 6:   **for**  $k \in [K]$  in parallel over all machines **do**
  - 7:      $(\mathbf{x}_k^{(t)}, \mathbf{y}_{[k]}^{(t)}) \leftarrow \text{Local-DiSPA}(\mathbf{x}^{(t-1)}, A\mathbf{y}^{(t-1)})$
  - 8:     Send  $A_{[k]}\mathbf{y}_{[k]}^{(t)}$  to the central node
  - 9:   **end for**
  - 10:    $A\mathbf{y}^{(t)} \leftarrow \sum_{k=1}^K A_{[k]}\mathbf{y}_{[k]}^{(t)}$  (on the central node)
  - 11:    $\mathbf{x}^{(t)} \leftarrow \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{C}^{(t)}(\mathbf{x})$  (on the central node)
  - 12:   Send  $\mathbf{x}^{(t)}$  and  $A\mathbf{y}^{(t)}$  back to each machine
  - 13: **end for**
  - 14: **Output:**  $(\mathbf{x}^{(T)}, \mathbf{y}^{(T)})$
- 

node. The goal of this internal procedure is to approximately solve each local subproblem. We refer this procedure as the inner loop, which is stated in Steps 6-9 of Algorithm 1. Thirdly, with the local updates  $A_{[k]}\mathbf{y}_{[k]}^{(t)}$ 's, the central node aggregates them to compute  $A\mathbf{y}^{(t)}$ , and solves the defined central problem  $\mathcal{C}^{(t)}(\mathbf{x})$  to get the update of  $\mathbf{x}^{(t)}$  (Steps 10-11 of Algorithm 1). Finally, the central node sends  $(\mathbf{x}^{(t)}, A\mathbf{y}^{(t)})$  back to each machine in order to start the next round of local optimization (Step 12), we refer Steps 5-13 as the outer loop of our algorithm.

Note that, on one hand, the essential information that the central problem requires can be represented by the vectors  $\{A_{[k]}\mathbf{y}_{[k]}^{(t)}\}$ 's of each local machine, which are of dimension of  $d$ . On the other hand, each local machine only requires  $(\mathbf{x}^{(t)}, A\mathbf{y}^{(t)})$ , which is of  $2d$  dimensions, from the central node.

## 4 CONVERGENCE ANALYSIS

Before we present our main convergence results, we make some assumptions on the objective functions and the local suboptimality on each local subproblem.

**Assumption 1.** Each  $\phi_i(\cdot)$  is convex and differentiable, and  $\phi_i(\cdot)$  is  $(1/\gamma)$ -smooth, i.e.,

$$|\phi'_i(a) - \phi'_i(b)| \leq (1/\gamma)|a - b|, \quad \forall a, b \in \mathbb{R}.$$

**Assumption 2.**  $g(\cdot)$  is  $\lambda$ -strongly convex i.e. for  $\forall x, y \in \mathbb{R}^d$  and  $g'(y) \in \partial g(y)$ ,

$$g(x) \geq g(y) + g'(y)(x - y) + \frac{\lambda}{2}\|x - y\|_2^2.$$

**Assumption 3.** There exists constants  $\Omega_{\mathbf{x}}, \Omega_{\mathbf{y}} > 0$  such that

$$\begin{aligned} \max_t \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 &\leq \Omega_{\mathbf{x}}, \\ \max_t \|\mathbf{y}^{(t)} - \mathbf{y}^*\|_2 &\leq \Omega_{\mathbf{y}}, \end{aligned}$$

where  $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$  is  $t$ -iteration update of Algorithm 1, and  $(\mathbf{x}^*, \mathbf{y}^*)$  is the optimal solution of (4).

**Assumption 4.** Each local machine  $k$  produces an approximate solution  $(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)})$  that satisfies

$$\mathbb{E} \left[ \left\| \mathbf{y}_{[k]}^{(t)} - \hat{\mathbf{y}}_{[k]}^{(t)} \right\|_2 \right] \leq \hat{\Theta} \left\| \mathbf{y}_{[k]}^{(t-1)} - \hat{\mathbf{y}}_{[k]}^{(t)} \right\|_2,$$

and

$$\begin{aligned} \mathbb{E}[(f(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) - f(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]}^{(t)}))] \\ \leq \tilde{\Theta} \left( f(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) - f(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]}^{(t-1)}) \right), \end{aligned}$$

where  $\hat{\Theta}, \tilde{\Theta} < 1$ , and  $(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)})$  is the optimal solution of local subproblem  $\mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]})$ .

**Assumption 5.** There exist constants  $c', c'' > 0$ , and the approximate solution  $\mathbf{y}^{(t)} = \sum_{k=1}^K \mathbf{y}_{[k]}^{(t)}$  satisfies

$$\mathbb{E}[\blacktriangle] \leq \Theta(f(\mathbf{x}^{(t-1)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)})),$$

and  $\blacktriangle$  is defined as

$$\blacktriangle = \sum_{k=1}^K (f(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) - f(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]}^{(t)})) + \frac{M}{n} \|\mathbf{y}^{(t)} - \hat{\mathbf{y}}^{(t)}\|_2,$$

where  $M = c'\Omega_{\mathbf{x}} + c''\Omega_{\mathbf{y}}$ , and  $\Theta \in (0, 1)$  is a pre-defined constant.

Note that Assumptions 4 and 5 imply that the local subproblems need to be solved approximately to some extent. As the local subproblems are saddle point problems, we can apply different kinds of algorithms for saddle point problems to solve them, such as (Chambolle and Pock, 2011, Yu et al., 2015, Zhang and Xiao, 2015b), which include both stochastic and non-stochastic optimization methods.

To analyse the convergence behavior of our distributed algorithm, we need to characterize the connection between local subproblems and central update on the central node. Based on the relationship between the central update  $\mathbf{x}^{(t)}$  and the local optimal solution  $(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)})$ , we can get the convergence guarantee of Algorithm 1. All proofs can be found in Appendix.

**Lemma 1.** *Suppose that Assumptions 1-5 hold. Let  $(\mathbf{x}^*, \mathbf{y}^*)$  be the unique saddle point of  $f(\mathbf{x}, \mathbf{y})$  defined in (4), and define*

$$\begin{aligned} \Delta^{(t)} = & f(\mathbf{x}^{(t)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t)}) + \frac{1 + \tau\lambda}{2\tau} \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 \\ & + \frac{2K^2 + 2K\sigma\gamma - 1}{4\sigma nK} \|\mathbf{y}^{(t)} - \mathbf{y}^*\|_2^2. \end{aligned}$$

*If the parameters  $\tau$  and  $\sigma$  are chosen such that  $\tau\sigma = \frac{1}{4R^2}$  and  $\sigma > \frac{1}{2K\gamma}$ , then for  $t \geq 1$ , the proposed DiSPA (Algorithm 1) achieves*

$$\mathbb{E}[\Delta^{(t)}] \leq \Theta \mathbb{E}[\Delta^{(t-1)}], \quad (7)$$

where  $\Theta = \max\left\{\frac{1}{1 + \frac{\sigma\gamma}{K} - \frac{1}{2K^2}}, \frac{1}{1 + \tau\lambda}\right\}$ .

**Theorem 1.** *Suppose that Assumptions 1-5 hold, and the parameters  $\tau = \gamma/(4R^2)$ ,  $\sigma = 1/\gamma$  and  $K \leq \kappa$ . In order for Algorithm 1 to obtain*

$$\mathbb{E}[\|\mathbf{x}^{(T)} - \mathbf{x}^*\|_2^2] \leq \epsilon,$$

*it suffices to have the number of communication iterations  $T$  satisfy*

$$T \geq \left(1 + \frac{4R^2}{\lambda\gamma}\right) \log\left(\frac{C}{\epsilon}\right),$$

where  $C = \frac{\Delta^{(0)}}{\frac{\lambda}{2}\left(1 + \frac{4R^2}{\lambda\gamma}\right)}$ .

Based on Lemma 1 and Theorem 1, we can get the communication complexity of Algorithm 1. By taking  $\tau = \frac{\gamma}{4R^2}$  and  $\sigma = \frac{1}{\gamma}$ , the communication complexity of Algorithm 1 is  $\tilde{O}(\kappa \log(1/\epsilon))$ . This does not match the optimal communication complexity. In the next section, we derive an accelerated version of DiSPA that attains the optimal communication complexity.

## 5 EXTENSION OF DISPA

In this section, we derive two extensions of DiSPA. The first extension applies a generic acceleration scheme proposed by Lin et al. (2015) to obtain an accelerated version of DiSPA. The second one extends our algorithm to handle non-smooth and non-strongly convex loss functions with convergence guarantee.

### 5.1 ACCELERATION

The communication complexity of our algorithm derived in Theorem 1 is  $\tilde{O}(\kappa \log(1/\epsilon))$  in order to achieve an  $\epsilon$ -suboptimal solution, which does not match the lower bound proved in (Arjevani and Shamir, 2015). Based on the ‘Catalyst’ acceleration scheme proposed by Lin et al. (2015), we develop the accelerated DiSPA to achieve

the optimal communication complexity. ‘Catalyst’ is a generic scheme for accelerating first-order optimization methods, similar to classical gradient descent schemes of Nesterov acceleration. It is a natural choice for distributed optimization algorithms that solve local sub-problem approximately in each communication iteration. Based on Catalyst acceleration, we modify the original objective function by adding a quadratic term

$$f^{(t)}(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}, \mathbf{y}) + \frac{\vartheta}{2} \|\mathbf{x} - \mathbf{z}^{(t-1)}\|_2^2,$$

where  $\vartheta > 0$  is a parameter defined in Catalyst,  $\mathbf{z}^{(t)}$  is obtained by an extrapolation step similar to the acceleration scheme in (Nesterov, 2013). With a carefully selected parameter  $\vartheta$ , we can develop an accelerated version of DiSPA described in Algorithm 2, which achieves the optimal communication complexity. We refer to the accelerated distributed saddle point algorithm as A-DiSPA.

**Theorem 2.** *Suppose Assumption 1-5 hold and assume  $\vartheta = \frac{4R^2}{\gamma} - \lambda > 0$ , if the parameters in Catalyst are chosen as  $T_s = \tilde{O}(1)$ ,  $q = \lambda/(\lambda + \vartheta)$ ,  $\alpha_0 = \sqrt{q}$ . The total communication iterations of Algorithm 2 for achieving  $\mathcal{P}(\mathbf{x}^{(T)}) - \mathcal{P}(\mathbf{x}^*) < \epsilon$  is  $\tilde{O}(\sqrt{\kappa} \log(\frac{1}{\epsilon}))$ .*

The ‘Catalyst’ acceleration scheme is applied widely to accelerate different kinds of optimization methods (Reddi et al., 2016). As the ‘acceleration’ quadratic term is related to the primal variable  $\mathbf{x}$  and the saddle point problem keeps the primal variable, it is convenient to apply this acceleration scheme to DiSPA and achieve the optimal communication complexity.

### 5.2 NON-SMOOTH OR/AND NON-STRONGLY CONVEX FUNCTIONS

The communication complexity bounds in Theorem 1 and Theorem 2 are developed under the Assumptions 1 and 2, which means that the derivative of  $\phi_i(\cdot)$  needs to be  $(1/\gamma)$ -Lipschitz continuous and  $g(\cdot)$  needs to be  $\lambda$ -strongly convex. For general loss functions in machine learning, both assumptions may fail, for example, when  $\phi_i(\cdot)$  is hinge loss and  $g(\cdot)$  is  $\ell_1$ -regularization term. Therefore, we aim to develop an extension of Algorithm 2 to deal with non-smooth and non-strongly convex loss functions.

To be concise, we only consider the case when  $\phi_i(\cdot)$  is non-smooth and  $g(\cdot)$  is non-strongly convex. The key idea here is to apply a perturbation term to analyze the non-smooth and non-strongly convex setting, which is commonly used in many optimization methods (Reddi et al., 2016, Zhang and Xiao, 2015b). Here we assume that neither  $\phi_i(\cdot)$  is smooth nor  $g(\cdot)$  is strongly convex.

---

**Algorithm 2** A-DiSPA

- 1: **Input:** Data points  $\{a_i\}_{i=1}^n$  distributed across  $K$  machines  $\{\mathcal{P}_k\}_{k=1}^K$ , parameters  $\tau, \sigma, \vartheta, \alpha_0 \in \mathbb{R}$ , initial primal variable  $\mathbf{x}^{(0)}$  generated on the central node, initial local dual variables  $\mathbf{y}^{(0)}$  generated on local machines, optimization algorithm DiSPA, number of inner iterations  $T_s$  of DiSPA.
  - 2: **Initialization:**  $q = \lambda/(\lambda + \vartheta)$ ,  $\mathbf{z}^{(0)} = \mathbf{x}^{(0)}$ ,  $t = 1$
  - 3: **while** stopping criterion is not satisfied **do**
  - 4:    $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) = \text{DiSPA}(f^{(t)}, \mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}, \tau, \sigma, T_s)$
  - 5:   Compute  $\alpha_t \in (0, 1)$  from equation
 
$$\alpha_t^2 = (1 - \alpha_t)\alpha_{t-1}^2 + q\alpha_t$$
  - 6:   Compute  $\mathbf{z}^{(t)} = \mathbf{x}^{(t)} + \beta_t(\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)})$ , where
 
$$\beta_t = \frac{\alpha_{t-1}(1 - \alpha_{t-1})}{\alpha_{k-1}^2 + \alpha_t}$$
  - 7:    $t = t + 1$
  - 8: **end while**
  - 9: **Output**  $\mathbf{x}^{(t)}$
- 

In particular, we consider the following modified saddle point function

$$f_\epsilon(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \left( y_i \langle a_i, \mathbf{x} \rangle - \left( \phi_i^*(y_i) + \frac{\epsilon}{2} y_i^2 \right) \right) + \left( g(\mathbf{x}) + \frac{\epsilon}{2} \|\mathbf{x}\|_2^2 \right), \quad (8)$$

where  $\epsilon > 0$  is a pre-defined scalar.  $f_\epsilon(\mathbf{x}, \mathbf{y})$  can be regarded as an approximation of the original saddle point problem  $f(\mathbf{x}, \mathbf{y})$ . We then obtain that both  $\phi_i^*(y_i) + \frac{\epsilon}{2} y_i^2$  and  $g(\mathbf{x}) + \frac{\epsilon}{2} \|\mathbf{x}\|_2^2$  are  $\epsilon$ -strongly convex. In the following corollary, we show that the perturbed function  $f_\epsilon(\mathbf{x}, \mathbf{y})$  is a good approximation of the origin function  $f(\mathbf{x}, \mathbf{y})$ .

**Corollary 1.** *Assume that  $\phi_i$  is convex and  $L_\phi$ -Lipschitz continuous, and  $g$  is convex. Let  $(\mathbf{x}^*, \mathbf{y}^*)$  be the unique saddle point of  $f(\mathbf{x}, \mathbf{y})$ , and  $(\mathbf{x}_\epsilon^*, \mathbf{y}_\epsilon^*)$  be the unique saddle point of  $f_\epsilon(\mathbf{x}, \mathbf{y})$ . Then we have*

$$\mathcal{P}(\mathbf{x}_\epsilon^*) - \mathcal{P}(\mathbf{x}^*) \leq \frac{\epsilon}{2} (\|\mathbf{x}^*\|_2^2 + L_\phi^2).$$

As a result, we could apply Algorithm 2 to the perturbed function  $f_\epsilon(\mathbf{x}, \mathbf{y})$ . By applying Theorem 2 and the relationship between  $f_\epsilon(\mathbf{x}, \mathbf{y})$  and  $f(\mathbf{x}, \mathbf{y})$ , we obtain that the communication complexity of Algorithm 2 for non-smooth and non-strongly convex functions is  $\tilde{O}\left(\frac{1}{\epsilon} \log\left(\frac{1}{\epsilon}\right)\right)$ , and the proof is similar to the one in (Zhang and Xiao, 2015b).

## 6 NUMERICAL EXPERIMENTS

In this section, we conduct numerical experiments of DiSPA and A-DiSPA on several real-world distributed datasets. We compare the performance of our algorithms with the state-of-the-art distributed optimization algorithm CoCoA+ (Ma et al., 2015).

For our experiments, we solve the standard binary classification task with datasets obtained from LIBSVM datasets (Chang and Lin, 2011). The statistics of the datasets are summarized in Table 1. Our goal is to minimize the regularized empirical risk with smoothed hinge loss:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{P}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \phi_i(a_i^\top \mathbf{x}) + \frac{\lambda}{2} \|\mathbf{x}\|_2^2.$$

For each task, the data point takes the form of  $(a_i, b_i)$ , where  $a_i$  is a feature vector,  $b_i \in \{1, -1\}$  is the corresponding class label and associated with loss function  $\phi_i$ . The smoothed hinge loss function (Shalev-Shwartz and Zhang, 2013)  $\phi_i$  is defined as

$$\phi_i(z) = \begin{cases} 0 & \text{if } b_i z \geq 1, \\ 1 - b_i z - \frac{\gamma}{2} & \text{if } b_i z \leq 1 - \gamma, \\ \frac{1}{2\gamma}(1 - b_i z)^2 & \text{otherwise,} \end{cases}$$

where  $\phi_i$  is  $(1/\gamma)$ -smooth. We set  $\gamma = 1$  in our experiments. The conjugate function of  $\phi_i$  is  $\phi_i^*(\beta) = b_i \beta + \frac{\gamma}{2} \beta^2$  for  $b_i \beta \in [-1, 0]$  and  $+\infty$  otherwise.

Table 1: Three Datasets for Numerical Experiments

DATASET	# SAMPLES $n$	# FEATURES $d$
RCV1	677,399	47,236
Realsim	72,309	20,958
Covtype	581,012	54

### 6.1 IMPLEMENTATION DETAILS

We implement DiSPA, A-DiSPA and CoCoA+ in Petuum (Xing et al., 2015). For DiSPA and A-DiSPA, we use SPDC (Zhang and Xiao, 2015b) as the local solver, and use SDCA (Shalev-Shwartz and Zhang, 2013) as the local solver for CoCoA+. When doing comparison among different algorithms, we run the same number of local iterations on local machines (i.e., iterations of SDCA or SPDC) before communication with the central node. We compare different algorithms based on communication iterations. From our results, we find that DiSPA converges faster when parameters  $\tau$  and  $\sigma$  are large, which is similar to the findings reported in (Johnson and Zhang, 2013, Zhang and Xiao, 2015a). We adopt

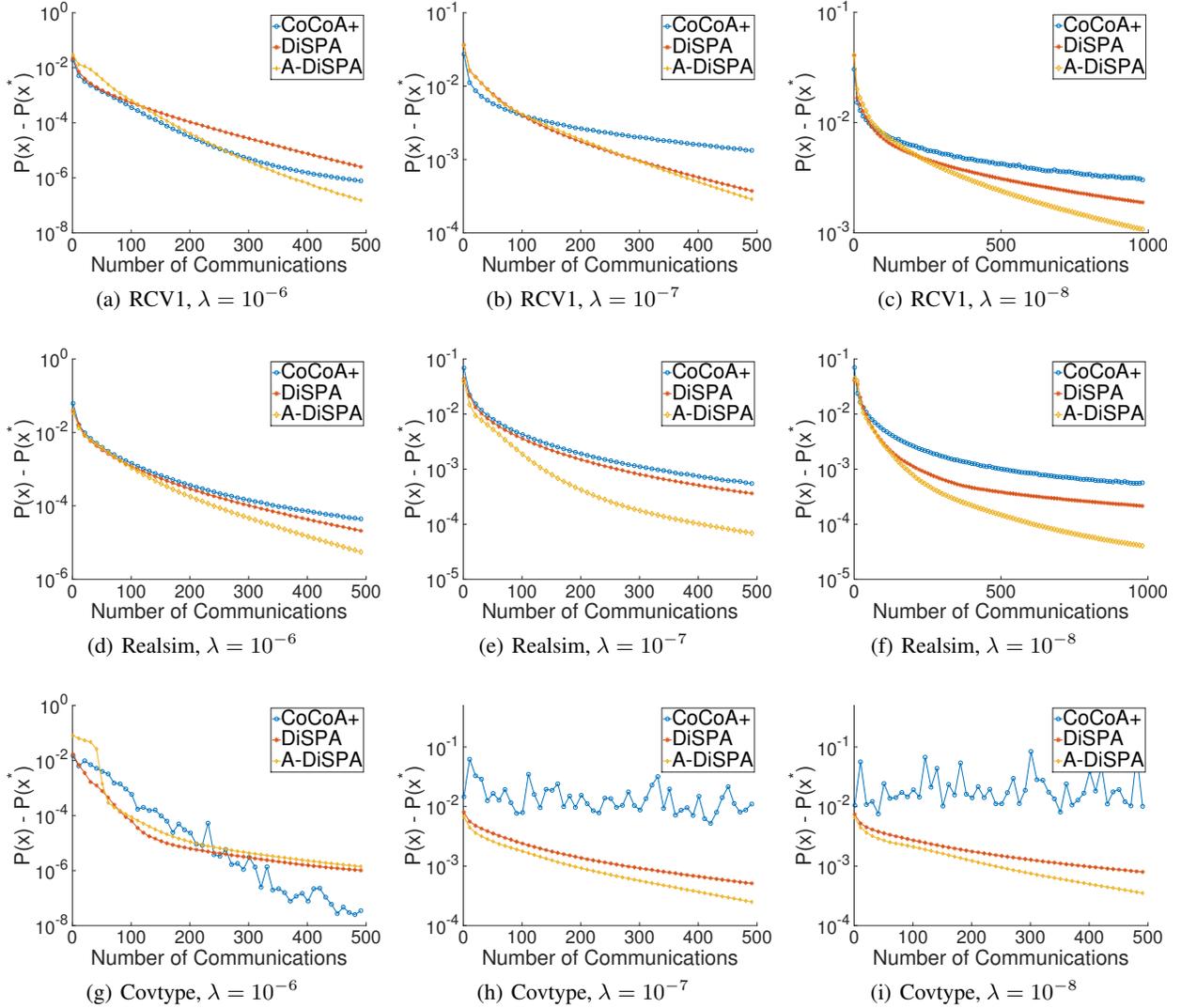


Figure 1: Comparison of DiSPA, A-DiSPA and CoCoA+ on three datasets: RCV1, Realsim, Covtype. The horizontal axis is the number of communications, and the vertical axis is the optimality gap ( $P(x^{(T)}) - P(x^*)$ ). Each plot contains a comparison of CoCoA+ (blue, circle), DiSPA (red, asterisk) and A-DiSPA (yellow, diamond). All the plots are shown on log-y scale. The number of local iterations performed in each communication iteration is 5,000 for Realsim, and 10,000 iterations for RCV1 and Covtype.

tuned parameters  $\tau$  and  $\sigma$  from a predefined range. For CoCoA+, we present the results of the selected optimal parameter  $\sigma'$ .

## 6.2 COMPARISON BETWEEN DISPA AND COCOA+

We compare DiSPA and CoCoA+ on three datasets: RCV1, Realsim, Covtype, across different values of the regularization parameter  $\lambda$ . To be specific, the regularization parameter  $\lambda$  is set to be  $\lambda \in \{10^{-6}, 10^{-7}, 10^{-8}\}$

on the three datasets. Both DiSPA and CoCoA+ are implemented on 16 machines ( $K = 16$ ). For both DiSPA and CoCoA+, the numbers of local iterations performed in each communication iteration are the same.

In Figure 1, the results show that the performance of DiSPA is comparable with CoCoA+. We can see that when  $\lambda = 10^{-6}$ , DiSPA has comparable convergence performance with regard to CoCoA+. On all the three datasets, when  $\lambda \in \{10^{-7}, 10^{-8}\}$ , DiSPA converges to the optimal solution faster and more stable than CoCoA+. We also notice that the performance of Co-

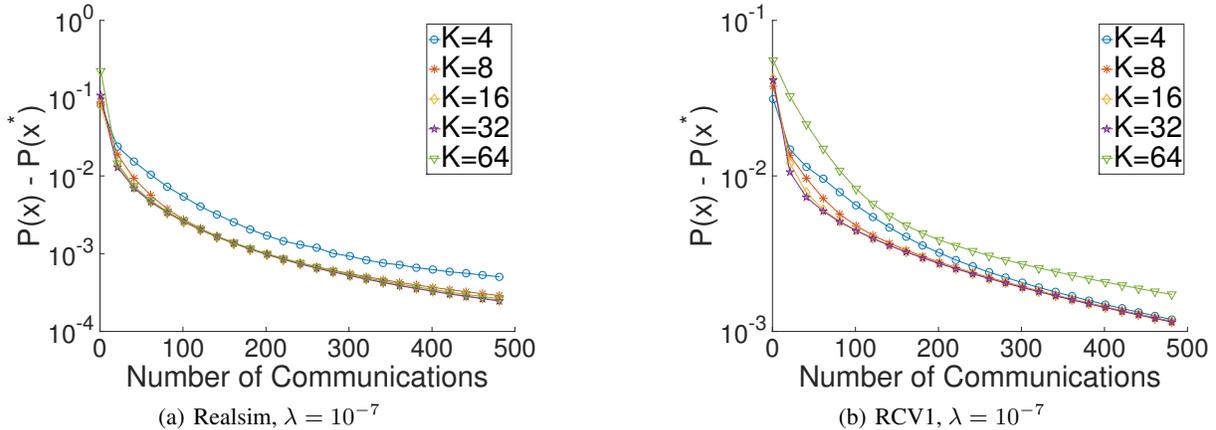


Figure 2: The performance of DiSPA with increasing number of machines,  $K \in \{4, 8, 16, 32, 64\}$ . The local iterations performed in each communication iteration is 1,000 for Realsim, and 5,000 iterations for RCV1.

CoA+ on Covtype is not stable compared to that on the other two datasets, especially when  $\lambda$  is very small, i.e.,  $\lambda = 10^{-7}$  or  $\lambda = 10^{-8}$ , which is consistent with the experimental results of CoCoA+ on Covtype in Reddi et al. (2016). In Section 6.3, we will show the superior performance of A-DiSPA on ill-conditioned problems.

### 6.3 COMPARISON OF DISPA AND A-DISPA

In Figure 1, we also compare A-DiSPA and DiSPA to show the acceleration obtained by applying Catalyst acceleration. We compare DiSPA and A-DiSPA on three datasets (RCV1, Realsim, Covtype) across different regularization parameter under the same settings as mentioned in Section 6.2. The results show that when problem is ill-conditioned, i.e.  $\kappa = R^2/(\lambda\gamma) \gg 1$ , A-DiSPA converges substantially faster than DiSPA. This confirms our theoretical analysis, as the communication complexity of A-DiSPA is  $\tilde{O}(\sqrt{\kappa} \log(\frac{1}{\epsilon}))$  and DiSPA is  $\tilde{O}(\kappa \log(\frac{1}{\epsilon}))$ .

### 6.4 SCALABILITY

Since scalability is an important metric of distributed optimization algorithms, in Figure 2, we study the scalability of DiSPA by observing the numerical performance with increasing number of machines. We conduct experiments on two datasets (RCV1, Realsim), and the machine number is set to be  $K \in \{4, 8, 16, 32, 64\}$ . For comparison, we keep the iterations and parameters  $\tau, \sigma$  the same for each dataset. The experiments show that DiSPA can scale effectively with number of machines, which confirms our theory in Section 4. We observe that performance slightly drops on RCV1 with 64 machines,

this is possibly because that local subproblems are solved with higher accuracy which may affect the effectiveness of aggregation in the central update.

## 7 CONCLUSION

In this work, we present a novel distributed optimization framework for solving the saddle point problem, which can be applied to solving a generic class of convex optimization problems. We provide the theoretical guarantee of our algorithms, and show that the accelerated algorithm can achieve the optimal communication complexity. We also extend our algorithm for solving non-smooth and non-strongly convex loss functions. The experimental results demonstrate that our algorithms obtain better performance compared to the state-of-the-art distributed optimization method.

### Acknowledgements

This work is supported by NTU Singapore Nanyang Assistant Professorship (NAP) grant M4081532.020, Singapore MOE AcRF Tier-2 grant MOE2016-T2-2-060, and MOE AcRF Tier-1 grant 2016-T1-001-159.

### References

- Y. Arjevani and O. Shamir. Communication complexity of distributed convex learning and optimization. In *Advances in Neural Information Processing Systems*, pages 1756–1764, 2015.
- S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers.

- Foundations and Trends® in Machine Learning*, 3(1): 1–122, 2011.
- A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 40(1):120–145, 2011.
- C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- T.-H. Chang, A. Nedic, and A. Scaglione. Distributed constrained optimization by consensus-based primal-dual perturbation method. *IEEE Transactions on Automatic Control*, 59(6):1524–1538, 2014.
- B. Dai, N. He, Y. Pan, B. Boots, and L. Song. Learning from conditional distributions via dual kernel embeddings. *arXiv preprint arXiv:1607.04579*, 2016.
- J. Duchi, M. I. Jordan, and B. McMahan. Estimation, optimization, and parallelism when data is sparse. In *Advances in Neural Information Processing Systems*, pages 2832–2840, 2013.
- M. Jaggi, V. Smith, M. Takáč, J. Terhorst, S. Krishnan, T. Hofmann, and M. I. Jordan. Communication-efficient distributed dual coordinate ascent. In *Advances in Neural Information Processing Systems*, pages 3068–3076, 2014.
- R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
- H. Lin, J. Mairal, and Z. Harchaoui. A universal catalyst for first-order optimization. In *Advances in Neural Information Processing Systems*, pages 3384–3392, 2015.
- C. Ma, V. Smith, M. Jaggi, M. Jordan, P. Richtárik, and M. Takac. Adding vs. averaging in distributed primal-dual optimization. In *Proceedings of The 32nd International Conference on Machine Learning*, pages 1973–1982, 2015.
- Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- S. J. Reddi, J. Konečný, P. Richtárik, B. Póczós, and A. Smola. Aide: Fast and communication efficient distributed optimization. *arXiv preprint arXiv:1608.06879*, 2016.
- S. Shalev-Shwartz and T. Zhang. Stochastic dual coordinate ascent methods for regularized loss minimization. *Journal of Machine Learning Research*, 14(Feb):567–599, 2013.
- O. Shamir, N. Srebro, and T. Zhang. Communication-efficient distributed optimization using an approximate newton-type method. In *Proceedings of The 31st International Conference on Machine Learning*, pages 1000–1008, 2014.
- W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin. On the linear convergence of the admm in decentralized consensus optimization. *IEEE Transactions on Signal Processing*, 62(7):1750–1761, 2014.
- V. Smith, S. Forte, M. I. Jordan, and M. Jaggi. L1-regularized distributed optimization: A communication-efficient primal-dual framework. *arXiv preprint arXiv:1512.04011*, 2015.
- E. P. Xing, Q. Ho, W. Dai, J. K. Kim, J. Wei, S. Lee, X. Zheng, P. Xie, A. Kumar, and Y. Yu. Petuum: A new platform for distributed machine learning on big data. *IEEE Transactions on Big Data*, 1(2):49–67, 2015.
- T. Yang. Trading computation for communication: Distributed stochastic dual coordinate ascent. In *Advances in Neural Information Processing Systems*, pages 629–637, 2013.
- A. W. Yu, Q. Lin, and T. Yang. Doubly stochastic primal-dual coordinate method for regularized empirical risk minimization with factorized data. *CoRR*, abs/1508.03390, 2015.
- Y. Zhang and L. Xiao. Disco: Distributed optimization for self-concordant empirical loss. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 362–370, 2015a.
- Y. Zhang and L. Xiao. Stochastic primal-dual coordinate method for regularized empirical risk minimization. In *Proceedings of The 32nd International Conference on Machine Learning*, pages 353–361, 2015b.
- Y. Zhang, M. J. Wainwright, and J. C. Duchi. Communication-efficient algorithms for statistical optimization. In *Advances in Neural Information Processing Systems*, pages 1502–1510, 2012.

## Appendix

### A Lemmas for Convergence Analysis of DiSPA

We first introduce the lemma that characterizes the optimality condition of local subproblem  $\mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]})$ .

**Lemma 2.** *Assume that each  $\phi_i(\cdot)$  is  $(1/\gamma)$ -smooth and  $g(\cdot)$  is  $\lambda$ -strongly convex,  $R = \max\{\|a_1\|_2, \dots, \|a_n\|_2\}$ . Let  $(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)})$  be the optimal solution of  $\mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]})$ ,  $k = 1, 2, \dots, K$ . Based on the strong convexity, we have*

$$\mathcal{L}_k^{(t)}(\mathbf{x}, \hat{\mathbf{y}}_{[k]}^{(t)}) \geq \mathcal{L}_k^{(t)}(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) + \left(\frac{1}{2\tau} + \frac{\lambda}{2}\right) \|\mathbf{x} - \hat{\mathbf{x}}_k^{(t)}\|^2, \quad \forall \mathbf{x} \in \mathbb{R}^d$$

$$\mathcal{L}_k^{(t)}(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]}) \leq \mathcal{L}_k^{(t)}(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) - \left(\frac{K}{2\sigma n} + \frac{\gamma}{2n}\right) \|\mathbf{y}_{[k]} - \hat{\mathbf{y}}_{[k]}^{(t)}\|^2, \quad \forall \mathbf{y}_{[k]} \in \mathbb{R}^n$$

*Proof.* Based on the definition of the saddle point, we can notice that  $-\mathcal{L}_k^{(t)}(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]})$  is a  $(\frac{K}{\sigma n} + \frac{\gamma}{n})$ -strongly convex function and minimized by  $\hat{\mathbf{y}}_{[k]}^{(t)}$ , which implies

$$-\mathcal{L}_k^{(t)}(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]}) \geq -\mathcal{L}_k^{(t)}(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) + \left(\frac{K}{2\sigma n} + \frac{\gamma}{2n}\right) \|\mathbf{y}_{[k]} - \hat{\mathbf{y}}_{[k]}^{(t)}\|^2, \quad \forall \mathbf{y}_{[k]} \in \mathbb{R}^n$$

Also notice that  $\mathcal{L}_k^{(t)}(\mathbf{x}, \hat{\mathbf{y}}_{[k]}^{(t)})$  is a  $(\frac{1}{\tau} + \lambda)$ -strongly convex function minimized by  $\hat{\mathbf{x}}_k^{(t)}$ , which implies

$$\mathcal{L}_k^{(t)}(\mathbf{x}, \hat{\mathbf{y}}_{[k]}^{(t)}) \geq \mathcal{L}_k^{(t)}(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) + \left(\frac{1}{2\tau} + \frac{\lambda}{2}\right) \|\mathbf{x} - \hat{\mathbf{x}}_k^{(t)}\|^2, \quad \forall \mathbf{x} \in \mathbb{R}^d$$

□

Based on the optimality condition of  $\mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]})$  and the central update  $\mathbf{x}^{(t)}$  on central worker, we can get the connection between local subproblems and central update on central worker.

**Lemma 3** (Relationship between local optimal solution and central update). *Assume that each  $\phi_i(\cdot)$  is  $(1/\gamma)$ -smooth and  $g(\cdot)$  is  $\lambda$ -strongly convex,  $R = \max\{\|a_1\|_2, \dots, \|a_n\|_2\}$ . Let  $(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)})$  be the optimal solution of  $\mathcal{L}_k^{(t)}(\mathbf{x}, \mathbf{y}_{[k]})$ ,  $k = 1, 2, \dots, K$ . Let  $\mathbf{x}^{(t)} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{C}^{(t)}(\mathbf{x})$ , it holds that*

$$\frac{(K-1)^2}{4\sigma n K} \|\mathbf{y}^{(t-1)} - \mathbf{y}^{(t)}\|_2^2 + \mathbf{\Lambda}' \geq \left(\frac{3}{4\tau} + \lambda\right) \sum_{k=1}^K \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2 \quad (9)$$

where  $\mathbf{\Lambda}' = \frac{1}{n} \sum_{k=1}^K (\hat{\mathbf{y}}_{[k]}^{(t)} - \mathbf{y}_{[k]}^{(t)})^\top A^\top (\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)})$  and  $1/\tau = 4\sigma R^2$ .

*Proof.* Based on lemma 2, we could get that for  $k = 1, 2, \dots, K$ ,

$$\begin{aligned} & \frac{1}{n} \sum_{j \in \mathcal{P}_k} \hat{y}_j^{(t)} \langle a_j, \mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)} \rangle + \frac{1}{n} \sum_{i \notin \mathcal{P}_k} y_i^{(t-1)} \langle a_i, \mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)} \rangle + g(\mathbf{x}^{(t)}) - g(\hat{\mathbf{x}}_k^{(t)}) + \frac{\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} \geq \\ & \frac{\|\hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \left(\frac{1}{2\tau} + \frac{\lambda}{2}\right) \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2 \end{aligned}$$

Since  $\mathbf{x}^{(t)}$  minimizes the function  $\mathcal{C}^{(t)}(\mathbf{x})$ , which means that for  $k = 1, 2, \dots, K$ , we have

$$\frac{1}{n} \sum_{i=1}^n y_i^{(t)} \langle a_i, \hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)} \rangle + g(\hat{\mathbf{x}}_k^{(t)}) - g(\mathbf{x}^{(t)}) + \frac{\|\hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} \geq \frac{\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \left(\frac{1}{2\tau} + \frac{\lambda}{2}\right) \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2$$

Sum up the above two inequalities, we can get

$$\begin{aligned}
& \frac{1}{n} \sum_{j \in \mathcal{P}_k} \hat{y}_j^{(t)} \langle a_j, \mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)} \rangle + \frac{1}{n} \sum_{i \notin \mathcal{P}_k} y_i^{(t-1)} \langle a_i, \mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)} \rangle + \frac{1}{n} \sum_{i=1}^n y_i^{(t)} \langle a_i, \hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)} \rangle \geq \left( \frac{1}{\tau} + \lambda \right) \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2 \\
& \frac{1}{n} \sum_{j \in \mathcal{P}_k} \left( \hat{y}_j^{(t)} - y_j^{(t)} \right) \langle a_j, \mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)} \rangle + \frac{1}{n} \sum_{i \notin \mathcal{P}_k} \left( y_i^{(t-1)} - y_i^{(t)} \right) \langle a_i, \mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)} \rangle \geq \left( \frac{1}{\tau} + \lambda \right) \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2 \\
& \frac{1}{n} (\hat{\mathbf{y}}_{[k]}^{(t)} - \mathbf{y}_{[k]}^{(t)})^\top A^\top (\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}) + \frac{1}{n} \sum_{l \neq k} (\mathbf{y}_{[l]}^{(t-1)} - \mathbf{y}_{[l]}^{(t)})^\top A^\top (\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}) \geq \left( \frac{1}{\tau} + \lambda \right) \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2
\end{aligned}$$

We need to upper bound the second term on the left-hand-side of the above inequality. Since  $\|a_i\|_2 \leq R$ , and we assume that  $1/\tau = 4\sigma R^2$ , then

$$\begin{aligned}
\left| \frac{1}{n} \sum_{l \neq k} (\mathbf{y}_{[l]}^{(t-1)} - \mathbf{y}_{[l]}^{(t)})^\top A^\top (\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}) \right| & \leq \frac{\|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2}{4\tau} + \frac{\|\sum_{l \neq k} A(\mathbf{y}_{[l]}^{(t-1)} - \mathbf{y}_{[l]}^{(t)})\|_2^2}{n^2/\tau} \\
& \leq \frac{\|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2}{4\tau} + \frac{(\sum_{i \notin \mathcal{P}_k} |y_i^{(t-1)} - y_i^{(t)}| \cdot \|a_i\|)^2}{4\sigma R^2 n^2} \\
& \leq \frac{\|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2}{4\tau} + \frac{(K-1)}{4\sigma n K} \sum_{l \neq k} \|\mathbf{y}_{[l]}^{(t-1)} - \mathbf{y}_{[l]}^{(t)}\|_2^2
\end{aligned}$$

Combine the above two inequalities,

$$\frac{(K-1)}{4\sigma n K} \sum_{l \neq k} \|\mathbf{y}_{[l]}^{(t-1)} - \mathbf{y}_{[l]}^{(t)}\|_2^2 + \frac{1}{n} (\hat{\mathbf{y}}_{[k]}^{(t)} - \mathbf{y}_{[k]}^{(t)})^\top A^\top (\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}) \geq \left( \frac{3}{4\tau} + \lambda \right) \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2$$

Sum up the above inequality over  $k = 1, 2, \dots, K$ , and denote  $\mathbf{\Lambda}' = \frac{1}{n} \sum_{k=1}^K (\hat{\mathbf{y}}_{[k]}^{(t)} - \mathbf{y}_{[k]}^{(t)})^\top A^\top (\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)})$ , we have

$$\frac{(K-1)^2}{4\sigma n K} \|\mathbf{y}^{(t-1)} - \mathbf{y}^{(t)}\|_2^2 + \mathbf{\Lambda}' \geq \left( \frac{3}{4\tau} + \lambda \right) \sum_{k=1}^K \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2$$

□

Lemma 3 shows that the distance between  $\mathbf{x}^{(t)}$  and  $\hat{\mathbf{x}}_k^{(t)}$  can be control by the update of  $\mathbf{y}^{(t)}$  in each iteration.

## B Proofs of Convergence for DiSPA and A-DiSPA

### B.1 Proof of Lemma 1

*Proof.* We start from characterizing the relationship between  $\mathbf{x}^{(t)}$  and  $\mathbf{x}^*$  after the  $t$ -update in Algorithm 1. According to the definition of  $\mathbf{x}^{(t)}$ , we have  $\mathcal{C}^{(t)}(\mathbf{x}^*) \geq \mathcal{C}^{(t)}(\mathbf{x}^{(t)}) + \left( \frac{1}{2\tau} + \frac{\lambda}{2} \right) \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2$ , i.e.,

$$\frac{1}{n} \sum_{i=1}^n y_i^{(t)} \langle a_i, \mathbf{x}^* - \mathbf{x}^{(t)} \rangle + g(\mathbf{x}^*) - g(\mathbf{x}^{(t)}) + \frac{\|\mathbf{x}^* - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} \geq \frac{\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \left( \frac{1}{2\tau} + \frac{\lambda}{2} \right) \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}^*\|_2^2 \quad (10)$$

We also could derive the inequality characterizing the relation between  $\hat{\mathbf{y}}^{(t)} = \sum_{k=1}^K \hat{\mathbf{y}}_{[k]}^{(t)}$  and  $\mathbf{y}^*$ . For  $k = 1, 2, \dots, K$ , we have

$$\hat{\mathbf{y}}_{[k]}^{(t)} = \arg \max_{\mathbf{y}_{[k]} \in \mathbb{R}^n} \mathcal{L}_k^{(t)}(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]})$$

Since  $\phi_i(\cdot)$  is  $(1/\gamma)$ -smooth, we have  $\phi_i^*(\cdot)$  is  $\gamma$ -strongly convex, therefore, for  $k = 1, 2, \dots, K$ ,  $\forall j \in \mathcal{P}_k$

$$(\hat{y}_j^{(t)} - y_j^*) \langle a_j, \hat{\mathbf{x}}_k^{(t)} \rangle + \left( \phi_j^*(y_j^*) - \phi_j^*(\hat{y}_j^{(t)}) \right) + \frac{K(y_j^* - y_j^{(t-1)})^2}{2\sigma} \geq \frac{K(\hat{y}_j^{(t)} - y_j^{(t-1)})^2}{2\sigma} + \left( \frac{K}{2\sigma} + \frac{\gamma}{2} \right) (\hat{y}_j^{(t)} - y_j^*)^2$$

Sum up the above inequality over  $j$ , we have

$$\sum_{j \in \mathcal{P}_k} (\hat{y}_j^{(t)} - y_j^*) \langle a_j, \hat{\mathbf{x}}_k^{(t)} \rangle + \sum_{j \in \mathcal{P}_k} (\phi_j^*(y_j^*) - \phi_j^*(\hat{y}_j^{(t)})) + \frac{K \|\mathbf{y}_{[k]}^* - \mathbf{y}_{[k]}^{(t-1)}\|_2^2}{2\sigma} \geq \frac{K \|\hat{\mathbf{y}}_{[k]}^{(t)} - \mathbf{y}_{[k]}^{(t-1)}\|_2^2}{2\sigma} + \left( \frac{K}{2\sigma} + \frac{\gamma}{2} \right) \|\mathbf{y}_{[k]}^* - \hat{\mathbf{y}}_{[k]}^{(t)}\|_2^2$$

Sum up the above inequality over  $k = 1, 2, \dots, K$  and multiplying both sides by  $1/n$ ,

$$\begin{aligned} & \frac{1}{n} \sum_{k=1}^K \sum_{j \in \mathcal{P}_k} (\hat{y}_j^{(t)} - y_j^*) \langle a_j, \hat{\mathbf{x}}_k^{(t)} \rangle + \frac{1}{n} \sum_{i=1}^n (\phi_i^*(y_i^*) - \phi_i^*(\hat{y}_i^{(t)})) + \frac{K \|\mathbf{y}^* - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} \\ & \geq \frac{K \|\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} + \left( \frac{K}{2\sigma n} + \frac{\gamma}{2n} \right) \|\mathbf{y}^* - \hat{\mathbf{y}}^{(t)}\|_2^2 \end{aligned} \quad (11)$$

In addition, we consider a combination of the saddle-point function values at different points, we have

$$\begin{aligned} & (f(\mathbf{x}^{(t)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^*)) + (f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t)})) \\ & = \left( \frac{1}{n} \sum_{i=1}^n y_i^* \langle a_i, \mathbf{x}^{(t)} \rangle - \frac{1}{n} \sum_{i=1}^n \phi_i^*(y_i^*) + g(\mathbf{x}^{(t)}) \right) - \left( \frac{1}{n} \sum_{i=1}^n y_i^{(t)} \langle a_i, \mathbf{x}^* \rangle - \frac{1}{n} \sum_{i=1}^n \phi_i^*(y_i^{(t)}) + g(\mathbf{x}^*) \right) \\ & = (g(\mathbf{x}^{(t)}) - g(\mathbf{x}^*)) + \frac{1}{n} \sum_{i=1}^n (\phi_i^*(y_i^{(t)}) - \phi_i^*(y_i^*)) + \left( \frac{1}{n} \sum_{i=1}^n y_i^* \langle a_i, \mathbf{x}^{(t)} \rangle - \frac{1}{n} \sum_{i=1}^n y_i^{(t)} \langle a_i, \mathbf{x}^* \rangle \right) \end{aligned}$$

Then we add (10) and (11) to the above equality, which implies

$$\begin{aligned} & \frac{1}{n} \sum_{k=1}^K \sum_{j \in \mathcal{P}_k} (\hat{y}_j^{(t)} - y_j^*) \langle a_j, \hat{\mathbf{x}}_k \rangle + \frac{1}{n} \sum_{i=1}^n y_i^{(t)} \langle a_i, \mathbf{x}^* - \mathbf{x}^{(t)} \rangle + \frac{1}{n} \sum_{i=1}^n y_i^* \langle a_i, \mathbf{x}^{(t)} \rangle - \frac{1}{n} \sum_{i=1}^n y_i^{(t)} \langle a_i, \mathbf{x}^* \rangle \\ & + \frac{1}{n} \sum_{i=1}^n (\phi_i^*(y_i^{(t)}) - \phi_i^*(\hat{y}_i^{(t)})) + \frac{\|\mathbf{x}^* - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K \|\mathbf{y}^* - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} \\ & \geq (f(\mathbf{x}^{(t)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^*)) + (f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t)})) \\ & + \frac{\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K \|\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} + \left( \frac{1}{2\tau} + \frac{\lambda}{2} \right) \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 + \left( \frac{K}{2\sigma n} + \frac{\gamma}{2n} \right) \|\mathbf{y}^* - \hat{\mathbf{y}}^{(t)}\|_2^2 \end{aligned}$$

which implies that

$$\begin{aligned} & \frac{1}{n} \sum_{k=1}^K \sum_{j \in \mathcal{P}_k} (\hat{y}_j^{(t)} - y_j^*) \langle a_j, \hat{\mathbf{x}}_k - \mathbf{x}^{(t)} \rangle + \frac{1}{n} \sum_{i=1}^n (\hat{y}_i^{(t)} - y_i^{(t)}) \langle a_i, \mathbf{x}^{(t)} \rangle + \frac{1}{n} \sum_{i=1}^n (\phi_i^*(y_i^{(t)}) - \phi_i^*(\hat{y}_i^{(t)})) \\ & + \frac{\|\mathbf{x}^* - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K \|\mathbf{y}^* - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} \\ & \geq (f(\mathbf{x}^{(t)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^*)) + (f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t)})) + \frac{\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K \|\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} \\ & + \left( \frac{1}{2\tau} + \frac{\lambda}{2} \right) \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 + \left( \frac{K}{2\sigma n} + \frac{\gamma}{2n} \right) \|\mathbf{y}^* - \hat{\mathbf{y}}^{(t)}\|_2^2 \end{aligned}$$

We need to upper bound the first term on the left-hand-side of the above inequality, assume that  $1/\tau = 4\sigma R^2$ , we have

$$\begin{aligned}
\left| \frac{1}{n} \sum_{j \in \mathcal{P}_k} (\hat{y}_j^{(t)} - y_j^*) \langle a_j, \hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)} \rangle \right| &= \left| \frac{1}{n} (\hat{\mathbf{y}}_{[k]}^{(t)} - \mathbf{y}_{[k]}^*)^\top A^\top (\hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)}) \right| \\
&\leq \frac{\|\hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)}\|_2^2}{4\tau} + \frac{\|A(\hat{\mathbf{y}}_{[k]}^{(t)} - \mathbf{y}_{[k]}^*)\|_2^2}{n^2/\tau} \\
&\leq \frac{\|\hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)}\|_2^2}{4\tau} + \frac{(\sum_{j \in \mathcal{P}_k} |\hat{y}_j^{(t)} - y_j^*| \cdot \|a_j\|)^2}{4\sigma R^2 n^2} \\
&\leq \frac{\|\hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)}\|_2^2}{4\tau} + \frac{\|\hat{\mathbf{y}}_{[k]}^{(t)} - \mathbf{y}_{[k]}^*\|_2^2}{4\sigma n K}
\end{aligned}$$

then we can get the upper bound

$$\frac{1}{n} \sum_{k=1}^K \sum_{j \in \mathcal{P}_k} (\hat{y}_j^{(t)} - y_j^*) \langle a_j, \hat{\mathbf{x}}_k - \mathbf{x}^{(t)} \rangle \leq \sum_{k=1}^K \frac{\|\hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)}\|_2^2}{4\tau} + \frac{\|\hat{\mathbf{y}}^{(t)} - \mathbf{y}^*\|_2^2}{4\sigma n K}$$

next we denote that

$$\mathbf{\Lambda}'' = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i^{(t)} - y_i^*) \langle a_i, \mathbf{x}^{(t)} \rangle + \frac{1}{n} \sum_{i=1}^n (\phi_i^*(y_i^{(t)}) - \phi_i^*(\hat{y}_i))$$

Combining the above inequality and equation, we derive that

$$\begin{aligned}
&\sum_{k=1}^K \frac{\|\hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)}\|_2^2}{4\tau} + \frac{\|\mathbf{x}^* - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K\|\mathbf{y}^* - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} + \mathbf{\Lambda}'' \\
&\geq \left( f(\mathbf{x}^{(t)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^*) \right) + \left( f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t)}) \right) + \frac{\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K\|\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} \\
&\quad + \left( \frac{1}{2\tau} + \frac{\lambda}{2} \right) \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 + \left( \frac{K}{2\sigma n} + \frac{\gamma}{2n} - \frac{1}{4\sigma n K} \right) \|\mathbf{y}^* - \hat{\mathbf{y}}^{(t)}\|_2^2
\end{aligned}$$

Based on the inequality (9) in lemma 3, we could get that

$$\begin{aligned}
&\sum_{k=1}^K \frac{\|\hat{\mathbf{x}}_k^{(t)} - \mathbf{x}^{(t)}\|_2^2}{4\tau} + \frac{(K-1)^2}{4\sigma n K} \|\mathbf{y}^{(t-1)} - \mathbf{y}^{(t)}\|_2^2 + \frac{\|\mathbf{x}^* - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K\|\mathbf{y}^* - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} + \mathbf{\Lambda}' + \mathbf{\Lambda}'' \\
&\geq \left( f(\mathbf{x}^{(t)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^*) \right) + \left( f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t)}) \right) + \frac{\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K\|\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} \\
&\quad + \left( \frac{3}{4\tau} + \lambda \right) \sum_{k=1}^K \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2 + \left( \frac{1}{2\tau} + \frac{\lambda}{2} \right) \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 + \left( \frac{K}{2\sigma n} + \frac{\gamma}{2n} - \frac{1}{4\sigma n K} \right) \|\mathbf{y}^* - \hat{\mathbf{y}}^{(t)}\|_2^2
\end{aligned}$$

Since we can rewrite  $\|\mathbf{y} - \hat{\mathbf{y}}^{(t)}\|_2^2 = \|\mathbf{y} - \mathbf{y}^{(t)}\|_2^2 + \|\mathbf{y}^{(t)} - \hat{\mathbf{y}}^{(t)}\|_2^2 + 2(\mathbf{y} - \mathbf{y}^{(t)})^\top (\mathbf{y}^{(t)} - \hat{\mathbf{y}}^{(t)})$ , then we denote that

$$\begin{aligned}
\mathbf{\Lambda}''' &= - \left( \frac{K}{2\sigma n} + \frac{\gamma}{2n} - \frac{1}{4\sigma n K} \right) \left( \|\mathbf{y}^{(t)} - \hat{\mathbf{y}}^{(t)}\|_2^2 + 2(\mathbf{y}^* - \mathbf{y}^{(t)})^\top (\mathbf{y}^{(t)} - \hat{\mathbf{y}}^{(t)}) \right) \\
&\quad - \frac{K}{2\sigma n} \left( \|\mathbf{y}^{(t)} - \hat{\mathbf{y}}^{(t)}\|_2^2 + 2(\mathbf{y}^{(t-1)} - \mathbf{y}^{(t)})^\top (\mathbf{y}^{(t)} - \hat{\mathbf{y}}^{(t)}) \right)
\end{aligned}$$

Based on the definition of  $\mathbf{\Lambda}'''$ , we derive that

$$\begin{aligned}
&\frac{\|\mathbf{x}^* - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K\|\mathbf{y}^* - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} + \mathbf{\Lambda}' + \mathbf{\Lambda}'' + \mathbf{\Lambda}''' \\
&\geq \left( f(\mathbf{x}^{(t)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^*) \right) + \left( f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t)}) \right) \\
&\quad + \left( \frac{1}{2\tau} + \frac{\lambda}{2} \right) \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 + \left( \frac{K}{2\sigma n} + \frac{\gamma}{2n} - \frac{1}{4\sigma n K} \right) \|\mathbf{y}^* - \mathbf{y}^{(t)}\|_2^2 + \mathbf{\Lambda}^\#
\end{aligned}$$

where we define

$$\mathbf{\Lambda}^\sharp = \frac{\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \left(\frac{1}{2\tau} + \lambda\right) \sum_{k=1}^K \|\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}\|_2^2 + \left(\frac{K}{2\sigma n} - \frac{(K-1)^2}{4\sigma n K}\right) \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|_2^2$$

Assume that  $\Theta \geq 0$ , then we can obtain

$$\begin{aligned} & \Theta \left( f(\mathbf{x}^{(t-1)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^*) \right) + \Theta \left( f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)}) \right) \\ & \frac{\|\mathbf{x}^* - \mathbf{x}^{(t-1)}\|_2^2}{2\tau} + \frac{K\|\mathbf{y}^* - \mathbf{y}^{(t-1)}\|_2^2}{2\sigma n} + \mathbf{\Lambda} \\ & \geq \left( f(\mathbf{x}^{(t)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^*) \right) + \left( f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t)}) \right) + \Theta \left( f(\mathbf{x}^{(t-1)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)}) \right) \\ & + \left( \frac{1}{2\tau} + \frac{\lambda}{2} \right) \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 + \left( \frac{K}{2\sigma n} + \frac{\gamma}{2n} - \frac{1}{4\sigma n K} \right) \|\mathbf{y}^* - \mathbf{y}^{(t)}\|_2^2 + \mathbf{\Lambda}^\sharp \end{aligned}$$

where  $\mathbf{\Lambda} = \mathbf{\Lambda}' + \mathbf{\Lambda}'' + \mathbf{\Lambda}'''$  and  $\mathbf{\Lambda}^\sharp > 0$ .

In order to get the convergence guarantee of our algorithm, we need to get the upper bound of  $\mathbf{\Lambda}$ , based on the definition of  $\mathbf{\Lambda}$ , we could get

$$\begin{aligned} \mathbf{\Lambda} & \leq \frac{1}{n} \sum_{k=1}^K (\hat{\mathbf{y}}_{[k]}^{(t)} - \mathbf{y}_{[k]}^{(t)})^\top A^\top (\mathbf{x}^{(t)} - \hat{\mathbf{x}}_k^{(t)}) + \frac{1}{n} (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)})^\top A^\top \mathbf{x}^{(t)} + \frac{1}{n} \sum_{i=1}^n \left( \phi_i^*(y_i^{(t)}) - \phi_i^*(\hat{y}_i) \right) \\ & + \left( \frac{K}{\sigma n} + \frac{\gamma}{n} - \frac{1}{2\sigma n K} \right) (\mathbf{y}^* - \mathbf{y}^{(t)})^\top (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}) + \frac{K}{\sigma n} (\mathbf{y}^{(t-1)} - \mathbf{y}^{(t)})^\top (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}) \end{aligned}$$

based on the assumption, then we could get that

$$\mathbf{\Lambda} \leq \frac{M}{n} \|\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}\|_2 + \sum_{k=1}^K \left( f(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) - f(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]}^{(t)}) \right)$$

where  $M = 4\sqrt{n}R\Omega_{\mathbf{x}} + \left(\frac{3K}{\sigma} + \gamma\right)\Omega_{\mathbf{y}} = c'\Omega_{\mathbf{x}} + c''\Omega_{\mathbf{y}}$  as defined in Assumption 5. Based on the assumption

$$\mathbb{E}[\mathbf{\Lambda}] \leq \hat{\Theta} \frac{M}{n} \|\mathbf{y}^{(t-1)} - \hat{\mathbf{y}}^{(t)}\|_2 + \tilde{\Theta} \sum_{k=1}^K \left( f(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) - f(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]}^{(t-1)}) \right)$$

where  $\hat{\Theta} < 1$ ,  $\tilde{\Theta} < 1$  are defined in Assumption 4, and define the parameter  $\Theta$  as

$$\Theta = \max \left\{ \frac{1}{1 + \left(\frac{\sigma\gamma}{K} - \frac{1}{2K^2}\right)}, \frac{1}{1 + \tau\lambda} \right\} \quad (12)$$

based on Assumption 5, and we could get that

$$\mathbb{E}[\mathbf{\Delta}] = \mathbb{E} \left[ \frac{M}{n} \|\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}\|_2 + \sum_{k=1}^K \left( f(\hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{y}}_{[k]}^{(t)}) - f(\hat{\mathbf{x}}_k^{(t)}, \mathbf{y}_{[k]}^{(t)}) \right) \right] \leq \Theta \left( f(\mathbf{x}^{(t-1)}, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)}) \right)$$

Then we could get our final conclusion

$$\mathbb{E}[\Delta^{(t)}] \leq \Theta \mathbb{E}[\Delta^{(t-1)}]$$

where we require that

$$\sigma > \frac{1}{2K\gamma}, \quad \tau\sigma = \frac{1}{4R^2} \quad (13)$$

□

## B.2 Proof of Theorem 1

*Proof.* By Lemma 1, for each  $t > 0$ , we have  $\mathbb{E}[\Delta^{(t)}] \leq \Theta^t \mathbb{E}[\Delta^{(0)}]$ , according to the definition of  $C$ , we have

$$\mathbb{E}[\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2] \leq \Theta^t C \quad (14)$$

In order to obtain  $\Theta^T C \leq \epsilon$ , which needs

$$T \geq -\frac{\log(C/\epsilon)}{\log(\Theta)}$$

Suppose the parameters  $\tau, \sigma$  are set as

$$\sigma = \frac{1}{\gamma}, \quad \tau = \frac{\gamma}{4R^2}$$

Then according to the definition of  $\Theta$

$$\Theta = \max \left\{ \frac{1}{1 + \left(\frac{\sigma\gamma}{K} - \frac{1}{2K^2}\right)}, \frac{1}{1 + \tau\lambda} \right\} = \max \left\{ \frac{1}{1 + \left(\frac{1}{K} - \frac{1}{2K^2}\right)}, \frac{1}{1 + (\lambda\gamma)/(4R^2)} \right\}$$

As long as  $\frac{1}{K} - \frac{1}{2K^2} \geq \frac{1}{K} - \frac{1}{2K} = \frac{1}{2K} \geq \frac{\lambda\gamma}{4R^2}$ , since we assume  $K \leq \kappa$ . Then we have

$$\Theta = \frac{1}{1 + (\lambda\gamma)/(4R^2)} = 1 - \frac{1}{1 + (4R^2)/(\lambda\gamma)}$$

Thus we can get

$$T \geq -\frac{\log(C/\epsilon)}{\log(\Theta)} = \frac{\log(C/\epsilon)}{-\log\left(1 - \frac{1}{1 + (4R^2)/(\lambda\gamma)}\right)}$$

where we apply the inequality  $\log(1 - x) \leq -x$  in the last inequality and we could get the conclusion.  $\square$

## B.3 Convergence guarantee of primal-dual gap

Next we derive the convergence rate of primal-dual gap based on Theorem 1.

**Lemma 4** (Yu et al. (2015)). *Suppose Assumption 1,2 holds, Let  $(\mathbf{x}^*, \mathbf{y}^*)$  is the unique saddle-point of  $f(\mathbf{x}, \mathbf{y})$ , and  $R = \max_{1 \leq i \leq n} \|a_i\|_2$ . Then for any point  $(\mathbf{x}, \mathbf{y}) \in \text{dom}(g) \times \text{dom}(\phi^*)$ , we have*

$$\mathcal{P}(\mathbf{x}) \leq f(\mathbf{x}, \mathbf{y}^*) + \frac{R^2}{2\gamma} \|\mathbf{x} - \mathbf{x}^*\|_2^2, \quad \mathcal{D}(\mathbf{y}) \geq f(\mathbf{x}^*, \mathbf{y}) - \frac{R^2}{2\lambda n} \|\mathbf{y} - \mathbf{y}^*\|_2^2 \quad (15)$$

**Corollary 2.** *Suppose Assumption A holds and the parameters  $\tau, \sigma, \Theta$  are set as in (12), (13). Then the iterates of Algorithm 1 satisfy*

$$\mathbb{E}[\mathcal{P}(\mathbf{x}^{(T)}) - \mathcal{D}(\mathbf{y}^{(T)})] \leq \epsilon$$

it suffices to have the number of communication iterations  $T$  satisfy,

$$T \geq \left(1 + \frac{4R^2}{\lambda\gamma}\right) \log\left(\left(1 + \frac{R^2}{\lambda\gamma}\right) \frac{C}{\epsilon}\right)$$

*Proof.* Based on Lemma 4, we have

$$\begin{aligned} \mathcal{P}(\mathbf{x}^{(t)}) - \mathcal{D}(\mathbf{y}^{(t)}) &= \mathcal{P}(\mathbf{x}^{(t)}) - \mathcal{P}(\mathbf{x}^*) + \mathcal{D}(\mathbf{y}^*) - \mathcal{D}(\mathbf{y}^{(t)}) \\ &\leq f(\mathbf{x}^{(t)}, \mathbf{y}^*) + \frac{R^2}{2\gamma} \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 - f(\mathbf{x}^*, \mathbf{y}^{(t)}) + \frac{R^2}{2\lambda n} \|\mathbf{y}^{(t)} - \mathbf{y}^*\|_2^2 \\ &\leq \left(1 + \frac{R^2}{\lambda\gamma}\right) \Delta^{(t)} \end{aligned}$$

Similar to Theorem 1, we could get the conclusion.  $\square$

#### B.4 Proof of Theorem 2

*Proof.* If the parameters  $\tau, \sigma$  in Algorithm 2 are set as

$$\tau = \frac{\gamma}{4R^2}, \quad \sigma = \frac{1}{\gamma}$$

then based on Lemma 1, we have

$$\Theta = \max \left\{ \frac{1}{1 + \left(\frac{\sigma\gamma}{K} - \frac{1}{2K^2}\right)}, \frac{1}{1 + \tau\lambda} \right\} = \max \left\{ \frac{1}{1 + \left(\frac{1}{K} - \frac{1}{2K^2}\right)}, \frac{1}{1 + \frac{\lambda\gamma}{4R^2}} \right\}$$

assume that  $K \geq 2$ , and  $\left(\frac{1}{K} - \frac{1}{2K^2}\right) \geq \frac{\lambda\gamma}{4R^2}$ , which means we focus on the situation that  $\kappa = \frac{R^2}{\lambda\gamma} \gg 1$ . Then we could get that

$$\Theta = \frac{1}{1 + \frac{\lambda\gamma}{4R^2}} = 1 - \frac{1}{1 + \frac{4R^2}{\lambda\gamma}}$$

Based on Corollary 2, we could get that

$$\mathcal{P}(\mathbf{x}^{(t)}) - \mathcal{P}(\mathbf{x}^*) \leq \mathcal{P}(\mathbf{x}^{(t)}) - \mathcal{D}(\mathbf{y}^{(t)}) \leq \left(1 + \frac{R^2}{\lambda\gamma}\right) \Delta^{(t)} \leq C' \Theta^t \left(\mathcal{P}(\mathbf{x}^{(0)}) - \mathcal{P}(\mathbf{x}^*)\right)$$

where  $C'$  is a constant and  $C' > 0$  and  $\mathcal{P}(\cdot)$  is defined in (1). If we apply DiSPA on  $f^{(t)}$  defined in Section 5.1, then based on Proposition 3.2 of Lin et al. (2015), the  $\tau_{\mathcal{M}}$  defined in Proposition 3.2 equals

$$\tau_{\mathcal{M}} = \frac{1}{1 + \frac{4R^2}{(\lambda + \vartheta)\gamma}}$$

where  $\vartheta$  is the parameter defined in Section 5.1. We could get that  $\tau_{\mathcal{M}} \approx \frac{(\lambda + \vartheta)\gamma}{4R^2}$  since  $\kappa \gg 1$ . Now apply Proposition 3.2 in Lin et al. (2015), we could get the global linear rate of convergence with the parameter  $\tau_{\mathcal{A},F}$ ,

$$\tau_{\mathcal{A},F} = \tilde{\mathcal{O}} \left( \tau_{\mathcal{M}} \sqrt{\lambda} / \sqrt{\lambda + \vartheta} \right) = \tilde{\mathcal{O}} \left( \frac{(\lambda + \vartheta)\gamma\sqrt{\lambda}}{4R^2\sqrt{\lambda + \vartheta}} \right)$$

if we take  $\vartheta$  as  $\vartheta = \frac{4R^2}{\gamma} - \lambda$ , then we can get that the communication complexity of the A-DiSPA for achieving  $\mathcal{P}(\mathbf{x}^{(T)}) - \mathcal{P}(\mathbf{x}^*) < \epsilon$  is  $\tilde{\mathcal{O}} \left( \sqrt{\kappa} \log \left( \frac{1}{\epsilon} \right) \right)$ .  $\square$