

# Distributed Momentum-based Frank-Wolfe Algorithm for Stochastic Optimization

Jie Hou, Xianlin Zeng, *Member, IEEE*, Gang Wang, *Member, IEEE*,  
Jian Sun, *Senior Member, IEEE*, and Jie Chen, *Fellow, IEEE*

**Abstract**—This paper considers distributed stochastic optimization, in which a number of agents cooperate to optimize a global objective function through local computations and information exchanges with neighbors over a network. Stochastic optimization problems are usually tackled by variants of projected stochastic gradient descent. However, projecting a point onto a feasible set is often expensive. The Frank-Wolfe (FW) method has well-documented merits in handling convex constraints, but existing stochastic FW algorithms are basically developed for centralized settings. In this context, the present work puts forth a distributed stochastic Frank-Wolfe solver, by judiciously combining Nesterov’s momentum and gradient tracking techniques for stochastic convex and nonconvex optimization over networks. It is shown that the convergence rate of the proposed algorithm is  $\mathcal{O}(k^{-\frac{1}{2}})$  for convex optimization, and  $\mathcal{O}(1/\log_2(k))$  for nonconvex optimization. The efficacy of the algorithm is demonstrated by numerical simulations against a number of competing alternatives.

**Index Terms**—Distributed Optimization, Frank-Wolfe Algorithms, Stochastic Optimization, Momentum-based Method.

## I. INTRODUCTION

**D**ISTRIBUTED stochastic optimization is a basic problem that arises widely in diverse engineering applications, including unmanned systems [1]–[3], distributed machine learning [4], and multi-agent reinforcement learning [5]–[7], to name a few. The goal is to minimize a shared objective function, which is defined as the expectation of a set of stochastic functions subject to general convex constraints, by means of local computations and information exchanges between working agents.

This paper considers a set  $\mathcal{N} = \{1, 2, \dots, n\}$  of working agents connected through a communication network  $\mathcal{G} =$

This work was supported in part by the National Key R&D Program of China under Grant 2021YFB1714800, the National Natural Science Foundation of China under Grants 62073035, 62173034, 61925303, 62088101, 61873033, the CAAI-Huawei MindSpore Open Fund, and the Chongqing Natural Science Foundation under Grant 2021ZX4100027. (*Corresponding author: Xianlin Zeng.*)

J. Hou and X. Zeng are with the Key Laboratory of Intelligent Control and Decision of Complex Systems, School of Automation, Beijing Institute of Technology, Beijing, 100081, China (E-mail: houjie@bit.edu.cn; xianlin.zeng@bit.edu.cn).

G. Wang and J. Sun are with the Key Laboratory of Intelligent Control and Decision of Complex Systems, School of Automation, Beijing Institute of Technology, Beijing, 100081, China and Beijing Institute of Technology Chongqing Innovation Center, Chongqing, 401120, China (E-mail: gangwang@bit.edu.cn; sunjian@bit.edu.cn).

J. Chen is with the School of Electronic and Information Engineering, Tongji University, Shanghai, 200082, China and also with the Key Laboratory of Intelligent Control and Decision of Complex Systems, School of Automation, Beijing Institute of Technology, Beijing, 100081, China (E-mail: chenjie@bit.edu.cn).

$(\mathcal{N}, \mathcal{E})$ , where  $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N}$  denotes the set of edges. Each agent  $i \in \mathcal{N}$  has a local objective function  $F_i(x) = \mathbb{E}[f_i(x, \xi^i)]$ , where  $f_i$  is a stochastic function involving strategy variable  $x \in \mathbb{R}^p$  and random variable  $\xi^i$  that follows an unknown distribution. The collective goal of all the agents is to find  $x^*$  that minimizes the average of all objective functions, i.e.,

$$\min_{x \in \mathcal{X}} F(x) := \frac{1}{n} \sum_{i=1}^n F_i(x) \quad (1)$$

where  $\mathcal{X} \subset \mathbb{R}^p$  is a convex and compact feasible set. Problems of the form (1) lie at the heart of machine learning and adaptive filtering, emerging in e.g., clustering, classification, energy management, and resource allocation [8]–[11].

A popular approach to solving problem (1) is the projected stochastic gradient descent (pSGD) [12]–[14]. In pSGD and its variants, the iteration variable is projected back onto  $\mathcal{X}$  after taking a step in the direction of negative stochastic gradient [15]–[20]. Such algorithms are efficient when the computational cost of performing the projection is low, e.g., projecting onto a hypercube or simplex. In many practical situations of interest, however, the cost of projecting onto  $\mathcal{X}$  can be high, e.g., dealing with a trace norm ball or a base polytope  $\mathcal{X}$  in submodular minimization [21].

An alternative for tackling problem (1) is the projection-free methods, including the Frank-Wolfe (FW) [22] and conditional gradient sliding [23]. In this paper, we focus on the FW algorithm, which is also known as conditional gradient method [22]. Classical FW methods circumvent the projection step by first solving a linear minimization subproblem over the constraint set  $\mathcal{X}$  to obtain a sort of conditional gradient  $\theta_k$ , which is followed by updating  $x_{k+1}$  through a convex combination of the current iteration variables  $x_k$  and  $\theta_k$ . On top of this idea, a number of modifications have been proposed to improve or accelerate the FW method in algorithm design or convergence analysis, see e.g., [8]–[10], [24]–[31].

Nonetheless, most existing projected stochastic gradient descent (pSGD) and Frank-Wolfe variants are designed for constrained centralized problems, and they cannot directly handle distributed problems. Therefore, it is necessary to develop distributed stochastic projection-free methods for problem (1). In addition, stochastic FW methods may not converge even in the centralized convex case, without increasing the batch size [8]. In this context, a natural question arises: *is it possible to develop a distributed FW method by using any fixed batch size for problem (1), while enjoying a convergence rate comparable to that of centralized stochastic FW methods?* In this paper, we answer this question affirmatively, by carefully designing a

distributed stochastic FW algorithm, which converges for any fixed batch size (can be as small as 1) and enjoys a comparable convergence rate as in the centralized stochastic setting.

### A. Related Works

**Projection-free stochastic** algorithms for addressing stochastic optimization problems were widely studied in recent years. When both the global function  $F$  and the feasibility set  $\mathcal{X}$  in (1) are convex, an online stochastic FW method using minibatches was proposed and shown to converge at a rate of  $\mathcal{O}(k^{-\frac{1}{4}})$  [28]. By progressively increasing the batch size per iteration, convergence rate of stochastic FW algorithms was improved to  $\mathcal{O}(k^{-\frac{1}{3}})$  in [29]. The work [8] further relaxed the requirement of increasing the batchsize by using a fixed small batch size along with some heuristics, while maintaining the convergence rate of  $\mathcal{O}(k^{-\frac{1}{3}})$ . Faster rate  $\mathcal{O}(k^{-\frac{1}{2}})$  was obtained by merging the Nesterov's momentum and classical FW method in [9]. When  $F$  becomes nonconvex, [10] proposed a stochastic FW variant, and established a convergence rate of  $\mathcal{O}(k^{-\frac{1}{4}})$  for handling nonconvex stochastic optimization problems. Lately, [31] studied nonconvex stochastic optimization on Riemannian manifolds and presented a projection-free stochastic algorithm which achieves the same convergence rate as in [10]. It is worth noting that the aforementioned FW methods are all centralized. Thus far, distributed projection-free stochastic algorithms have rarely been studied.

**Distributed FW** methods play an important role in distributed convex and nonconvex optimization, a sample of which can be found in [32]–[38]. In the deterministic case, a distributed FW algorithm was developed in [33] for a class of nonconvex optimization problems. The work [34] further devised distributed FW algorithms, and showed convergence rate of  $\mathcal{O}(k^{-1})$  for convex optimization and  $\mathcal{O}(k^{-\frac{1}{2}})$  for nonconvex optimization. For submodular maximization, [37] proposed two distributed algorithms for deterministic and stochastic optimization, and obtained the convergence rate of  $\mathcal{O}(k^{-\frac{1}{3}})$  for stochastic optimization. The convergence rate in [37] was improved to  $\mathcal{O}(k^{-\frac{1}{2}})$  [38] by using variance reduction techniques and gradient tracking strategies.

Although considerable results have been reported for distributed FW in deterministic settings, they cannot be directly applied in and/or generalized to stochastic settings. The reason is twofold: i) FW may diverge due to the non-vanishing variance in gradient estimates; and, ii) the desired convergence rate of FW for stochastic optimization is not guaranteed to be comparable to pSGD, even for the centralized setting.

To address these challenges, the present paper puts forth a distributed stochastic version of the celebrated FW algorithm for stochastic optimization over networks. The main idea behind our proposal is a judicious combination of the recursive momentum [39] and the Nesterov's momentum [40]. On the theory side, it is shown that the proposed algorithm can not only attenuate the noise in gradient approximation, but also achieve a convergence guarantee comparable to pSGD in convex case. Comparison of the proposed algorithm in context is provided in Table I.

TABLE I  
CONVERGENCE RATE FOR STOCHASTIC OPTIMIZATION

Reference	Setting	Projection-free	Function	Rate
RSA [16]	centralized	unconstrained	smooth convex	$\mathcal{O}(k^{-\frac{1}{2}})$
RSG [19]	centralized	no	smooth nonconvex	$\mathcal{O}(k^{-\frac{1}{2}})$
SPPDM [13]	distributed	unconstrained	nonsmooth nonconvex	$\mathcal{O}(k^{-\frac{1}{2}})$
OFW [28]	centralized	yes	smooth $Q$ -Lipschitz convex	$\mathcal{O}(k^{-\frac{1}{4}})$
SFW [8]	centralized	yes	smooth convex	$\mathcal{O}(k^{-\frac{1}{3}})$
MSHFW [9]	centralized	yes	smooth convex	$\mathcal{O}(k^{-\frac{1}{2}})$
NSFW [10]	centralized	yes	smooth $Q$ -Lipschitz nonconvex	$\mathcal{O}(k^{-\frac{1}{4}})$
SRFW [31]	centralized	yes	smooth $Q$ -Lipschitz nonconvex	$\mathcal{O}(k^{-\frac{1}{4}})$
<b>This Work</b>	distributed	yes	smooth convex	$\mathcal{O}(k^{-\frac{1}{2}})$
			smooth nonconvex	$\mathcal{O}(\frac{1}{\log_2(k)})$

\* The function  $f_i(x, \xi^i)$  in (1) is  $Q$ -Lipschitz if  $\|\nabla f_i(x, \xi^i)\| \leq Q$  for all  $\xi^i$ , where  $Q$  is a positive constant [10];  $k$  denotes the number of iterations.

### B. Our contributions

In succinct form, the contributions of this work are summarized as follows.

- 1) We propose a projection-free algorithm, referred to as the distributed momentum-based Frank-Wolfe (DMFW), for convex and nonconvex stochastic optimization over networks. Compared with the centralized FW methods [8]–[10], [28], [31], DMFW is considerably different in algorithm design and convergence analysis.
- 2) For convex objective functions, we establish a convergence rate of  $\mathcal{O}(k^{-\frac{1}{2}})$  for DMFW, which matches that of distributed pSGD [16] and is even faster than those of centralized FW algorithms in [8], [28].
- 3) For nonconvex objective functions, we establish a convergence rate of  $\mathcal{O}(1/\log_2(k))$  for DMFW, which, to the authors' best knowledge, marks the first FW's convergence rate result for distributed nonconvex stochastic optimization.

## II. PRELIMINARIES AND ALGORITHM DESIGN

### A. Notation and preliminaries

Let  $\mathbb{R}$  denote the set of real numbers, and  $\mathbb{R}^p$  the set of  $p$ -dimensional real vectors;  $\langle \cdot \rangle$  denotes the inner product;  $(\cdot)^T$  represents the transpose;  $\|x\|$  denotes the  $l_2$  norm (Euclidean norm) of vector  $x$ , and  $\|x\|_q$  ( $q \in [1, +\infty)$ ) symbols the  $l_q$  norm of vector  $x$ ;  $\max\{\cdot\}$  denotes the maximum element in set  $\{\cdot\}$ ;  $\lceil \cdot \rceil$  is the ceiling operation;  $\mathbb{E}[\cdot]$  is the expectation operator;  $\mathbb{E}_k[\cdot]$  is the conditional expectation  $\mathbb{E}[\cdot | \mathcal{F}_k]$  on the sigma field  $\mathcal{F}_k$  which contains all types of randomness up to iteration  $k$ ;  $C = [c_{ij}]_{n \times n}$  is the weighted adjacency matrix of

graph  $\mathcal{G}(\mathcal{N}, \mathcal{E})$ . For  $\forall i, j \in \mathcal{N}$ , if  $(i, j) \in \mathcal{E}$ , then  $c_{ij} > 0$ , and  $c_{ij} = 0$  otherwise.

Consider a differentiable function  $F : \mathbb{R}^p \rightarrow \mathbb{R}$ , whose gradient is  $\nabla F(\cdot)$ . The function is  $L$ -smooth over a convex set  $\mathcal{X}$  if

$$F(x) - F(y) \leq \langle \nabla F(y), x - y \rangle + \frac{L}{2} \|x - y\|^2, \quad \forall x, y \in \mathcal{X}$$

where  $L > 0$  is a constant. The function  $F$  is said to be convex over a convex set  $\mathcal{X}$  if  $F(x) - F(y) \geq \langle \nabla F(y), x - y \rangle$  for all  $x, y \in \mathcal{X}$ .

## B. Algorithm design

To solve problem (1), we propose a distributed momentum-based Frank-Wolfe algorithm, which is summarized in Algorithm 1.

---

### Algorithm 1 Distributed Momentum-based Frank-Wolfe

---

**Input:** number of iterations  $K$ , initial condition  $x_1^i \in \mathcal{X}$ , and  $y_1^i = \nabla f_i(\hat{x}_1^i, \xi_1^i) = s_1^i$  for  $\forall i \in \mathcal{N}$ .

1: **for all**  $k = 1, 2, \dots, K$  **do**

2: *Average consensus:*

$$\hat{x}_k^i = \sum_{j \in \mathcal{N}_i} c_{ij} x_k^j \quad (2)$$

where  $\mathcal{N}_i$  is the set of neighbors of node  $i$ .

3: *Momentum update:*

$$y_k^i = (1 - \gamma_k) y_{k-1}^i + \nabla f_i(\hat{x}_k^i, \xi_k^i) - (1 - \gamma_k) \nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i) \quad (3)$$

where  $\gamma_k \in (0, 1]$  is a step size.

4: *Gradient tracking:*

$$s_k^i = \sum_{j \in \mathcal{N}_i} c_{ij} s_{k-1}^j + y_k^i - y_{k-1}^i \quad (4)$$

$$p_k^i = \sum_{j \in \mathcal{N}_i} c_{ij} s_k^j \quad (5)$$

5: *Frank-Wolfe step:*

$$\theta_k^i \in \operatorname{argmin}_{\phi \in \mathcal{X}} \langle p_k^i, \phi \rangle \quad (6)$$

$$x_{k+1}^i = \hat{x}_k^i + \eta_k (\theta_k^i - \hat{x}_k^i) \quad (7)$$

where  $\eta_k \in (0, 1]$  is a step size.

6: **end for**

7: **return**  $x_{k+1}^i$  for all  $i \in \mathcal{N}$ .

---

*Average consensus:* We employ the average consensus (AC) protocol [41]–[44], in which an agent takes a weighted average of the values from its neighbors according to  $C$ .

*Momentum update:* Because the distribution of  $\xi^i$  in (1) is unknown, we can only have access to stochastic gradients of  $F_i(x)$ , that is, for a given  $x \in \mathcal{X}$  and randomly sampled  $\xi^i$ , the oracle returns  $\nabla f_i(x, \xi^i)$ , which is assumed to be an unbiased estimate of  $\nabla F_i(x)$ . It is well known that the naive stochastic implementation of Frank-Wolfe by replacing  $\nabla F_i(x)$  with  $\nabla f_i(x, \xi^i)$ , may diverge due to the non-vanishing variance of

$\nabla f_i(x, \xi^i)$ . To address this issue, we generalize the recursive momentum in [39] to distributed stochastic optimization.

*Gradient tracking:* Inspired by the gradient tracking method in [45], [46], which reuses the global gradient  $p_{k-1}^i$  from the last iteration, agent  $i$  at iteration  $k$  approximates the global gradient via (4) and (5). The initialization  $s_1^i = \nabla f_i(\hat{x}_1^i, \xi_1^i)$  is set for  $\forall i \in \mathcal{N}$ .

*Frank-Wolfe step:* A feasible direction  $\theta_k^i$  is obtained by minimizing its correlation with  $p_i(k)$  over  $\mathcal{X}$  in (6). Subsequently, the variable  $x_{k+1}^i$  is generated as a convex combination of  $\hat{x}_k^i$  and  $\theta_k^i$ .

**Remark 1.** *There are two mechanisms for information exchanging with neighbors in DMFW: (i) **average consensus**; and (ii) **gradient tracking**.*

*In the **average consensus** step, agent  $i$  approximates the average iteration by exchanging the latest iteration information with its neighbors. In the **gradient tracking** step, agent  $i$  approximates the global gradient by weighted averaging  $s_{k-1}^j$  and  $y_k^i - y_{k-1}^i$ , which is an estimate of local gradient difference.*

**Remark 2.** *Compared with the existing distributed solutions [34], [45]–[47], Algorithm 1 shares very similar structures: **global consensus** steps plus **local adaptation** steps.*

***global consensus:** Algorithm 1 realizes the distributed update by exploiting a twofold consensus-based mechanism to: (i) enforce an agreement among the agents' estimates  $\hat{x}_k^i$ ; and (ii) dynamically track the gradient of the whole cost function through an auxiliary variable  $s_k^i$ .*

***local adaptation:** Agent  $i$  approximates the local gradient  $y_k^i$  and updates its variable  $x_{k+1}^i$  independently via local learning process by using  $\hat{x}_k^i$  and  $s_k^i$  obtained from the **global consensus** steps.*

**Remark 3.** *It is worth mentioning that Algorithm 1 works with a single stochastic gradient (i.e., with batch size as small as 1), unlike the methods in [37], which requires increasing the batch sizes as the number of iterations  $k$  grows.*

## III. MAIN RESULTS

In this section, we establish the convergence results of the proposed algorithm for convex and nonconvex problems, respectively. Before providing the results, we outline some standing assumptions and facts.

### A. Assumptions and facts

**Assumption 1** (Weight rule). *The weighted adjacency matrix  $C$  is a doubly stochastic matrix, i.e., the row sum and the column sum of  $C$  are all 1.*

Assumption 1 indicates that for each round of the *Average Consensus* step of Algorithm 1, the agent takes a weighted average of the values from its neighbors according to  $C$ .

**Assumption 2** (Connectivity). *The network  $\mathcal{G}$  is connected.*

If Assumptions 1 and 2 hold, the magnitude of the second largest eigenvalue of the weighted adjacency matrix  $C$ , denoted by  $\lambda$ , is strictly less than one, i.e.,  $|\lambda| < 1$  [34]. The following fact holds for any doubly stochastic matrix  $C$ .

**Fact 1** Let  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x^i$  and  $\hat{x}^i = \sum_{j=1}^n c_{ij} x^j$ . Then, the following inequality holds

$$\left( \sum_{i=1}^n \|\hat{x}^i - \bar{x}\|^2 \right)^{\frac{1}{2}} \leq |\lambda| \left( \sum_{i=1}^n \|x^i - \bar{x}\|^2 \right)^{\frac{1}{2}}. \quad (8)$$

If Assumptions 1 and 2 hold, Fact 1 implies that each *Average Consensus* update brings the iteration variables closer to their average  $\bar{x}$ . For convenience, we define  $k_0$  to be the smallest positive integer such that  $|\lambda| \leq [k_0/(k_0 + 1)]^2$ . Clearly,  $k_0 = \lceil (|\lambda|^{-\frac{1}{2}} - 1)^{-1} \rceil$ . The next assumption is on the constraint set of  $\mathcal{X}$ , which is a standard requirement for analyzing FW methods.

**Assumption 3** (Constraint set).  $\mathcal{X}$  is a convex and compact set with diameter  $D$ , i.e., there is some constant  $D > 0$  such that  $\|x - x'\| \leq D$  for all  $x, x' \in \mathcal{X}$ .

**Assumption 4** ( $L$ -smoothness). Functions  $F_i(x)$  and  $f_i(x, \xi^i)$  are  $L$ -smooth with respect to  $x$  for all  $\xi^i$  and  $i \in \mathcal{N}$ .

**Assumption 5** (Bounded stochastic gradients). The variance of the stochastic gradient  $\nabla f_i(x, \xi^i)$  ( $x \in \mathcal{X}$ ,  $i \in \mathcal{N}$ ) is bounded, that is,  $\mathbb{E}[\|\nabla F_i(x) - \nabla f_i(x, \xi^i)\|^2] \leq \delta^2$ .

Assumption 5 is standard for stochastic FW algorithms. The bound will be frequently used for convergence analysis.

**Fact 2** Suppose Assumptions 3 and 5 hold. There exists a constant  $G > 0$  such that  $\mathbb{E}[\|\nabla f_i(x, \xi^i)\|^2] \leq G^2$  and  $\mathbb{E}[\|\nabla F_i(x, \xi^i)\|] \leq G$ .

*Proof.* It follows from the Jensen's inequality that

$$\begin{aligned} \mathbb{E}[\|\nabla F_i(x) - \nabla f_i(x, \xi^i)\|^2] &\geq (\mathbb{E}[\|\nabla f_i(x, \xi^i) - \nabla F_i(x)\|])^2 \\ &\geq \|\mathbb{E}[\nabla f_i(x, \xi^i)] - \nabla F_i(x)\|^2. \end{aligned} \quad (9)$$

Thus, it follows from (9) and Assumption 5 that  $\|\mathbb{E}[\nabla f_i(x, \xi^i)] - \nabla F_i(x)\| \leq \delta^2$ , i.e.,  $\|\mathbb{E}[\nabla f_i(x, \xi^i)] - \nabla F_i(x)\| \leq \delta$ . In addition, we also obtain that  $\mathbb{E}[\nabla f_i(x, \xi^i)]$  is bounded because  $\nabla F_i(x)$  is bounded followed by Assumption 3. In the meanwhile, we have  $\|\nabla F_i(x)\|^2 - 2\mathbb{E}[\nabla f_i(x, \xi^i)]^\top \nabla F_i(x) + \mathbb{E}[\|\nabla f_i(x, \xi^i)\|^2] \leq \delta^2$  from (9). Hence,  $\mathbb{E}[\|\nabla f_i(x, \xi^i)\|^2]$  has an upper bound, which implies that there is a scalar  $G_i$  that  $\mathbb{E}[\|\nabla f_i(x, \xi^i)\|^2] \leq G_i^2$ . Because  $G_i^2 \geq \mathbb{E}[\|\nabla f_i(x, \xi^i)\|^2] \geq (\mathbb{E}[\|\nabla f_i(x, \xi^i)\|])^2$ , it has  $(\mathbb{E}[\|\nabla f_i(x, \xi^i)\|])^2 \leq G_i^2$ , that is  $\mathbb{E}[\|\nabla f_i(x, \xi^i)\|] \leq G_i$ . Let  $G := \max_{i \in \mathcal{N}} \{G_i\}$ . We have that  $(\mathbb{E}[\|\nabla f_i(x, \xi^i)\|])^2 \leq G^2$  and  $\mathbb{E}[\|\nabla f_i(x, \xi^i)\|] \leq G$ .  $\square$

**Remark 4.** This paper makes weaker assumptions on objective functions. Specifically, compared with [28], [10] and [31], we do not require the  $Q$ -Lipschitz continuity of  $f_i(x, \xi^i)$  in (1), i.e.,  $\|\nabla f_i(x, \xi^i)\| \leq Q$  for all  $\xi^i$ , where  $Q$  is a positive constant.

### B. Convergence rate for convex stochastic optimization

This subsection is dedicated to the performance analysis of Algorithm 1. Let us start by defining the following auxiliary vectors

$$\bar{x}_k := \frac{1}{n} \sum_{i=1}^n x_k^i, \quad \bar{y}_k := \frac{1}{n} \sum_{i=1}^n y_k^i, \quad \bar{P}_k := \frac{1}{n} \sum_{i=1}^n \nabla F_i(\hat{x}_k^i).$$

We begin our analysis by characterizing the behavior of  $\{\hat{x}_k^i\}$  for all  $i \in \mathcal{N}$  in the next lemma.

**Lemma 1.** Suppose Assumptions 1-3 hold. Let  $\eta_k = \frac{2}{k+2}$ . Then, for any  $i \in \mathcal{N}$  and  $k \geq 1$ , we have

$$\|\hat{x}_k^i - \bar{x}_k\| \leq \frac{2C_1}{k+2}, \quad (10)$$

where  $C_1 = k_0 \sqrt{n}D$ .

The proof is presented in Appendix VI-A. Lemma 1 shows that  $\|\hat{x}_k^i - \bar{x}_k\| = \mathcal{O}(1/k)$ , which implies that  $\|\hat{x}_k^i - \bar{x}_k\|$  converges to zero as  $k \rightarrow \infty$ . By selecting appropriate step sizes  $\gamma_k$  and  $\eta_k$ , we establish the boundedness of  $\|p_k^i - \bar{y}_k\|^2$  for all  $i \in \mathcal{N}$  in the next lemma.

**Lemma 2.** Suppose Assumptions 1-5 hold. Choose the step sizes  $\gamma_k = \frac{2}{k+1}$  and  $\eta_k = \frac{2}{k+2}$ . Then, for any  $i \in \mathcal{N}$  and  $k \geq 1$ , it holds

$$\mathbb{E}[\|p_k^i - \bar{y}_k\|^2] \leq \frac{4C_2}{(k+2)^2}, \quad (11)$$

where  $C_2 = k_0^3(4n)^{k_0}n(12L^2(D+2C_1)^2 + 12(G^2 + \hat{\psi}))$ ,  $\hat{\psi} = \max_{i \in \mathcal{N}} \{\|y_1^i\|^2, 4L(D+2C_1)\psi + 4G\psi + 8G^2 + 8L^2(D+2C_1)^2\}$ ,  $\psi = \max_{i \in \mathcal{N}} \{\|y_1^i\|, 2G + 2L(D+2C_1)\}$ .

The proof is presented in Appendix VI-C. In order to prove the convergence of Algorithm 1, we provide the following lemma.

**Lemma 3.** Suppose Assumptions 1-5 hold. Then,

(a) the conditional expectation of  $\|\bar{P}_k - \bar{y}_k\|^2$  satisfies

$$\begin{aligned} \mathbb{E}_k[\|\bar{P}_k - \bar{y}_k | \mathcal{F}_k\|^2] &\leq (1 - \gamma_k) \|\bar{P}_{k-1} - \bar{y}_{k-1}\|^2 \\ &\quad + 6L^2(D+2C_1)^2 \eta_{k-1}^2 + 3\gamma_k^2 \delta^2 \end{aligned}$$

for any  $k \geq 2$ ,

(b) taking the step sizes  $\gamma_k = \frac{2}{k+1}$  and  $\eta_k = \frac{2}{k+2}$ , the expectation of  $\|\bar{P}_k - \bar{y}_k\|^2$  satisfies

$$\mathbb{E}[\|\bar{P}_k - \bar{y}_k\|^2] \leq \frac{C_3}{k+2} \quad (12)$$

for any  $k \geq 1$ , where  $C_3 = 24L^2(D+2C_1)^2 + 12\delta^2$ .

The proof is presented in Appendix VI-D. Lemma 3 asserts that the expectation of  $\|\bar{P}_k - \bar{y}_k\|^2$  converges to zero as  $k \rightarrow \infty$ . Leveraging Lemma 2 and Lemma 3, the boundedness of  $\mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|^2]$  can be derived in the next lemma.

**Remark 5.** Consider the error  $\epsilon_k = \bar{P}_k - \bar{y}_k$ , which measures the difference incurred by using  $y_k^i$  as the update direction instead of the correct yet unknown direction  $\nabla F_i(\hat{x}_k^i)$  for each agent  $i$ . Lemma 3 suggests that  $\mathbb{E}[\|\epsilon_k\|^2]$  decreases over iterations, that is, the noise of the stochastic gradient approximation diminishes as the number of iterations increases.

**Remark 6.** The conditions in Lemma 3 can be guaranteed by choosing  $\gamma_k = A/(k+t_0)$  and  $\eta_k = B/(k+t_0+1)$  according to Lemma 5, where  $A > 1$ ,  $B \geq 0$  and  $t_0$  is a constant. Typical choices are  $\gamma_k = 2/(k+1)$  and  $\eta_k = 2/(k+2)$ .

**Lemma 4.** *Suppose Assumptions 1-5 hold. Take the step sizes  $\gamma_k = \frac{2}{k+1}$  and  $\eta_k = \frac{2}{k+2}$ . Then, for any  $i \in \mathcal{N}$  and  $k \geq 1$ , we have*

$$\mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|^2] \leq \frac{12L^2C_1^2 + 3C_3 + 12C_2}{k+2}. \quad (13)$$

The proof is presented in Appendix VI-E. Making use of Lemma 4, the convergence rate of Algorithm 1 is established.

**Theorem 1.** *Suppose Assumptions 1-5 hold. The function  $F$  is convex. Choose the step sizes  $\gamma_k = \frac{2}{k+1}$  and  $\eta_k = \frac{2}{k+2}$ . Then, for any  $k \geq 1$ , it holds*

$$\mathbb{E}[F(\bar{x}_k)] - F(x^*) \leq \frac{C_4}{(k+3)^{\frac{1}{2}}}, \quad (14)$$

where  $C_4 = \frac{\max\{\sqrt{3}(F(\bar{x}_1) - F(x^*)), 2LD^2 + 2D\sqrt{12L^2C_1^2 + 3C_3 + 12C_2}\}}{C_4k^{-\alpha}}$ .

The proof is presented in Appendix VI-F.

**Remark 7.** *Theorem 1 implies that the expected suboptimality  $\mathbb{E}[F(\bar{x}_k)] - F(x^*)$  of the iteration variables generated by DMFW converges to zero at least at a sublinear rate of  $\mathcal{O}(k^{-\frac{1}{2}})$ , coinciding with that of MSHFW [9]. By using the Markov inequality, we can also obtain that  $\mathbb{P}\{F(\bar{x}_k) - F(x^*) \geq C_4k^{-\alpha}\} \leq \frac{\mathbb{E}[F(\bar{x}_k)] - F(x^*)}{C_4k^{-\alpha}} \leq k^\alpha/(k+3)^{\frac{1}{2}} = \mathcal{O}(k^{\alpha-\frac{1}{2}})$ , where  $0 < \alpha < \frac{1}{2}$ . That is,  $\mathbb{P}\{F(\bar{x}_k) - F(x^*) < C_4k^{-\alpha}\} \geq 1 - k^\alpha/(k+3)^{\frac{1}{2}}$ . Hence,  $\mathbb{P}\{F(\bar{x}_k) - F(x^*) < C_4k^{-\alpha}\} = 1$  when  $k \rightarrow \infty$ , which indicates that  $F(\bar{x}_k)$  converges to  $F(x^*)$  with probability 1.*

### C. Convergence rate for nonconvex optimization

This subsection provides the convergence rate of the proposed DMFW for (1) with nonconvex objective functions. To show the convergence performance of DMFW for nonconvex case, we introduce the FW-gap, which is defined as

$$g_k = \max_{x \in \mathcal{X}} \langle \nabla F(\bar{x}_k), \bar{x}_k - x \rangle. \quad (15)$$

According to (15), the variable  $\bar{x}_k$  is a stationary point to (1) if  $g_k = 0$  [34]. Hence,  $g_k$  can be regarded as a measure of the stationarity of variable  $\bar{x}_k$ . Since the set  $\mathcal{X}$  is compact, we assume that the set of stationary points  $\mathcal{X}^*$  of (1) is nonempty and the function  $F(x)$  is finite on  $\mathcal{X}^*$ . The convergence result is presented in the following theorem.

**Theorem 2.** *Consider the DMFW algorithm. Suppose Assumptions 1-5 hold.  $F$  is possibly nonconvex. If  $\gamma_k = \frac{2}{k+1}$  and  $\eta_k = \frac{2}{k+2}$ , the FW-gap satisfies*

$$\mathbb{E} \left[ \min_{k \in [1, K]} g_k \right] \leq \frac{1}{\log_2(K) - 1} \left( \mathbb{E}[F(\bar{x}_1) - F(x^*)] + 4LD^2 + 2D\beta\sqrt{12L^2C_1^2 + 3C_3 + 12C_2} \right),$$

where  $K = 2^m$  ( $m \in \mathbb{Z}_+$ ) and  $\beta$  is a constant such that  $\sum_{k=1}^{2^m} \frac{2}{(k+2)^{1.5}} \leq \beta$ .

The proof is presented in Appendix VI-G.

**Remark 8.** *Theorem 2 indicates that the convergence rate of the proposed DMFW is  $\mathcal{O}(1/\log_2(k))$  for nonconvex objective*

*functions. It is worth mentioning that the obtained result is novel compared with previous nonconvex stochastic FW studies, even in centralized setting. For instance, [10] and [31] proposed centralized stochastic FW algorithms for stochastic nonconvex problems by using minibatch method. The proposed DMFW only requires selecting a small fixed batch data (as small as 1) to compute the stochastic gradient to achieve convergence. In addition, the proposed DMFW relaxes the assumption of  $Q$ -Lipschitz continuity of  $f_i(x, \xi^i)$  with respects to  $x$ , which is required in e.g. [10] and [31].*

**Remark 9.** *There are two challenges in the convergence analysis of the proposed algorithm. On one hand, our method needs to deal with the consensus error among different agents compared with the solutions in [39] and [9]. On the other hand, the introduced local gradient by using recursive momentum makes it more challenging to handle the consensus error compared with the method in [37].*

## IV. NUMERICAL TESTS

### A. Binary classification

In this subsection, several numerical experiments are provided on the binary classification with an  $l_2$ -norm ball constraint (i.e.,  $\mathcal{X} = \{x \mid \|x\| \leq \frac{D}{2}\}$ ). We solve the problem with Algorithm 1 (DMFW) and compare it against SFW [8], MSHFW [9] and DeFW [34] as baselines. We consider two cases where objective functions are convex and nonconvex, respectively. We use three different public datasets, which are summarized in Table II. In the experiment for SFW, MSHFW and DMFW, we compute a stochastic gradient by using 1% of data per iteration; in the experiment for deterministic algorithm DeFW, we use full data to compute the gradient. Distributed algorithms DeFW and DMFW are applied over a connected network  $\mathcal{G}$  of 5 agents with a doubly stochastic adjacency matrix  $C$  to solve the problem. The communication topology is demonstrated in Fig. 1.

#### 1) Binary classification with convex objective functions:

We consider a popular logistic regression for binary classification with convex objective functions as follows:

$$\min_{x \in \mathcal{X}} F(x) = \frac{1}{n} \sum_{i=1}^n F_i(x),$$

$$F_i(x) = \frac{1}{m_i} \sum_{j=1}^{m_i} \ln[1 + \exp(-b_j \langle a_j, x \rangle)],$$

where  $(a_i, b_i)$  represents the (feature, label) pair of datum  $i$ ,  $m_i$  is the total number of training samples of agent  $i$ , and  $n$  is the total number of agents over network  $\mathcal{G}$ . Here, we set  $n = 5$  for DeFW and DMFW. Note that SFW and MSHFW are centralized algorithms, so  $n = 1$  for these two algorithms. As benchmarks, we choose SFW with step sizes  $\gamma_k = 2/(k+8)$  and  $\rho_k = 4/(k+8)^{\frac{2}{3}}$  in [8], MSHFW with step sizes  $\gamma_k = 2/(k+1)$  and  $\eta_k = 2/(k+2)$  in [9], DeFW with step size  $\gamma_k = 2/(k+1)$  in [34]. The step sizes of the algorithm DMFW are  $\gamma_k = 2/(k+1)$  and  $\eta_k = 2/(k+2)$ . The constraint is selected as  $\mathcal{X} = \{x \mid \|x\| \leq 5\}$ .

As described in Lemma 3, the proposed algorithm can attenuate the noise in gradient approximation, i.e.,  $\mathbb{E}[\|\bar{P}_k -$

TABLE II  
REAL DATA FOR BLACK-BOX BINARY CLASSIFICATION

datasets	#samples	#features	#classes
<i>covtype.binary</i>	581012	54	2
<i>a9a</i>	32561	123	2
<i>w8a</i>	64700	300	2

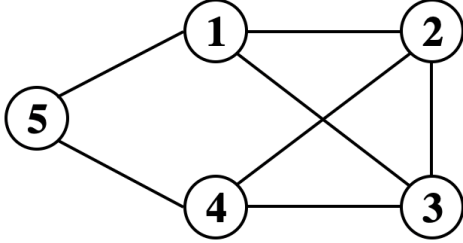


Fig. 1. Multi-agent communication topology.

$\bar{y}_k\|^2] = \mathcal{O}(1/k)$ . This is verified by our simulation result in Fig. 1. It can be observed that  $\|\bar{P}_k - \bar{y}_k\|^2$  generally decreases with the increase of iterations, which is consistent with our theoretical analysis.

We evaluate the algorithm in terms of the FW-gap, which is defined as

$$g_k = \max_{x \in \mathcal{X}} \langle \nabla F(\bar{x}_k), \bar{x}_k - x \rangle.$$

For different datasets, Fig. 2 shows the FW-gap of SFW, MSHFW, DeFW and DMFW. It is shown that DMFW has comparable convergence performance compared with the distributed deterministic algorithm DeFW, although DMFW uses less data. This implies that the local gradient estimate  $y_k^i$  in DMFW may be a better candidate for approximating the gradient  $\nabla F_i(x)$  comparing to the unbiased gradient estimate  $\nabla f_i(x, \xi^i)$ . Comparing stochastic methods (MSHFW, SFW and DMFW), DMFW and MSHFW outperform SFW, especially on dataset *w8a*, which is consistent with the theoretical result.

**Remark 10.** *Stochastic algorithms only require a certain amount of data to be obtained randomly at each time to achieve convergence, unlike deterministic algorithms, which need to obtain the whole data at one time prior to the start of the algorithms. Hence, finite-sum problem can also be solved by stochastic algorithms. In Test A, the random variable  $\xi^i$  denotes training examples obtained randomly by agent  $i$  at each time, we use 1% of data to estimate a stochastic gradient  $\nabla f_i(x, \xi^i)$  for stochastic algorithms at each iteration, not the full gradient in finite-sum setting.*

2) *Binary Classification with Nonconvex Objective Functions:* We consider a binary classification with nonconvex objective function as follows:

$$\begin{aligned} \min_{x \in \mathcal{X}} F(x) &= \frac{1}{n} \sum_{i=1}^n F_i(x), \\ F_i(x) &= \frac{1}{m_i} \sum_{i=1}^{m_i} \frac{1}{1 + \exp(b_i \langle a_i, x \rangle)} + \lambda_1 \|x\|^2 \end{aligned} \quad (16)$$

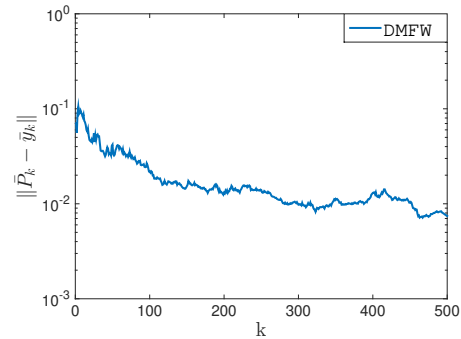


Fig. 2. The error  $\|\bar{P}_k - \bar{y}_k\|$  of DMFW on *a9a* dataset.

where  $(a_i, b_i)$ ,  $m_i$ ,  $n$  and the constraint  $\mathcal{X}$  have the same definitions as those in IV-A1, and  $\lambda_1 = 5 \times 10^{-6}$ . It is obvious that the objective function  $F$  is a nonconvex function. As benchmarks, we choose DeFW with step size  $\gamma_k = \frac{1}{k^{0.5}}$  as mentioned in [34]. Algorithms SFW and MSHFW use the same step sizes in Section IV-A1. The step sizes of the algorithm DMFW are  $\gamma_k = \frac{2}{k+1}$  and  $\eta_k = \frac{2}{k+2}$ . Note that SFW and MSHFW are not proved to converge for the nonconvex problem (16). We implement these two algorithms only for comparison purpose.

Fig. 3 shows the FW-gap of SFW, MSHFW, DeFW and DMFW for solving nonconvex problem (16). From the results, it can be observed that stochastic algorithms (SFW, DMFW and MSHFW) perform better compared to deterministic algorithm DeFW in all tested datasets. This indicates that the stochastic FW algorithms are more efficient than the deterministic FW algorithms in solving nonconvex problems (16). Comparing DMFW with centralized algorithms MSHFW and SFW, DMFW slightly outperforms SFW, especially in datasets *a9a*, but is slower than MSHFW.

### B. Stochastic ridge regression

In this subsection, several numerical experiments are conducted for the stochastic ridge regression. The constraint sets considered include  $l_1$ -,  $l_2$ - and  $l_{\frac{5}{4}}$ - norm balls. We solve the problem with Algorithm 1 (DMFW) and compare it against SFW [8] and MSHFW [9] as baselines. DMFW is applied over a connected network  $\mathcal{G}$  of  $n$  agents with a doubly stochastic adjacency matrix  $C$  to solve the problem. The stochastic ridge regression is as follows:

$$\begin{aligned} \min_{x \in \mathcal{X}} F(x) &= \frac{1}{n} \sum_{i=1}^n F_i(x), \\ F_i(x) &= \mathbb{E}_{a_i, b_i} [(a_i^T x - b_i)^2 + \lambda_1 \|x\|^2], \end{aligned}$$

where  $\lambda_1 = 5 \times 10^{-6}$  is a penalty parameter,  $n$  is the total number of agents over network  $\mathcal{G}$ ,  $(a_i, b_i)$  represents the (feature, label) pair of datum  $i$ . We assume that each  $a_i \in [0.3, 0.4]^p$  is uniformly distributed, and  $b_i$  is chosen according to  $b_i = a_i^T z_i + \sigma_i$ , where  $z_i$  is a predefined parameter evenly distributed in  $[0, 10]^p$ , and  $\sigma_i$  is independent Gaussian noise with mean value of 0 and variance of 1. Given

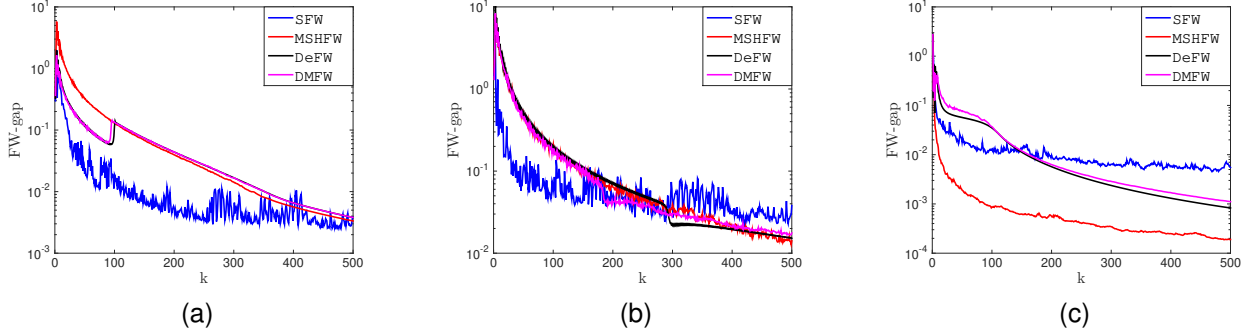


Fig. 3. The comparison between SFW, MSHFW, DeFW and DMFW on three datasets. (a) *covtype.binary* dataset. (b) *a9a* dataset. (c) *w8a* dataset.

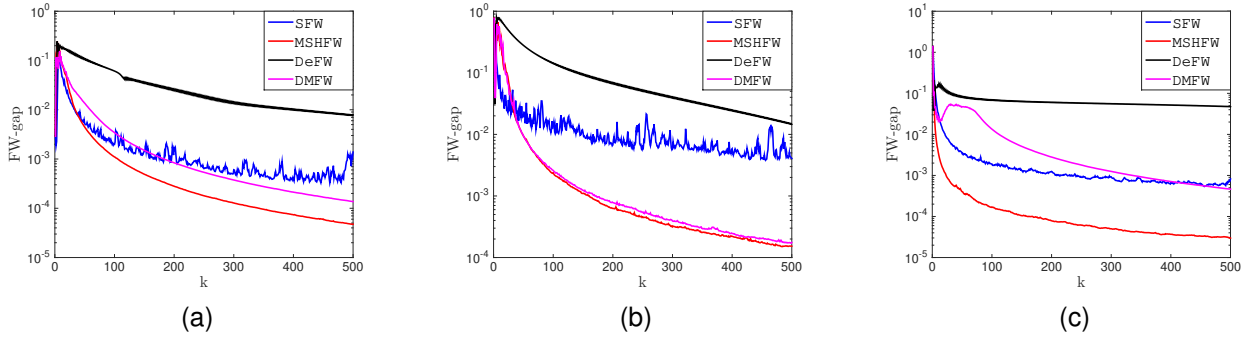


Fig. 4. The comparison between SFW, MSHFW, DeFW and DMFW on three datasets. (a) *covtype.binary* dataset. (b) *a9a* dataset. (c) *w8a* dataset.

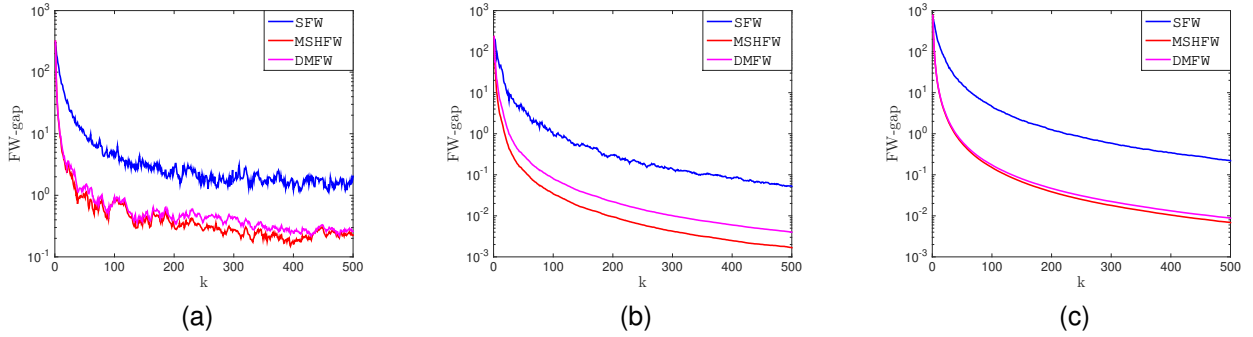


Fig. 5. The comparison between SFW, MSHFW and DMFW. (a)  $l_1$  norm ball constraint. (b)  $l_2$  norm ball constraint. (c)  $l_{\frac{5}{4}}$  norm ball constraint.

a pair  $(a_i, b_i)$ , agent  $i$  can compute an estimated gradient of  $F_i(x)$ :  $\nabla f_i(x, a_i, b_i) = 2(a_i^T x - b_i)a_i + 2\lambda_1 x$ . Choose  $p = 50$  and  $n = 50$ . In the experiments, three constraint sets  $\{x \mid \|x\|_1 \leq 5\}$ ,  $\{x \mid \|x\| \leq 5\}$  and  $\{x \mid \|x\|_{\frac{5}{4}} \leq 5\}$  are considered. As benchmarks, we choose SFW with step sizes  $\gamma_k = 2/(k+8)$  and  $\rho_k = 4/(k+8)^{\frac{2}{3}}$  as mentioned in [8], MSHFW with step sizes  $\gamma_k = 2/(k+1)$  and  $\eta_k = 2/(k+2)$  as mentioned in [9]. The step sizes of the algorithm DMFW are  $\gamma_k = 2/(k+1)$  and  $\eta_k = 2/(k+2)$ .

Fig. 4 shows the convergence performances of stochastic algorithms SFW, MSHFW and DMFW with the same parameters. In the simulation, SFW converges slower than DMFW and MSHFW in all cases, which is consistent with the theoretical results. DMFW performs comparable convergence performance to centralized algorithm MSHFW when the

constraint set are  $l_2$ - and  $l_{\frac{5}{4}}$ - norm balls.

## V. CONCLUSIONS

This paper proposed a distributed stochastic Frank-Wolfe algorithm by combining Nesterov's momentum with gradient tracking technique for stochastic optimization problems. As far as we know, this is the first distributed Frank-Wolfe algorithm for solving stochastic optimization problems. For convex objective functions, the proposed method achieves the convergence rate of  $\mathcal{O}(k^{-\frac{1}{2}})$ , coinciding with that of the centralized stochastic algorithms. In addition, the proposed algorithm achieves the convergence rate of  $\mathcal{O}(1/\log_2(k))$  for nonconvex problems under weaker assumptions than existing ones. The efficacy of the proposed algorithm was tested on binary classification problems with convex and nonconvex

objective functions. In our future research, it is interesting to consider stochastic optimization with nonsmooth objective functions by using the conditional gradient method.

## VI. APPENDIX

### A. Proof of Lemma 1

Before proving Lemma 1, we first give the following technical lemmas.

The following lemma is presented in Lemma 2 of the supplementary material for [9].

**Lemma 5.** Let  $\phi_k$  be a sequence of real numbers satisfying

$$\phi_k = \left(1 - \frac{A}{(k+t_0)^{r_1}}\right)\phi_{k-1} + \frac{B}{(k+t_0)^{r_2}},$$

for some  $r_1 \in [0, 1]$  such that  $r_1 \leq r_2 \leq 2r_1$ ,  $A > 1$  and  $B \geq 0$ . Then,  $\phi_k$  converges to zero at the following rate:

$$\phi_k \leq \frac{H}{(k+t_0+1)^{r_2-r_1}},$$

where  $H = \max\{\phi_0(t_0+1)^{r_2-r_1}, \frac{B}{A-1}\}$ .

**Lemma 6.** Consider Algorithm 1. Suppose Assumption 1 holds. For all  $k = 1, 2, \dots, K$ , the following relationships are established.

(a)  $\frac{1}{n} \sum_{i=1}^n s_{k+1}^i = \bar{y}_{k+1}$ ;

(b)  $\bar{x}_{k+1} = (1 - \eta_k)\bar{x}_k + \eta_k\bar{\theta}_k$ , where  $\bar{\theta}_k = \frac{1}{n} \sum_{i=1}^n \theta_k^i$ .

*Proof.* (a) From (4) of Algorithm 1, we have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n s_{k+1}^i &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^n c_{ij} s_k^j + y_{k+1}^i - y_k^i \right) \\ &= \frac{1}{n} \sum_{i=1}^n s_k^i + \bar{y}_{k+1} - \bar{y}_k \\ &= \frac{1}{n} \sum_{i=1}^n y_1^i - \bar{y}_1 + \bar{y}_{k+1} = \bar{y}_{k+1}, \end{aligned}$$

where the second equality is due to the fact that matrix  $C$  is doubly stochastic. Hence,  $\frac{1}{n} \sum_{i=1}^n s_{k+1}^i = \bar{y}_{k+1}$ .

(b) According to the definitions of  $\bar{x}_k$  and  $x_k^i$ , we have

$$\begin{aligned} \bar{x}_{k+1} &= \frac{1}{n} \sum_{i=1}^n \left[ (1 - \eta_k) \sum_{j=1}^n c_{ij} x_k^j + \eta_k \theta_k^i \right] \\ &= \frac{(1 - \eta_k)}{n} \sum_{i=1}^n x_k^i + \frac{\eta_k}{n} \sum_{i=1}^n \theta_k^i \\ &= (1 - \eta_k)\bar{x}_k + \eta_k\bar{\theta}_k, \end{aligned}$$

where the first equality is due to the fact that matrix  $C$  is doubly stochastic. Hence,  $\bar{x}_{k+1} = (1 - \eta_k)\bar{x}_k + \eta_k\bar{\theta}_k$ .  $\square$

Then we give the proof of Lemma 1.

*Proof.* We use mathematical induction to prove this lemma. It follows from the properties of Euclidean norm that  $\|\hat{x}_k^i - \bar{x}_k\|$

$\bar{x}_k\| \leq \max_{i \in \mathcal{N}} \|\hat{x}_k^i - \bar{x}_k\| \leq \left( \sum_{i=1}^n \|\hat{x}_k^i - \bar{x}_k\|^2 \right)^{\frac{1}{2}}$ . We next prove the following inequality by using induction on  $k$ ,

$$\left( \sum_{i=1}^n \|\hat{x}_k^i - \bar{x}_k\|^2 \right)^{\frac{1}{2}} \leq \frac{2C_1}{k+2} = \frac{2k_0\sqrt{n}D}{k+2} = C_1\eta_k. \quad (17)$$

It is obvious that inequality (17) holds for  $k = 1$  to  $k = k_0 - 2$ .

For induction step, we assume that (17) holds for some  $k \geq k_0 - 2$ . According to Lemma 6 (b) and (7), definitions of  $\hat{x}_k^i$  and  $\bar{x}_k$ , we have

$$\begin{aligned} &\sum_{i=1}^n \|\hat{x}_{k+1}^i - \bar{x}_{k+1}\|^2 \\ &= \sum_{i=1}^n \left\| \sum_{j=1}^n c_{ij}(1 - \eta_k)\hat{x}_k^j + \sum_{j=1}^n c_{ij}\eta_k\theta_k^j - (1 - \eta_k)\bar{x}_k - \eta_k\bar{\theta}_k \right\|^2 \\ &= \sum_{i=1}^n \left\| \sum_{j=1}^n c_{ij} \left[ (1 - \eta_k) \sum_{h=1}^n c_{jh}x_k^h + \eta_k\theta_k^j \right] - \frac{1}{n} \sum_{j=1}^n \left[ (1 - \eta_k) \sum_{h=1}^n c_{jh}x_k^h + \eta_k\theta_k^j \right] \right\|^2 \\ &\leq |\lambda|^2 \sum_{i=1}^n \left\| (1 - \eta_k)(\hat{x}_k^i - \bar{x}_k) + \eta_k(\theta_k^i - \bar{\theta}_k) \right\|^2, \end{aligned} \quad (18)$$

where  $\lambda$  is the second largest eigenvalue of  $C$ ; the last inequality follows from (8). Now we only need prove that  $\sum_{i=1}^n \|(1 - \eta_k)(\hat{x}_k^i - \bar{x}_k) + \eta_k(\theta_k^i - \bar{\theta}_k)\|^2$  has an upper bound. We have

$$\begin{aligned} &\sum_{i=1}^n \|(1 - \eta_k)(\hat{x}_k^i - \bar{x}_k) + \eta_k(\theta_k^i - \bar{\theta}_k)\|^2 \\ &\stackrel{(a)}{\leq} \sum_{i=1}^n \left[ (1 - \eta_k)^2 \|\hat{x}_k^i - \bar{x}_k\|^2 + \eta_k^2 D^2 + 2\eta_k(1 - \eta_k)D \|\hat{x}_k^i - \bar{x}_k\| \right] \\ &\stackrel{(b)}{\leq} \sum_{i=1}^n \left( \|\hat{x}_k^i - \bar{x}_k\|^2 + 2\eta_k D \|\hat{x}_k^i - \bar{x}_k\| + \eta_k^2 D^2 \right) \\ &\stackrel{(c)}{\leq} C_1^2 \eta_k^2 + n\eta_k^2 D^2 + 2\eta_k D \sqrt{n} \sqrt{\sum_{i=1}^n \|\hat{x}_k^i - \bar{x}_k\|^2} \\ &\leq C_1^2 \eta_k^2 + n\eta_k^2 D^2 + 2D\eta_k^2 \sqrt{n} C_1 \\ &= \eta_k^2 (C_1 + C_1 k_0^{-1})^2 \\ &= \left( \frac{k_0 + 1}{k_0} C_1 \eta_k \right)^2, \end{aligned} \quad (19)$$

where (a) follows from Assumption 3; (b) follows from  $1 - \eta_k \leq 1$ ; (c) is due to  $\sum_{i=1}^n \|\hat{x}_k^i - \bar{x}_k\| \leq \sqrt{n} \sqrt{\sum_{i=1}^n \|\hat{x}_k^i - \bar{x}_k\|^2}$  and the induction hypothesis (17). Substituting (19),  $\lambda \leq \left( \frac{k_0}{k_0+1} \right)^2$  and  $\eta_k = \frac{2}{k+2}$  into (18), we arrive at

$$\begin{aligned} &\sum_{i=1}^n \|\hat{x}_{k+1}^i - \bar{x}_{k+1}\|^2 \\ &\leq \left( \frac{2(k_0 + 1)}{k_0(k+2)} \left( \frac{k_0}{k_0+1} \right)^2 C_1 \right)^2 \\ &\leq \left( \frac{2(k+2)}{(k+2+1)(k+2)} C_1 \right)^2 = C_1^2 \eta_{k+1}^2 \end{aligned} \quad (20)$$



where the second inequality is because of the monotonically increasing property of function  $g(v) = v/(1+v)$  with respect to  $v$  over  $[0, \infty)$ . It follows from (20) that  $\sum_{i=1}^n \|\hat{x}_{k+1}^i - \bar{x}_{k+1}\|_2 \leq C_1 \eta_{k+1}$ , that is, (17) holds for iteration  $k+1$ . Hence,  $\|\hat{x}_k^i - \bar{x}_k\| \leq 2C_1/(k+2)$  for all  $k \geq 1$ .  $\square$

### B. Technical Lemmas

**Lemma 7.** *Suppose Assumptions 1-3 hold. Let  $\eta_k = \frac{2}{k+2}$ . Then, for any  $i \in \mathcal{N}$  and  $k \geq 1$ , we have*

$$\|\hat{x}_{k+1}^i - \hat{x}_k^i\| \leq \frac{2(D+2C_1)}{k+2}, \quad (21)$$

where  $C_1 = k_0 \sqrt{nD}$ .

*Proof.* It follows from the definition of  $\hat{x}_k^i$  that

$$\begin{aligned} & \|\hat{x}_{k+1}^i - \hat{x}_k^i\| \\ & \leq \sum_{j=1}^n c_{ij} (\|x_{k+1}^j - \hat{x}_k^j\| + \|\hat{x}_k^j - \hat{x}_k^i\|) \\ & \stackrel{(a)}{=} \sum_{j=1}^n c_{ij} (\|\eta_k \theta_k^j - \eta_k \hat{x}_k^j\| + \|\hat{x}_k^j - \bar{x}_k + \bar{x}_k - \hat{x}_k^i\|) \\ & \stackrel{(b)}{\leq} \sum_{j=1}^n c_{ij} (\|\hat{x}_k^j - \bar{x}_k\| + \|\hat{x}_k^i - \bar{x}_k\|) + \sum_{j=1}^n c_{ij} \|\eta_k (\theta_k^j - \hat{x}_k^j)\| \\ & \stackrel{(c)}{\leq} \sum_{j=1}^n c_{ij} (\eta_k D + 2C_1 \eta_k) = \frac{2(D+2C_1)}{k+2}, \end{aligned} \quad (22)$$

where (a) follows from (7); (b) holds for the triangle inequality; (c) is because of Assumption 3 and Lemma 1.  $\square$

**Lemma 8.** *Suppose Assumptions 1-5 hold. Choose the step sizes  $\gamma_k = \frac{2}{k+1}$  and  $\eta_k = \frac{2}{k+2}$ . Then, for any  $i \in \mathcal{N}$  and  $k \geq 1$ , the variable  $y_k^i$  of Algorithm 1 satisfies the following bound:*

$$\mathbb{E}[\|y_k^i\|] \leq \psi, \quad (23)$$

where  $\psi = \max_{i \in \mathcal{N}} \{\|y_1^i\|, 2G + 2L(D+2C_1)\}$ .

*Proof.* It is obvious that (23) holds when  $k=1$ . Next let's discuss the case when  $k \geq 2$ . It follows from (3) of Algorithm 1 that

$$\begin{aligned} & \mathbb{E}[\|y_k^i\|] \\ & = \mathbb{E}[\|(1-\gamma_k)y_{k-1}^i + \nabla f_i(\hat{x}_k^i, \xi_k^i) - (1-\gamma_k)\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|] \\ & \leq (1-\gamma_k)\mathbb{E}[\|y_{k-1}^i\|] + \mathbb{E}[\|\nabla f_i(\hat{x}_k^i, \xi_k^i) - (1-\gamma_k)\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|] \\ & \leq (1-\gamma_k)\mathbb{E}[\|y_{k-1}^i\|] + \gamma_k \mathbb{E}[\|\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|] \\ & \quad + \mathbb{E}[\|\nabla f_i(\hat{x}_k^i, \xi_k^i) - \nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|] \\ & \stackrel{(a)}{\leq} (1-\gamma_k)\mathbb{E}[\|y_{k-1}^i\|] + G\gamma_k + L\mathbb{E}[\|\hat{x}_k^i - \hat{x}_{k-1}^i\|] \\ & \stackrel{(b)}{\leq} (1-\gamma_k)\mathbb{E}[\|y_{k-1}^i\|] + G\gamma_k + L(D+2C_1)\eta_{k-1} \\ & = (1-\frac{2}{k+1})\mathbb{E}[\|y_{k-1}^i\|] + \frac{2G+2L(D+2C_1)}{k+1} \end{aligned} \quad (24)$$

where (a) is due to the  $L$ -smooth property of function  $f_i$  and Fact 2; (b) follows from (21). By using Lemma 5 ( $r_1=1$ ,  $r_2=1$ ,  $A=2$  and  $B=2G+2L(D+2C_1)$ ), we obtain  $\mathbb{E}[\|y_k^i\|] \leq \psi$ , where  $\psi = \max_{i \in \mathcal{N}} \{\|y_1^i\|, 2G+2L(D+2C_1)\}$ .  $\square$

**Lemma 9.** *Suppose Assumptions 1-5 hold. Choose the step sizes  $\gamma_k = \frac{2}{k+1}$  and  $\eta_k = \frac{2}{k+2}$ . Then, for any  $i \in \mathcal{N}$  and  $k \geq 1$ , it holds*

$$\mathbb{E}[\|y_k^i\|^2] \leq \hat{\psi}, \quad (25)$$

where  $\hat{\psi} = \max_{i \in \mathcal{N}} \{\|y_1^i\|^2, 4L(D+2C_1)\psi + 4G\psi + 8G^2 + 8L^2(D+2C_1)^2\}$ .

*Proof.* It is obvious that (25) holds when  $k=1$ . Next let's discuss the case when  $k \geq 2$ . From (3) of Algorithm 1, we have

$$\begin{aligned} & \|y_k^i\|^2 \leq (1-\gamma_k)^2 \|y_{k-1}^i\|^2 + 2(1-\gamma_k) \|\nabla f_i(\hat{x}_k^i, \xi_k^i) \\ & \quad - (1-\gamma_k) \nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\| \|y_{k-1}^i\| + \|\nabla f_i(\hat{x}_k^i, \xi_k^i) \\ & \quad - (1-\gamma_k) \nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|^2 \\ & \stackrel{(a)}{\leq} \|\nabla f_i(\hat{x}_k^i, \xi_k^i) - \nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i) + \gamma_k \nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|^2 \\ & \quad + (1-\gamma_k)^2 \|y_{k-1}^i\|^2 + 2(1-\gamma_k) (\|\nabla f_i(\hat{x}_k^i, \xi_k^i) \\ & \quad - \nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\| + \gamma_k \|\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|) \|y_{k-1}^i\| \\ & \stackrel{(b)}{\leq} 2L^2 \|\hat{x}_k^i - \hat{x}_{k-1}^i\|^2 + 2\gamma_k^2 \|\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|^2 + (1-\gamma_k) \|y_{k-1}^i\|^2 \\ & \quad + 2(L\|\hat{x}_k^i - \hat{x}_{k-1}^i\| + \gamma_k \|\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|) \|y_{k-1}^i\| \\ & \stackrel{(c)}{\leq} (1-\gamma_k) \|y_{k-1}^i\|^2 + 2\gamma_k^2 \|\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|^2 + 2[L(D \\ & \quad + 2C_1)\eta_{k-1} + \gamma_k \|\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|] \|y_{k-1}^i\| \\ & \quad + 2L^2(D+2C_1)^2 \eta_{k-1}^2 \end{aligned} \quad (26)$$

where (a) is due to the triangle inequality; (b) holds because of the fact  $1-\gamma_k \leq 1$  and the  $L$ -smooth property of function  $f_i$ ; (c) follows from (21). Taking the conditional expectation of (26) on  $\mathcal{F}_k$ , we obtain

$$\begin{aligned} \mathbb{E}_k[\|y_k^i\|^2] & \leq (1-\gamma_k) \|y_{k-1}^i\|^2 + 2L(D+2C_1)\eta_{k-1} \|y_{k-1}^i\| \\ & \quad + 2\gamma_k \mathbb{E}_k[\|\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|] \|y_{k-1}^i\| \\ & \quad + 2\gamma_k^2 \mathbb{E}_k[\|\nabla f_i(\hat{x}_{k-1}^i, \xi_{k-1}^i)\|^2] \\ & \quad + 2L^2(D+2C_1)^2 \eta_{k-1}^2 \\ & \leq (1-\gamma_k) \|y_{k-1}^i\|^2 + 2L(D+2C_1)\eta_{k-1} \|y_{k-1}^i\| + 2\gamma_k^2 G^2 \\ & \quad + 2\gamma_k G \|y_{k-1}^i\| + 2L^2(D+2C_1)^2 \eta_{k-1}^2 \end{aligned} \quad (27)$$

where the last inequality follows from Fact 2. Taking the full expectation of (27) and choosing the step sizes  $\gamma_k = \frac{2}{k+1}$ ,  $\eta_k = \frac{2}{k+2}$ , we arrive at

$$\begin{aligned} \mathbb{E}[\|y_k^i\|^2] & \leq (1-\gamma_k) \mathbb{E}[\|y_{k-1}^i\|^2] + 2L(D+2C_1)\eta_{k-1} \mathbb{E}[\|y_{k-1}^i\|] \\ & \quad + 2\gamma_k G \mathbb{E}[\|y_{k-1}^i\|] + 2\gamma_k^2 G^2 + 2L^2(D+2C_1)^2 \eta_{k-1}^2 \\ & \leq (1-\gamma_k) \mathbb{E}[\|y_{k-1}^i\|^2] + 2L(D+2C_1)\psi \eta_{k-1} + 2\gamma_k G \psi \\ & \quad + 2\gamma_k^2 G^2 + 2L^2(D+2C_1)^2 \eta_{k-1}^2 \\ & \leq (1-\frac{2}{k+1}) \mathbb{E}[\|y_{k-1}^i\|^2] + \frac{4L(D+2C_1)\psi + 4G\psi}{k+1} \\ & \quad + \frac{8G^2 + 8L^2(D+2C_1)^2}{k+1} \end{aligned}$$

where the second inequality follows from (23) of Lemma 8. By using Lemma 5 ( $r_1=1$ ,  $r_2=1$ ,  $A=2$  and  $B=4L(D+2C_1)\psi + 4G\psi + 8G^2 + 8L^2(D+2C_1)^2$ ), we obtain  $\mathbb{E}[\|y_k^i\|^2] \leq \hat{\psi}$ , where  $\hat{\psi} = \max_{i \in \mathcal{N}} \{\|y_1^i\|^2, 4L(D+2C_1)\psi + 4G\psi + 8G^2 + 8L^2(D+2C_1)^2\}$ .  $\square$

### C. Proof of Lemma 2

*Proof.* It follows from the properties of Euclidean norm that  $\mathbb{E}[\|p_k^i - \bar{y}_k\|^2] \leq \mathbb{E}[\sum_{i=1}^n \|p_k^i - \bar{y}_k\|^2]$ . Then, proving inequality  $\mathbb{E}[\|p_k^i - \bar{y}_k\|^2] \leq 4C_2/(k+2)^2$  can be transformed into proving the inequality

$$\begin{aligned} & \mathbb{E}\left[\sum_{i=1}^n \|p_k^i - \bar{y}_k\|^2\right] \\ & \leq C_2 \eta_k^2 = \frac{k_0^3 (4n)^{k_0+1} (12L^2(D+2C_1)^2 + 12(G^2 + \hat{\psi}))}{(k+2)^2} \end{aligned} \quad (28)$$

by using induction on  $k$ . We first prove that (28) holds for all  $1 \leq k \leq k_0 - 2$ . According to (4) and (5) of Algorithm 1, we have  $p_k^i = s_{k+1}^i + y_k^i - y_{k+1}^i$ . Therefore,

$$\begin{aligned} & \mathbb{E}\left[\sum_{i=1}^n \|p_k^i - \bar{y}_k\|^2\right] \\ & \leq \mathbb{E}\left[\sum_{i=1}^n \left(4\|s_{k+1}^i\|^2 + 4\|y_k^i\|^2 + 4\|y_{k+1}^i\|^2 + \frac{4}{n^2} \left\|\sum_{j=1}^n y_k^j\right\|^2\right)\right] \\ & \leq 4 \sum_{i=1}^n \mathbb{E}[\|s_{k+1}^i\|^2] + 8n\hat{\psi} + 4n\hat{\psi}, \end{aligned} \quad (29)$$

where the first inequality holds for the fact that  $\|\sum_{i=1}^n z_i\|^2 \leq n \sum_{i=1}^n \|z_i\|^2$  ( $z_i$  is an arbitrary vector for  $\forall i \in \{1, 2, \dots, n\}$ ) and the last inequality is due to (25). The first term of RHS of (29) can be written as

$$\begin{aligned} \mathbb{E}[\|s_{k+1}^i\|^2] &= \mathbb{E}\left[\left\|\sum_{j=1}^n c_{ij} s_k^j + y_{k+1}^i - y_k^i\right\|^2\right] \\ &\leq 3\mathbb{E}\left[\left\|\sum_{j=1}^n c_{ij} s_k^j\right\|^2 + \|y_{k+1}^i\|^2 + \|y_k^i\|^2\right] \\ &\stackrel{(a)}{\leq} 3\left(\mathbb{E}\left[\left\|\sum_{j=1}^n c_{ij} s_k^j\right\|^2\right] + 2\hat{\psi}\right) \\ &\leq 3\left(n\mathbb{E}\left[\sum_{j=1}^n \|c_{ij} s_k^j\|^2\right] + 2\hat{\psi}\right) \\ &\stackrel{(b)}{\leq} 3n \sum_{j=1}^n c_{ij} \mathbb{E}[\|s_k^j\|^2] + 6\hat{\psi} \\ &\stackrel{(c)}{\leq} (3n)^k G^2 + 6k(3n)^{(k-1)} \hat{\psi} \end{aligned} \quad (30)$$

where (a) follows from (25); (b) is due to the fact that  $\|\sum_{i=1}^n z_i\|^2 \leq n \sum_{i=1}^n \|z_i\|^2$  ( $z_i$  is an arbitrary vector for  $\forall i \in \{1, 2, \dots, n\}$ ) and the fact that  $c_{ij} \leq 1$  for all  $i, j \in \mathcal{N}$ ; (c) is because  $\mathbb{E}[\|s_1^j\|^2] = \mathbb{E}[\|\nabla f_j(\hat{x}_1^j, \xi_1^j)\|^2] \leq G^2$ . Now substituting (30) into (29), we arrive at

$$\begin{aligned} & \mathbb{E}\left[\sum_{i=1}^n \|p_k^i - \bar{y}_k\|^2\right] \\ & \leq 4n(3n)^k G^2 + 24nk(3n)^{(k-1)} \hat{\psi} + 8n\hat{\psi} + 4n\hat{\psi} \\ & \leq (4n)^{k+1} G^2 + 6k(4n)^k \hat{\psi} + 12n\hat{\psi}, \end{aligned}$$

that is,  $\mathbb{E}[\sum_{i=1}^n \|p_k^i - \bar{y}_k\|^2] \leq (4n)^{k_0-1} G^2 + 6(k_0 - 2)(4n)^{k_0-2} \hat{\psi} + 12n\hat{\psi}$  holds for all  $1 \leq k \leq k_0 - 2$ . Hence, the inequality (28) is obviously true for  $k = 1$  to  $k = k_0 - 2$ .

For induction step, we assume that (28) holds for some  $k \geq k_0 - 2$ . Define the slack variable  $\Delta y_{k+1}^i := y_{k+1}^i - y_k^i$ . Then, we observe that  $s_{k+1}^i = \Delta y_{k+1}^i + p_k^i$  due to the definition of  $s_{k+1}^i$  and  $p_{k+1}^i$ . From (8) and Lemma 6 (a), we have

$$\begin{aligned} & \mathbb{E}\left[\sum_{i=1}^n \|p_{k+1}^i - \bar{y}_{k+1}\|^2\right] \\ & \leq \mathbb{E}\left[|\lambda|^2 \sum_{i=1}^n \|\Delta y_{k+1}^i + p_k^i - \bar{y}_{k+1}\|^2\right]. \end{aligned} \quad (31)$$

Define  $\Delta Y_{k+1} := \bar{y}_{k+1} - \bar{y}_k$ . Then,

$$\begin{aligned} & \sum_{i=1}^n \|\Delta y_{k+1}^i + p_k^i - \bar{y}_{k+1}\|^2 \\ &= \sum_{i=1}^n \|p_k^i - \bar{y}_k + \Delta y_{k+1}^i - \Delta Y_{k+1}\|^2 \\ &\leq \sum_{i=1}^n [\|p_k^i - \bar{y}_k\|^2 + \|\Delta y_{k+1}^i - \Delta Y_{k+1}\|^2 \\ &\quad + 2\|\Delta y_{k+1}^i - \Delta Y_{k+1}\| \|p_k^i - \bar{y}_k\|]. \end{aligned} \quad (32)$$

According to the definition of  $\Delta y_{k+1}^i$ , we get

$$\begin{aligned} & \mathbb{E}[\|\Delta y_{k+1}^i\|^2] = \mathbb{E}[\|y_{k+1}^i - y_k^i\|^2] \\ &= \mathbb{E}[\|\nabla f_i(\hat{x}_{k+1}^i, \xi_{k+1}^i) - \nabla f_i(\hat{x}_k^i, \xi_{k+1}^i) \\ &\quad + \gamma_{k+1}(\nabla f_i(\hat{x}_k^i, \xi_{k+1}^i) - y_k^i)\|^2] \\ &\stackrel{(a)}{\leq} 3\mathbb{E}[\|\nabla f_i(\hat{x}_{k+1}^i, \xi_{k+1}^i) - \nabla f_i(\hat{x}_k^i, \xi_{k+1}^i)\|^2] \\ &\quad + 3\gamma_{k+1}^2 \mathbb{E}[\|\nabla f_i(\hat{x}_k^i, \xi_{k+1}^i)\|^2] + 3\gamma_{k+1}^2 \mathbb{E}[\|y_k^i\|^2] \\ &\stackrel{(b)}{\leq} 3L^2 \mathbb{E}[\|\hat{x}_{k+1}^i - \hat{x}_k^i\|^2] + 3\gamma_{k+1}^2 (G^2 + \hat{\psi}) \\ &\stackrel{(c)}{\leq} 3L^2 (D + 2C_1)^2 \eta_k^2 + 3(G^2 + \hat{\psi}) \gamma_{k+1}^2, \end{aligned} \quad (33)$$

where (a) is due to the fact that  $\|\sum_{i=1}^n z_i\|^2 \leq n \sum_{i=1}^n \|z_i\|^2$  ( $z_i$  is an arbitrary vector for  $\forall i \in \{1, 2, \dots, n\}$ ); (b) holds because of the  $L$ -smooth property of function  $f_i(x, \xi)$ , Fact 3 and (25); (c) follows from (21). It follows from the definition of  $\Delta Y_{k+1}$  that

$$\begin{aligned} & \mathbb{E}[\|\Delta y_{k+1}^i - \Delta Y_{k+1}\|^2] \\ &= \mathbb{E}\left[\left\|\Delta y_{k+1}^i - \frac{1}{n} \sum_{i=1}^n \Delta y_{k+1}^i\right\|^2\right] \\ &= \mathbb{E}\left[\left\|\left(1 - \frac{1}{n}\right) \Delta y_{k+1}^i - \frac{1}{n} \sum_{j \neq i} \Delta y_{k+1}^j\right\|^2\right] \\ &\leq 2\left(1 - \frac{1}{n}\right) \mathbb{E}[\|\Delta y_{k+1}^i\|^2] + \frac{2}{n} \sum_{j \neq i} \mathbb{E}[\|\Delta y_{k+1}^j\|^2] \\ &\stackrel{(a)}{\leq} 4\left(1 - \frac{1}{n}\right) [3L^2 (D + 2C_1)^2 \eta_k^2 + 3(G^2 + \hat{\psi}) \gamma_{k+1}^2] \\ &\stackrel{(b)}{\leq} [12L^2 (D + 2C_1)^2 + 12(G^2 + \hat{\psi})] \eta_k^2, \end{aligned} \quad (34)$$

where (a) follows from (33); (b) is by the fact  $1 - \frac{1}{n} \leq 1$ . Substituting (34) into (32) and taking the full expectation of (32), we obtain

$$\mathbb{E}\left[\sum_{i=1}^n \|\Delta y_{k+1}^i + p_k^i - \bar{y}_{k+1}\|^2\right]$$

$$\begin{aligned}
&\stackrel{(a)}{\leq} \mathbb{E} \left[ \sum_{i=1}^n \|p_k^i - \bar{y}_k\|^2 \right] + n[12L^2(D + 2C_1)^2 + 12(G^2 \\
&\quad + \hat{\psi})]\eta_k^2 + 2 \sum_{i=1}^n (\mathbb{E}[\|\Delta y_{k+1}^i - \Delta Y_{k+1}\|^2])^{\frac{1}{2}} \\
&\quad \times (\mathbb{E}[\|p_k^i - \bar{y}_k\|^2])^{\frac{1}{2}} \\
&\stackrel{(b)}{\leq} C_2\eta_k^2 + n[12L^2(D + 2C_1)^2 + 12(G^2 + \hat{\psi})]\eta_k^2 \\
&\quad + 2n\eta_k^2 \sqrt{[12L^2(D + 2C_1)^2 + 12(G^2 + \hat{\psi})]C_2} \\
&\stackrel{(c)}{\leq} \eta_k^2 \left( C_2 + n^2[12L^2(D + 2C_1)^2 + 12(G^2 + \hat{\psi})] \right) \\
&\quad + 2n \sqrt{[12L^2(D + 2C_1)^2 + 12(G^2 + \hat{\psi})]C_2} \\
&\leq \eta_k^2 \left( \sqrt{C_2} + \frac{\sqrt{C_2}}{k_0} \right)^2 \\
&= \left( \frac{k_0 + 1}{k_0} \sqrt{C_2} \eta_k \right)^2, \tag{35}
\end{aligned}$$

where (a) is due to the Hölder's inequality ( $\mathbb{E}[|XY|] \leq (\mathbb{E}[|X|^2])^{\frac{1}{2}} (\mathbb{E}[|Y|^2])^{\frac{1}{2}}$ ); (b) follows from (34) and the induction hypothesis (28); (c) is due to the fact  $n \geq 1$ . Substituting (35),  $|\lambda| \leq \left(\frac{k_0}{k_0+1}\right)^2$  and  $\eta_k = \frac{2}{k+2}$  into (31), we obtain

$$\begin{aligned}
\mathbb{E} \left[ \sum_{i=1}^n \|p_{k+1}^i - \bar{y}_{k+1}\|^2 \right] &\leq \left( \frac{2k_0}{(k+2)(k_0+1)} \sqrt{C_2} \right)^2 \\
&\leq \left( \frac{2(k+2)}{(k+3)(k+2)} \sqrt{C_2} \right)^2 = C_2\eta_{k+1}^2,
\end{aligned}$$

where the last inequality is because of the monotonically increasing property of function  $g(v) = v/(1+v)$  with respect to  $v$  over  $[0, \infty)$ . Hence, (28) holds for iteration  $k+1$ . Therefore,  $\mathbb{E}[\|p_k^i - \bar{y}_k\|^2] \leq 4C_2/(k+2)^2$  holds for all  $k \geq 1$ .  $\square$

### D. Proof of Lemma 3

*Proof.* (a) It follows from the definition of  $\bar{y}_k$  that

$$\begin{aligned}
\bar{P}_k - \bar{y}_k &= \bar{P}_k - \frac{1}{n} \sum_{i=1}^n y_k^i \\
&= \bar{P}_k - (1 - \gamma_k)\bar{y}_{k-1} - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i) \\
&\quad + \frac{1 - \gamma_k}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_{k-1}^i, \xi_k^i). \tag{36}
\end{aligned}$$

Introducing  $(1 - \gamma_k)\bar{P}_{k-1}$  into (36) and taking norm square of (36), we arrive at

$$\begin{aligned}
\|\bar{P}_k - \bar{y}_k\|^2 &= \|(1 - \gamma_k)(\bar{P}_{k-1} - \bar{y}_{k-1}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i) + \bar{P}_k \\
&\quad + (1 - \gamma_k) \left( \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_{k-1}^i, \xi_k^i) - \bar{P}_{k-1} \right)\|^2. \tag{37}
\end{aligned}$$

Taking the conditional expectation of (37) on  $\mathcal{F}_k$ , we obtain

$$\mathbb{E}_k[\|\bar{P}_k - \bar{y}_k\|^2]$$

$$\begin{aligned}
&= (1 - \gamma_k)^2 \|\bar{P}_{k-1} - \bar{y}_{k-1}\|^2 + 2(1 - \gamma_k)(\bar{P}_{k-1} - \bar{y}_{k-1})^\top \left( \mathbb{E}_k \left[ \bar{P}_k \right. \right. \\
&\quad \left. \left. - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i) \right] + (1 - \gamma_k) \mathbb{E}_k \left[ \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_{k-1}^i, \xi_k^i) \right. \right. \\
&\quad \left. \left. - \bar{P}_{k-1} \right] \right) + \mathbb{E}_k \left[ \left\| \bar{P}_k - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i) + (1 - \right. \right. \\
&\quad \left. \left. \gamma_k) \left( \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_{k-1}^i, \xi_k^i) - \bar{P}_{k-1} \right) \right\|^2 \right] \\
&= (1 - \gamma_k)^2 \|\bar{P}_{k-1} - \bar{y}_{k-1}\|^2 + \mathbb{E}_k \left[ \left\| \bar{P}_k - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i) \right. \right. \\
&\quad \left. \left. + (1 - \gamma_k) \left( \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_{k-1}^i, \xi_k^i) - \bar{P}_{k-1} \right) \right\|^2 \right], \tag{38}
\end{aligned}$$

where the last equality is because  $\mathbb{E}_k[\bar{P}_k - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i)] = \mathbb{E}_k[\frac{1}{n} \sum_{i=1}^n \nabla F_i(\hat{x}_k^i) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_k[\nabla F_i(\hat{x}_k^i) - \nabla f_i(\hat{x}_k^i, \xi_k^i)] = 0$ . Adding and subtracting  $\gamma_k(\bar{P}_k - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i))$  in the second term of RHS of (38), we have

$$\begin{aligned}
&\mathbb{E}_k \left[ \left\| (1 - \gamma_k) \left( \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_{k-1}^i, \xi_k^i) - \bar{P}_{k-1} + \bar{P}_k \right. \right. \right. \\
&\quad \left. \left. - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i) \right) + \gamma_k \left( \bar{P}_k - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i) \right) \right\|^2 \right] \\
&\stackrel{(a)}{\leq} \frac{3(1 - \gamma_k)^2}{n} \sum_{i=1}^n \mathbb{E}_k[\|\nabla f_i(\hat{x}_{k-1}^i, \xi_k^i) - \nabla f_i(\hat{x}_k^i, \xi_k^i)\|^2] \\
&\quad + 3(1 - \gamma_k)^2 \left\| \frac{1}{n} \sum_{i=1}^n (\nabla F_i(\hat{x}_k^i) - \nabla F_i(\hat{x}_{k-1}^i)) \right\|^2 \\
&\quad + 3\gamma_k^2 \mathbb{E}_k \left[ \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(\hat{x}_k^i, \xi_k^i) - \bar{P}_k \right\|^2 \right] \\
&\stackrel{(b)}{\leq} \frac{3(1 - \gamma_k)^2}{n} \sum_{i=1}^n \mathbb{E}_k[L^2 \|\hat{x}_k^i - \hat{x}_{k-1}^i\|^2] + 3\gamma_k^2 \delta^2 \\
&\quad + \frac{3(1 - \gamma_k)^2 L^2}{n} \sum_{i=1}^n \|\hat{x}_k^i - \hat{x}_{k-1}^i\|^2 \\
&\stackrel{(c)}{\leq} 6L^2(D + 2C_1)^2 \eta_{k-1}^2 + 3\gamma_k^2 \delta^2, \tag{39}
\end{aligned}$$

where (a) is due to the fact that  $\|\sum_{i=1}^n z_i\|^2 \leq n \sum_{i=1}^n \|z_i\|^2$  ( $z_i$  is an arbitrary vector for  $\forall i \in \{1, 2, \dots, n\}$ ) and  $\bar{P}_k = \frac{1}{n} \sum_{i=1}^n \nabla F_i(\hat{x}_k^i)$ ; (b) is due to Assumption 5, the  $L$ -smooth property of function  $f_i$  and  $F_i$ ; (c) is obtained by (21) and the fact  $(1 - \gamma_k)^2 \leq 1$ . Substituting (39) into (38), we obtain

$$\begin{aligned}
\mathbb{E}_k[\|\bar{P}_k - \bar{y}_k\|^2] &\leq (1 - \gamma_k) \|\bar{P}_{k-1} - \bar{y}_{k-1}\|^2 \\
&\quad + 6L^2(D + 2C_1)^2 \eta_{k-1}^2 + 3\gamma_k^2 \delta^2. \tag{40}
\end{aligned}$$

(b) It is obvious that  $\mathbb{E}[\|\bar{P}_1 - \bar{y}_1\|^2] = 0 \leq \frac{C_3}{k+2}$ , that is, (12) holds when  $k = 1$ . Now, let's discuss the case when  $k \geq 2$ . Taking full expectation of (40) and choosing the step sizes  $\gamma_k = \frac{2}{k+1}$ ,  $\eta_k = \frac{2}{k+2}$ , we have

$$\mathbb{E}[\|\bar{P}_k - \bar{y}_k\|^2]$$

$$\leq (1 - \frac{2}{k+1})\mathbb{E}[\|\bar{P}_{k-1} - \bar{y}_{k-1}\|^2] + \frac{24L^2(D + 2C_1)^2 + 12\delta^2}{(k+1)^2}.$$

By using Lemma 5 ( $r_1 = 1$ ,  $r_2 = 2$ ,  $A = 2$  and  $B = 24L^2(D + 2C_1)^2 + 12\delta^2$ ), we obtain  $\mathbb{E}[\|\bar{P}_k - \bar{y}_k\|^2] \leq \frac{C_3}{k+2}$ , where  $C_3 = 24L^2(D + 2C_1)^2 + 12\delta^2$ .  $\square$

### E. Proof of Lemma 4

*Proof.* Adding and subtracting  $\bar{P}_k$  and  $\bar{y}_k$  to  $\|\nabla F(\bar{x}_k) - p_k^i\|^2$ , we have

$$\begin{aligned} & \mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|^2] \\ &= \mathbb{E}[\|\nabla F(\bar{x}_k) - \bar{P}_k + \bar{P}_k - \bar{y}_k + \bar{y}_k - p_k^i\|^2] \\ &\leq 3\mathbb{E}[\|\nabla F(\bar{x}_k) - \bar{P}_k\|^2] + 3\mathbb{E}[\|\bar{P}_k - \bar{y}_k\|^2] \\ &\quad + 3\mathbb{E}[\|\bar{y}_k - p_k^i\|^2], \end{aligned} \quad (41)$$

where the inequality is due to the fact  $\|\sum_{i=1}^n z_i\|^2 \leq n \sum_{i=1}^n \|z_i\|^2$ ,  $z_i (i \in \{1, 2, \dots, n\})$  is an arbitrary vector. The first term of RHS of (41) can be written as

$$\begin{aligned} & 3\mathbb{E}[\|\nabla F(\bar{x}_k) - \bar{P}_k\|^2] \\ &= 3\mathbb{E}\left[\left\|\frac{1}{n} \sum_{i=1}^n \nabla F_i(\bar{x}_k) - \frac{1}{n} \sum_{i=1}^n \nabla F_i(\hat{x}_k^i)\right\|^2\right] \\ &\leq 3\mathbb{E}\left[\left(\frac{1}{n} \sum_{i=1}^n \|\nabla F_i(\bar{x}_k) - \nabla F_i(\hat{x}_k^i)\|\right)^2\right] \\ &\stackrel{(a)}{\leq} 3\frac{L^2}{n} \sum_{i=1}^n \mathbb{E}[\|\bar{x}_k - \hat{x}_k^i\|^2] \stackrel{(b)}{\leq} 3L^2 C_1^2 \eta_k^2, \end{aligned} \quad (42)$$

where (a) is due to the  $L$ -smooth property of function  $F_i$ ; (b) follows from (10). Substituting (42), (12) and (11) into (41), we obtain

$$\begin{aligned} \mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|^2] &\leq \frac{12L^2 C_1^2}{(k+2)^2} + \frac{3C_3}{k+2} + \frac{12C_2}{(k+2)^2} \\ &\leq \frac{12L^2 C_1^2 + 3C_3 + 12C_2}{k+2}. \end{aligned}$$

### F. Proof of Theorem 1

*Proof.* Since the function  $F(x)$  is  $L$ -smooth, we have

$$\begin{aligned} F(\bar{x}_{k+1}) &\leq F(\bar{x}_k) + \nabla^T F(\bar{x}_k)(\bar{x}_{k+1} - \bar{x}_k) + \frac{L}{2} \|\bar{x}_{k+1} - \bar{x}_k\|^2 \\ &= F(\bar{x}_k) + \eta_k \nabla^T F(\bar{x}_k)(\bar{\theta}_k - \bar{x}_k) + \frac{L\eta_k^2}{2} \|\bar{\theta}_k - \bar{x}_k\|^2 \\ &\leq F(\bar{x}_k) + \eta_k \nabla^T F(\bar{x}_k)(\bar{\theta}_k - \bar{x}_k) + \frac{L\eta_k^2}{2} D^2, \end{aligned} \quad (43)$$

where the first equation is due to Lemma 6 (b) and the last inequality follows from Assumption 3. From the definition of  $\bar{\theta}_k$ , the second term of RHS of (43) can be written as

$$\begin{aligned} \nabla^T F(\bar{x}_k)(\bar{\theta}_k - \bar{x}_k) &= \frac{1}{n} \sum_{i=1}^n \nabla^T F(\bar{x}_k)(\theta_k^i - \bar{x}_k) \\ &= \frac{1}{n} \sum_{i=1}^n (\nabla F(\bar{x}_k) - p_k^i + p_k^i)^T (\theta_k^i - \bar{x}_k) \end{aligned}$$

$$\begin{aligned} & \stackrel{(a)}{\leq} \frac{1}{n} \sum_{i=1}^n [(\nabla F(\bar{x}_k) - p_k^i)^T (\theta_k^i - \bar{x}_k) + p_k^{iT} (x^* - \bar{x}_k)] \\ &= \frac{1}{n} \sum_{i=1}^n [(\nabla F(\bar{x}_k) - p_k^i)^T (\theta_k^i - x^* + x^* - \bar{x}_k) \\ &\quad + p_k^{iT} (x^* - \bar{x}_k)] \\ &= \frac{1}{n} \sum_{i=1}^n [(\nabla F(\bar{x}_k) - p_k^i)^T (\theta_k^i - x^*) + \nabla^T F(\bar{x}_k)(x^* - \bar{x}_k)] \\ &\stackrel{(b)}{\leq} \frac{1}{n} \sum_{i=1}^n \|\nabla F(\bar{x}_k) - p_k^i\| \|\theta_k^i - x^*\| + F(x^*) - F(\bar{x}_k) \\ &\leq \frac{1}{n} \sum_{i=1}^n \|\nabla F(\bar{x}_k) - p_k^i\| D + F(x^*) - F(\bar{x}_k), \end{aligned} \quad (44)$$

where (a) follows from the optimality of  $\theta_i$  in (6); (b) is due to the property of convex function  $F$ , i.e.,  $F(x^*) - F(\bar{x}_k) \geq \nabla^T F(\bar{x}_k)(x^* - \bar{x}_k)$ . Substituting (44) into (43), we arrive at

$$\begin{aligned} F(\bar{x}_{k+1}) &\leq F(\bar{x}_k) + \eta_k \left( \frac{1}{n} \sum_{i=1}^n \|\nabla F(\bar{x}_k) - p_k^i\| D \right. \\ &\quad \left. + F(x^*) - F(\bar{x}_k) \right) + \frac{L\eta_k^2}{2} D^2 \\ &= (1 - \eta_k) F(\bar{x}_k) + \frac{\eta_k}{n} \sum_{i=1}^n \|\nabla F(\bar{x}_k) - p_k^i\| D \\ &\quad + \eta_k F(x^*) + \frac{L\eta_k^2}{2} D^2. \end{aligned} \quad (45)$$

Subtracting  $F(x^*)$  from both sides of inequality (45), we arrive at

$$\begin{aligned} & F(\bar{x}_{k+1}) - F(x^*) \\ &\leq (1 - \eta_k)(F(\bar{x}_k) - F(x^*)) + \frac{L\eta_k^2}{2} D^2 \\ &\quad + \frac{\eta_k}{n} \sum_{i=1}^n \|\nabla F(\bar{x}_k) - p_k^i\| D. \end{aligned} \quad (46)$$

$\square$

Taking the expectation and using the Jensen's inequality on the last term of (46), that is,  $\mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|] = \sqrt{(\mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|]^2)} \leq \sqrt{\mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|^2]}$ , we get

$$\begin{aligned} & \mathbb{E}[F(\bar{x}_{k+1})] - F(x^*) \\ &\leq (1 - \eta_k)(\mathbb{E}[F(\bar{x}_k)] - F(x^*)) + \frac{L\eta_k^2}{2} D^2 \\ &\quad + \frac{\eta_k}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|] D \\ &\leq (1 - \eta_k)(\mathbb{E}[F(\bar{x}_k)] - F(x^*)) + \frac{L\eta_k^2}{2} D^2 \\ &\quad + \frac{\eta_k}{n} \sum_{i=1}^n D \sqrt{\mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|^2]}. \end{aligned} \quad (47)$$

Substituting  $\gamma_k = \frac{2}{k+1}$ ,  $\eta_k = \frac{2}{k+2}$  and (13) into (47), we have

$$\begin{aligned} & \mathbb{E}[F(\bar{x}_{k+1})] - F(x^*) \\ &\leq (1 - \frac{2}{k+1})(\mathbb{E}[F(\bar{x}_k)] - F(x^*)) + \frac{2LD^2}{(k+2)^2} \end{aligned}$$

$$\begin{aligned}
& + \frac{2D\sqrt{12L^2C_1^2 + 3C_3 + 12C_2}}{(k+2)(k+2)^{\frac{1}{2}}} \\
& \leq (1 - \frac{2}{k+2})(\mathbb{E}[F(\bar{x}_k)] - F(x^*)) \\
& + \frac{2LD^2 + 2D\sqrt{12L^2C_1^2 + 3C_3 + 12C_2}}{(k+2)^{\frac{3}{2}}}. \quad (48)
\end{aligned}$$

By using Lemma 5 ( $r_1 = 1$ ,  $r_2 = \frac{3}{2}$ ,  $A = 2$  and  $B = 2LD^2 + 2D\sqrt{12L^2C_1^2 + 3C_3 + 12C_2}$ ), we obtain  $\mathbb{E}[F(\bar{x}_{k+1})] - F(x^*) \leq \frac{C_4}{(k+3)^{\frac{1}{2}}}$ , where  $C_4 = \max\{\sqrt{3}(F(\bar{x}_1) - F(x^*)), 2LD^2 + 2D\sqrt{12L^2C_1^2 + 3C_3 + 12C_2}\}$ .  $\square$

### G. Proof of Theorem 2

*Proof.* It follows from the definition of FW-gap (15) that

$$g_k = \max_{x \in \mathcal{X}} \langle \nabla F(\bar{x}_k), \bar{x}_k - x \rangle = \langle \nabla F(\bar{x}_k), \bar{x}_k - \hat{v}_k \rangle, \quad (49)$$

where

$$\hat{v}(k) \in \operatorname{argmin}_{v \in \mathcal{X}} \langle v, \nabla F(\bar{x}_k) \rangle. \quad (50)$$

According to the  $L$ -smooth property of  $F$ , we can write

$$\begin{aligned}
F(\bar{x}_{k+1}) & \leq F(\bar{x}_k) + \nabla^T F(\bar{x}_k)(\bar{x}_{k+1} - \bar{x}_k) + \frac{L}{2} \|\bar{x}_{k+1} - \bar{x}_k\|^2 \\
& \stackrel{(a)}{=} F(\bar{x}_k) + \frac{\eta_k}{n} \sum_{i=1}^n (\nabla F(\bar{x}_k) + p_k^i - p_k^i)^T (\theta_k^i - \bar{x}_k) \\
& + \frac{L\eta_k^2}{2} \|\bar{\theta}_k - \bar{x}_k\|^2 \\
& \stackrel{(b)}{\leq} F(\bar{x}_k) + \frac{\eta_k}{n} \sum_{i=1}^n [(p_k^i - \nabla F(\bar{x}_k) + \nabla F(\bar{x}_k))^T (\hat{v}_k - \bar{x}_k) \\
& + (\nabla F(\bar{x}_k) - p_k^i)^T (\theta_k^i - \bar{x}_k)] + \frac{L\eta_k^2}{2} \|\bar{\theta}_k - \bar{x}_k\|^2 \\
& \stackrel{(c)}{\leq} F(\bar{x}_k) + \frac{2\eta_k}{n} \sum_{i=1}^n \|\nabla F(\bar{x}_k) - p_k^i\| D \\
& - \eta_k g_k + \frac{LD^2}{2} \eta_k^2, \quad (51)
\end{aligned}$$

where (a) is because of Lemma 6 (b); (b) is due to the fact  $\theta_k^i = \operatorname{argmin}_{\phi \in \mathcal{X}} \langle p_k^i, \phi \rangle$  in (6) and  $\hat{v}(k)$  is defined in (50); (c) follows from (49) and the Assumption 3. Taking the expectation on both sides of (51), we arrive at

$$\begin{aligned}
\mathbb{E}[F(\bar{x}_{k+1})] & \leq \mathbb{E}[F(\bar{x}_k)] + \frac{2\eta_k}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|] D \\
& - \eta_k \mathbb{E}[g_k] + \frac{LD^2}{2} \eta_k^2 \\
& \stackrel{(a)}{\leq} \mathbb{E}[F(\bar{x}_k)] + \frac{2\eta_k}{n} \sum_{i=1}^n D \sqrt{\mathbb{E}[\|\nabla F(\bar{x}_k) - p_k^i\|^2]} \\
& - \eta_k \mathbb{E}[g_k] + \frac{LD^2}{2} \eta_k^2 \\
& \stackrel{(b)}{\leq} \mathbb{E}[F(\bar{x}_k)] + \frac{2\eta_k D \sqrt{12L^2C_1^2 + 3C_3 + 12C_2}}{(k+2)^{\frac{1}{2}}} \\
& - \eta_k \mathbb{E}[g_k] + \frac{LD^2}{2} \eta_k^2, \quad (52)
\end{aligned}$$

where (a) is due to the Jensen's inequality; (b) follows from (13). Summing from  $k = 1$  to  $k = K$  on both sides of (52), we get

$$\begin{aligned}
\mathbb{E} \left[ \sum_{k=1}^K \eta_k g_k \right] & \leq \mathbb{E} \left[ \sum_{k=1}^K (F(\bar{x}_k) - F(\bar{x}_{k+1})) \right] + \sum_{k=1}^K \left( \frac{LD^2}{2} \eta_k^2 \right. \\
& \left. + \frac{2\eta_k D \sqrt{12L^2C_1^2 + 3C_3 + 12C_2}}{(k+2)^{\frac{1}{2}}} \right) \\
& = \mathbb{E}[F(\bar{x}_1) - F(\bar{x}_{K+1})] + \sum_{k=1}^K \left( \frac{LD^2}{2} \eta_k^2 \right. \\
& \left. + \frac{2\eta_k D \sqrt{12L^2C_1^2 + 3C_3 + 12C_2}}{(k+2)^{\frac{1}{2}}} \right) \\
& \leq \mathbb{E}[F(\bar{x}_1) - F(x^*)] + \sum_{k=1}^K \left( \frac{LD^2}{2} \eta_k^2 \right. \\
& \left. + \frac{2\eta_k D \sqrt{12L^2C_1^2 + 3C_3 + 12C_2}}{(k+2)^{\frac{1}{2}}} \right), \quad (53)
\end{aligned}$$

where the last inequality arises from the optimality of  $x^*$ . According to the property of  $p$  series, we have  $m-1 \leq \sum_{k=1}^{2^m} \frac{2}{k+2}$ ,  $\sum_{k=1}^{2^m} (\frac{2}{k+2})^2 \leq 8$ , and there exists  $\beta \in \mathbb{R}$  such that  $\sum_{k=1}^{2^m} \frac{2}{(k+2)^{1.5}} \leq \beta$ . Define  $g_a := \min_{k \in [1, K]} g_k$ . Substituting  $K = 2^m$  and  $\eta_k = \frac{2}{k+2}$  into (53), we arrive at

$$\begin{aligned}
\mathbb{E}[(m-1)g_a] & \leq \mathbb{E} \left[ \sum_{k=1}^{2^m} \frac{2g_a}{k+2} \right] \\
& \leq \mathbb{E}[F(\bar{x}_1) - F(x^*)] + \frac{LD^2}{2} \sum_{k=1}^{2^m} \left( \frac{2}{k+2} \right)^2 \\
& + 2D \sqrt{12L^2C_1^2 + 3C_3 + 12C_2} \sum_{k=1}^{2^m} \frac{2}{(k+2)^{1.5}} \\
& \leq 4LD^2 + 2D\beta \sqrt{12L^2C_1^2 + 3C_3 + 12C_2} \\
& + \mathbb{E}[F(\bar{x}_1) - F(x^*)]. \quad (54)
\end{aligned}$$

Hence,

$$\begin{aligned}
\mathbb{E}[g_a] & \leq \frac{1}{m-1} \left( \mathbb{E}[F(\bar{x}_1) - F(x^*)] + 4LD^2 \right. \\
& \left. + 2D\beta \sqrt{12L^2C_1^2 + 3C_3 + 12C_2} \right) \\
& = \frac{1}{\log_2(K) - 1} \left( \mathbb{E}[F(\bar{x}_1) - F(x^*)] + 4LD^2 \right. \\
& \left. + 2D\beta \sqrt{12L^2C_1^2 + 3C_3 + 12C_2} \right), \quad (55)
\end{aligned}$$

where the equality is due to  $K = 2^m$ .  $\square$

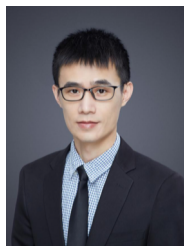
### REFERENCES

- [1] J. Chen, J. Sun, and G. Wang, "From unmanned systems to autonomous intelligent systems," *Eng.*, vol. 12, no. 5, pp. 16–19, 2022.
- [2] Z. Jiang, Q. Jia, and X. Guan, "On large action space in EV charging scheduling optimization," *Sci. China Inf. Sci.*, vol. 65, no. 122201, pp. <https://doi.org/10.1007/s11432-020-3106-7>, 2022.
- [3] R. Yang, L. Liu, and G. Feng, "An overview of recent advances in distributed coordination of multi-agent systems," *Unmanned Syst.*, vol. 10, no. 03, pp. 307–325, 2022.

- [4] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proc. COMPSTAT'2010*. Heidelberg: Physica-Verlag HD, 2010, pp. 177–186.
- [5] P. Sun, Z. Guo, G. Wang, J. Lan, and Y. Hu, "MARVEL: Enabling controller load balancing in software-defined networks with multi-agent reinforcement learning," *Comput. Netw.*, vol. 177, p. 107230, 2020.
- [6] G. Wang, S. Lu, G. B. Giannakis, G. Tesauro, and J. Sun, "Decentralized TD tracking with linear function approximation and its finite-time analysis," in *Adv. Neural Inf. Process. Syst.*, 2020.
- [7] D. Wang, Z. Wang, and Z. Wu, "Distributed convex optimization for nonlinear multi-agent systems disturbed by a second-order stationary process over a digraph," *Sci. China Inf. Sci.*, vol. 65, 2022.
- [8] A. Mokhtari, H. Hassani, and A. Karbasi, "Stochastic conditional gradient methods: From convex minimization to submodular maximization," *J. Mach. Learn. Res.*, vol. 21, no. 105, pp. 1–49, 2020.
- [9] Z. Akhtar and K. Rajawat, "Momentum based projection free stochastic optimization under affine constraints," in *American Control Conf.*, 2021, pp. 2619–2624.
- [10] S. J. Reddi, S. Sra, B. Póczos, and A. Smola, "Stochastic Frank-Wolfe methods for nonconvex optimization," in *Annual Allerton Conf. Commun., Control, Comput.*, 2016, pp. 1244–1251.
- [11] Z. Gong, Y. Xu, and D. Luo, "UAV cooperative air combat maneuvering confrontation based on multi-agent reinforcement learning," *Unmanned Syst.*, pp. 1–14, 2022.
- [12] S. Pu, A. Olshevsky, and I. C. Paschalidis, "Asymptotic network independence in distributed stochastic optimization for machine learning: examining distributed and centralized stochastic gradient descent," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 114–122, 2020.
- [13] Z. Wang, J. Zhang, T.-H. Chang, J. Li, and Z.-Q. Luo, "Distributed stochastic consensus optimization with momentum for nonconvex nonsmooth problems," *IEEE Trans. Signal Process.*, vol. 69, pp. 4486–4501, 2021.
- [14] G. Wang, G. B. Giannakis, and J. Chen, "Learning ReLU networks on linearly separable data: Algorithm, optimality, and generalization," *IEEE Trans. Signal Process.*, vol. 67, no. 9, pp. 2357–2370, 2019.
- [15] Y. Dai and Y. Weng, "Synchronous parallel block coordinate descent method for nonsmooth convex function minimization," *J. Syst. Sci. Complex.*, vol. 33, no. 2, pp. 345–365, 2020.
- [16] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, "Robust stochastic approximation approach to stochastic programming," *SIAM J. Optim.*, vol. 19, no. 4, pp. 1574–1609, 2009.
- [17] Y. Zhang, Y. Lou, Y. Hong, and L. Xie, "Distributed projection-based algorithms for source localization in wireless sensor networks," *IEEE Trans. Wireless Commun.*, vol. 14, no. 6, pp. 3131–3142, 2015.
- [18] X.-F. Wang, Y. Hong, X.-M. Sun, and K.-Z. Liu, "Distributed optimization for resource allocation problems under large delays," *IEEE Trans. Ind. Electron.*, vol. 66, no. 12, pp. 9448–9457, 2019.
- [19] S. Ghadimi and G. Lan, "Stochastic first- and zeroth-order methods for nonconvex stochastic programming," *SIAM J. Optim.*, vol. 23, no. 4, pp. 2341–2368, 2013.
- [20] A. Agrawal, S. Barratt, and S. Boyd, "Learning convex optimization models," *IEEE/CAA J. Autom. Sinica*, vol. 8, no. 8, pp. 1355–1364, 2021.
- [21] S. Fujishige and S. Isotani, "A submodular function minimization algorithm based on the minimum-norm base," *Pacific J. Opt.*, vol. 7, no. 1, pp. 3–17, 2011.
- [22] M. Frank and P. Wolfe, "An algorithm for quadratic programming," *Nav. Res. Logist.*, vol. 3, no. 1–2, pp. 95–110, 1956.
- [23] G. Lan and Y. Zhou, "Conditional gradient sliding for convex optimization," *SIAM J. Optim.*, vol. 26, no. 2, 2016.
- [24] D. S. Kalhan, A. Singh Bedi, A. Koppel, K. Rajawat, H. Hassani, A. K. Gupta, and A. Banerjee, "Dynamic online learning via Frank-Wolfe algorithm," *IEEE Trans. Signal Process.*, vol. 69, pp. 932–947, 2021.
- [25] T. Kerdreux, A. d'Aspremont, and S. Pokutta, "Projection-free optimization on uniformly convex sets," in *Proc. Artif. Intell. Statist.*, vol. 130, 2021, pp. 19–27.
- [26] B. Li, M. Coutiño, G. B. Giannakis, and G. Leus, "A momentum-guided Frank-Wolfe algorithm," *IEEE Trans. Signal Process.*, vol. 69, pp. 3597–3611, 2021.
- [27] Y. Zhang, B. Li, and G. B. Giannakis, "Accelerating Frank-Wolfe with weighted average gradients," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2021, pp. 5529–5533.
- [28] E. Hazan and S. Kale, "Projection-free online learning," in *Proc. 29th Int. Conf. Mach. Learn.*, 2012, pp. 521–528.
- [29] E. Hazan and H. Luo, "Variance-reduced and projection-free stochastic optimization," in *Proc. Int. Conf. Mach. Learn.*, 2016, pp. 1263–1271.
- [30] K. Bekiroglu, S. Srinivasan, E. Png, R. Su, and C. Lagoa, "Recursive approximation of complex behaviours with IoT-data imperfections," *IEEE/CAA J. Autom. Sinica*, vol. 7, no. 3, pp. 656–667, 2020.
- [31] M. Weber, "Projection-free nonconvex stochastic optimization on Riemannian manifolds," *IMA J. Numer. Anal.*, 2021.
- [32] A. Bellet, Y. Liang, A. B. Garakani, M. Balcan, and F. Sha, "A distributed Frank-Wolfe algorithm for communication-efficient sparse learning," in *Proc. SIAM Int. Conf. Data Mining*, 2015.
- [33] H.-T. Wai, A. Scaglione, J. Lafond, and E. Moulines, "A projection-free decentralized algorithm for non-convex optimization," in *IEEE Global Conf. Signal and Inf. Process.*, 2016, pp. 475–479.
- [34] H.-T. Wai, J. Lafond, A. Scaglione, and E. Moulines, "Decentralized Frank-Wolfe algorithm for convex and nonconvex problems," *IEEE Trans. Autom. Control*, vol. 62, no. 11, pp. 5522–5537, 2017.
- [35] G. Chen, P. Yi, and Y. Hong, "Distributed optimization with projection-free dynamics," *arXiv:2105.02450*, 2021.
- [36] L. Zhang, G. Wang, D. Romero, and G. B. Giannakis, "Randomized block Frank-Wolfe for convergent large-scale learning," *IEEE Trans. on Signal Processing*, vol. 65, no. 24, pp. 6448–6461, 2017.
- [37] J. Xie, C. Zhang, Z. Shen, C. Mi, and H. Qian, "Decentralized gradient tracking for continuous DR-submodular maximization," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2019.
- [38] H. Gao, H. Xu, and S. Vucetic, "Sample efficient decentralized stochastic Frank-Wolfe methods for continuous DR-Submodular maximization," in *Proc. Int. Joint Conf. Artif. Intell.*, 2021.
- [39] A. Cutkosky and F. Orabona, "Momentum-based variance reduction in non-convex SGD," in *Proc. Adv. Neural Inf. Process. Syst.*, 2019, pp. 15 210–15 219.
- [40] Y. E. Nesterov, "A method for solving the convex programming problem with convergence rate  $\mathcal{O}(1/k^2)$ ," *Dokl. Akad. Nauk SSSR*, vol. 269, no. 3, pp. 543–547, 1983.
- [41] J. N. Tsitsiklis, "Problems in decentralized decision making and computation," *Ph.D. dissertation, Dept. Elect. Eng. Comput. Sci., MIT, Boston, MA, USA*, 1984.
- [42] Y. Zhang, Y. Lou, and Y. Hong, "An approximate gradient algorithm for constrained distributed convex optimization," *IEEE/CAA J. Autom. Sinica*, vol. 1, no. 1, pp. 61–67, 2014.
- [43] X. Ren, D. Li, Y. Xi, and H. Shao, "Distributed subgradient algorithm for multi-agent optimization with dynamic stepsize," *IEEE/CAA J. Autom. Sinica*, vol. 8, no. 8, pp. 1451–1464, 2021.
- [44] J. Gao, X.-W. Liu, Y.-H. Dai, Y. Huang, and P. Yang, "Achieving geometric convergence for distributed optimization with barzilai-borwein step sizes," *Sci. China Inf. Sci.*, vol. 65, no. 149204, pp. <https://doi.org/10.1007/s11432-020-3256-x>, 2022.
- [45] P. D. Lorenzo and G. Scutari, "Next: In-network nonconvex optimization," *IEEE Trans. Signal Inf. Process. Netw.*, vol. 2, no. 2, pp. 120–136, 2016.
- [46] Z. Li, B. Liu, and Z. Ding, "Consensus-based cooperative algorithms for training over distributed data sets using stochastic gradients," *IEEE Trans. Neural Netw. Learn. Syst.*, pp. 1–11, 2021.
- [47] X. Ma, P. Yi, and J. Chen, "Distributed gradient tracking methods with finite data rates," *J. Syst. Sci. Complex.*, vol. 34, no. 5, pp. 1927–1952, 2021.



**Jie Hou** received the bachelor's degree in computer science and technology from School of Computer and Information Engineering, Luoyang Institute of Science and Technology, China, in 2017, the master's degree in computer science and technology from School of Computer Science and Software Engineering, Tiangong University, China, in 2020. She is currently pursuing the Ph.D. degree in control science and engineering with the School of Automation, Beijing Institute of Technology, Beijing, China. Her current research interests include distributed optimization and stochastic optimization.



**Xianlin Zeng** (S'12-M'15) received the B.S. and M.S. degrees in Control Science and Engineering from the Harbin Institute of Technology, Harbin, China, in 2009 and 2011, respectively, and the Ph.D. degree in Mechanical Engineering from the Texas Tech University in 2015. He is currently an associate professor in the Key Laboratory of Intelligent Control and Decision of Complex Systems, School of Automation, Beijing Institute of Technology, Beijing, China. His current research interests include distributed optimization, distributed control,

and distributed computation of network systems.



**Jie Chen** (M'09-SM'12-F'19) received his B.Sc., M.Sc., and the Ph.D. degrees in control theory and control engineering from the Beijing Institute of Technology, Beijing, China, in 1986, 1996, and 2001, respectively. From 1989 to 1990, he was a visiting scholar at the California State University, Long Beach, California, USA. From 1996 to 1997, he was a research fellow with the School of Engineering at the University of Birmingham, Birmingham, UK. He is a Professor with the School of Automation, Beijing Institute of Technology, where he serves as

the Director of the Key Laboratory of Intelligent Control and Decision of Complex Systems. He also serves as the President of Tongji University, Shanghai, China. His research interests include multiagent systems, multiobjective optimization and decision, and constrained nonlinear control.

Prof. Chen is currently the Editor-in-Chief of Unmanned Systems, Autonomous Intelligent Systems, and Journal of Systems Science and Complexity. He has served on the editorial boards of several journals, including the IEEE Transactions on Cybernetics, International Journal of Robust and Nonlinear Control, and Science China Information Sciences. He is a Fellow of IEEE, IFAC, and a member of the Chinese Academy of Engineering.



**Gang Wang** (M'18) received a B.Eng. degree in Automatic Control in 2011, and a Ph.D. degree in Control Science and Engineering in 2018, both from the Beijing Institute of Technology, Beijing, China. He also received a Ph.D. degree in Electrical and Computer Engineering from the University of Minnesota, Minneapolis, USA, in 2018, where he stayed as a postdoctoral researcher until July 2020. Since August 2020, he has been a professor with the School of Automation at the Beijing Institute of Technology.

His research interests focus on the areas of signal processing, control, and reinforcement learning with applications to cyber-physical systems and multi-agent systems. He was the recipient of the Best Paper Award from the Frontiers of Information Technology & Electronic Engineering (FITEE) in 2021, the Excellent Doctoral Dissertation Award from the Chinese Association of Automation in 2019, the Best Student Paper Award from the 2017 European Signal Processing Conference, and the Best Conference Paper at the 2019 IEEE Power & Energy Society General Meeting. He is currently on the editorial board of Signal Processing, IEEE Open Journal of Control Systems, and IEEE Transactions on Signal and Information Processing over Networks.



**Jian Sun** received his B.Sc. degree from the Department of Automation and Electric Engineering, Jilin Institute of Technology, Changchun, China, in 2001, the M.Sc. degree from the Changchun Institute of Optics, Fine Mechanics and Physics, Chinese Academy of Sciences (CAS), Changchun, China, in 2004, and the Ph.D. degree from the Institute of Automation, CAS, Beijing, China, in 2007.

He was a Research Fellow with the Faculty of Advanced Technology, University of Glamorgan, Pontypridd, U.K., from 2008 to 2009. He was a Post-

Doctoral Research Fellow with the Beijing Institute of Technology, Beijing, from 2007 to 2010. In 2010, he joined the School of Automation, Beijing Institute of Technology, where he has been a Professor since 2013. His current research interests include networked control systems, time-delay systems, and security of cyber-physical systems.

Dr. Sun is an editorial board member of the IEEE Transactions on Systems, Man and Cybernetics: Systems, the Journal of Systems Science & Complexity, and Acta. Automatica Sinica.