# Distributed Stochastic Consensus Optimization with Momentum for Nonconvex Nonsmooth Problems

Zhiguo Wang, Jiawei Zhang, Tsung-Hui Chang, Jian Li and Zhi-Quan Luo

*Abstract*—While many distributed optimization algorithms have been proposed for solving smooth or convex problems over the networks, few of them can handle non-convex and non-smooth problems. Based on a proximal primal-dual approach, this paper presents a new (stochastic) distributed algorithm with Nesterov momentum for accelerated optimization of non-convex and non-smooth problems. Theoretically, we show that the proposed algorithm can achieve an $\epsilon$-stationary solution under a constant step size with $\mathcal{O}(1/\epsilon^2)$ computation complexity and $\mathcal{O}(1/\epsilon)$ communication complexity. When compared to the existing gradient tracking based methods, the proposed algorithm has the same order of computation complexity but lower order of communication complexity. To the best of our knowledge, the presented result is the first stochastic algorithm with the $\mathcal{O}(1/\epsilon)$ communication complexity for non-convex and non-smooth problems. Numerical experiments for a distributed non-convex regression problem and a deep neural network based classification problem are presented to illustrate the effectiveness of the proposed algorithms.

*Index Terms*—Distributed optimization, stochastic optimization, momentum, non-convex and non-smooth optimization.

## I. INTRODUCTION

Recently, motivated by large-scale machine learning [1] and mobile edge computing [2], many signal processing applications involve handling very large datasets [3] that are processed over networks with distributed memories and processors. Such signal processing and machine learning problems are usually formulated as a multi-agent distributed optimization problem [4]. In particular, many of the applications can be formulated as the following finite sum problem

$$\min_x \ \sum_{i=1}^N \Big( f_i(x) + r_i(x) \Big), \qquad (1)$$

where $N$ is the number of agents, $x \in \mathbb{R}^n$ contains the model parameters to be learned, $f_i(x) : \mathbb{R}^n \to \mathbb{R}$ is a closed and smooth (possibly nonconvex) loss function, and $r_i(x)$ is a convex and possibly non-smooth regularization term. Depending on how the data are acquired, there are two scenarios for problem (1) [5].

Zhiguo Wang is with College of Mathematics, Sichuan University, Chengdu, Sichuan 610064, China (e-mail: wangzhiguo@scu.edu.cn). This work was done in the Chinese University of Hong Kong, Shenzhen.

Jiawei Zhang, Tsung-Hui Chang, and Zhi-Quan (Tom) Luo are with the Chinese University of Hong Kong, Shenzhen 518172, China and also with Shenzhen Research Institute of Big Data , Shenzhen, Guangdong Province 518172, China (e-mail: jiaweizhang2@link.cuhk.edu.cn; tsunghui.chang@ieee.org; luozq@cuhk.edu.cn). Corresponding author: Zhi-Quan (Tom) Luo.

Jian Li is with Department of Electrical and Computer Engineering, University of Florida (e-mail: li@dsp.ufl.edu).

- Offline/Batch learning: the agents are assumed to have the complete local dataset. Specifically, the local cost functions can be written as

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_i^j(x), \ i = 1, \ldots, N, \qquad (2)$$

where $f_i^j(x)$ is the cost for the $j$-th data sample at the $i$-th agent, and $m$ is the total number of local samples. When $m$ is not large, each agent $i$ may compute the full gradient of $f_i(x)$ for deterministic parameter optimization.

- Online/Streaming learning: when the data samples follow certain statistical distribution and are acquired by the agents in an online/streaming fashion, one can define $f_i(x)$ as the following expected cost

$$f_i(x) = \mathbb{E}_{\xi \sim \mathcal{B}_i}[f_i(x, \xi)], i = 1, \ldots, N, \qquad (3)$$

where $\mathcal{B}_i$ denotes the data distribution at agent $i$, and $f_i(x, \xi)$ is the cost function of a random data sample $\xi$. Under the online setting, only a stochastic estimate $G_i(x, \xi)$ for $\nabla f_i(x)$ can be obtained by the agent and stochastic optimization methods can be used. Note that if the agent is not able to compute the full gradient in the batch setting, a stochastic gradient estimate by mini-batch data samples can be obtained and the problem is solved in a similar fashion by stochastic optimization.

These two settings for local cost functions are popularly used in many machine learning models including deep learning and empirical risk minimization problems [5]. For both scenarios, many distributed optimization methods have been developed for solving problems (1).

Specifically, for batch learning and under convex or strongly convex assumptions, algorithms such as the distributed sub-gradient method [6], EXTRA [7], PG-EXTRA [8] and primal-dual based methods including the alternating direction method of multipliers (ADMM) [1], [4], [9] and the UDA in [10] are proposed. For non-convex problems, the authors in [11] studied the convergence of proximal decentralized gradient descent (DGD) method with a diminishing step size. Based on the successive convex approximation (SCA) technique and the gradient tracking (GT) method, the authors in [12] proposed a network successive convex approximation (NEXT) algorithm for (1), and it is extended to more general scenarios with time varying networks and stronger convergence analysis results [13], [14]. In [15], based on an inexact augmented Lagrange method, a proximal primal-dual algorithm (Prox-PDA) is developed for (1) with smooth and non-convex $f_i(x)$ and without $r_i(x)$. A near-optimal algorithm xFilter

is further proposed in [16] that can achieve the computation complexity lower bound of first-order distributed optimization algorithms. To handle non-convex and non-smooth problems with polyhedral constraints, the authors of [17], [18] proposed a proximal augmented Lagrangian (AL) method for solving (1) by introducing a proximal variable and an exponential averaging scheme.

For streaming learning, the stochastic proximal gradient consensus method based on ADMM is proposed in [19] to solve (1) with convex objective functions. For non-convex problems, the decentralized parallel stochastic gradient descent (D-PSGD) [20] is applied to (1) (without $r_i(x)$) for training large-scale neural networks, and the convergence rate is analyzed. The analysis of D-PSGD relies on an assumption that $\frac{1}{N} \sum_{i=1}^{N} ||\nabla f_i(x) - \nabla f(x)||^2$ is bounded, which implies that the variance of data distributions across the agents should be controlled. In [21], the authors proposed an improved D-PSGD algorithm, called $D^2$, which removes such assumption and is less sensitive to the data variance across agents. However, $D^2$ requires a restrictive assumption on the eigenvalue of the mixing matrix. This assumption is relaxed by the GNSD algorithm in [22], which essentially is a stochastic counterpart of the GT algorithm in [14]. We should emphasize here that the algorithms in [20], [21], [22] can only handle smooth problems without constraints and regularization terms. The work [23] proposed a multi-agent projected stochastic gradient decent (PSGD) algorithm for (1) but $r_i(x)$ is limited to the indicator function of compact convex sets. Besides, there is no convergence rate analysis in [23].

In this paper, we develop a new distributed stochastic optimization algorithm for the non-convex and non-smooth problem (1). The proposed algorithm is inspired by the proximal AL framework in [17] and has three new features. First, the proposed algorithm is a stochastic distributed algorithm that can be used either for streaming/online learning or batch/offline learning with mini-batch stochastic gradients. Second, the proposed algorithm can handle problem (1) with non-smooth terms that have a polyhedral epigraph, which is more general than [17], [18]. Third, the proposed algorithm incorporates the Nesterov momentum technique for fast convergence. The Nesterov momentum technique has been applied for accelerating the convergence of distributed optimization. For example, in [24], [25], the distributed gradient descent methods with the Nesterov momentum are proposed, and are shown to achieve the optimal iteration complexity for convex problems. In practice, since SGD with momentum often can converge faster, it is also commonly used to train deep neural networks [26], [27]. We note that [24], [25], [26], [27] are for smooth problems. To the best of our knowledge, the Nesterov momentum technique has not been used for distributed non-convex and non-smooth optimization.

Our contributions are summarized as follows.

- We propose a new stochastic proximal primal dual algorithm with momentum (SPPDM) for non-convex and non-smooth problem (1) under the online/streaming setting. For the offline/batch setting where the full gradients of the local cost functions are available, SPPDM reduces to a deterministic algorithm, named the PPDM algorithm.

- We show that the proposed SPPDM and PPDM can achieve an $\epsilon$-stationary solution of (1) under a constant step size with computation complexities of $\mathcal{O}(1/\epsilon^2)$ and $\mathcal{O}(1/\epsilon)$, respectively, while both have a communication complexity of $\mathcal{O}(1/\epsilon)$. The convergence analysis neither requires assumption on the boundedness of $\frac{1}{N} \sum_{i=1}^{N} ||\nabla f_i(x) - \nabla f(x)||^2$ nor on the eigenvalues of the mixing matrix.

- As shown in Table I, the proposed SPPDM/PPDM algorithms have the same order of computation complexity as the existing methods and lower order of communication complexity when compared with the existing GT based methods.

- Numerical experiments for a distributed non-convex regression problem and a deep neural network (DNN) based classification problem show that the proposed algorithms outperforms the existing methods.

**Notation:** We denote $\mathbf{I}_n$ as the $n$ by $n$ identity matrix and $\mathbf{1}$ as the all-one vector, i.e., $\mathbf{1} = [1, \ldots, 1]^\top$. $\langle \mathbf{a}, \boldsymbol{b} \rangle$ represents the inner product of vectors $\mathbf{a}$ and $\boldsymbol{b}$, $||\mathbf{a}||$ is the Euclidean norm of vector $\mathbf{a}$ and $||\mathbf{a}||_1$ is the $\ell_1$-norm of vector $\mathbf{a}$; $\otimes$ denotes the Kronecker product. For a matrix $\mathbf{A}$, $\sigma_A > 0$ denotes its largest singular value. $\text{diag}\{a_1, \ldots, a_N\}$ denotes a diagonal matrix with $a_1, \ldots, a_N$ being the diagonal entries while $\text{diag}\{\mathbf{A}_1, \ldots, \mathbf{A}_N\}$ denotes a block diagonal matrices with each $\mathbf{A}_i$ being the $i$th block diagonal matrix. $[\mathbf{A}]_{ij}$ represents the element of $\mathbf{A}$ in the $i$th row and $j$th column.

For problem (1), we denote $\mathbf{x} = [x_1^\top, \ldots, x_N^\top]^\top \in \mathbb{R}^{Nn}$, $f(\mathbf{x}) = \sum_{i=1}^{N} f_i(x_i)$, and $r(\mathbf{x}) = \sum_{i=1}^{N} r_i(x_i)$. The gradient of $f(\cdot)$ at $\mathbf{x}$ is denoted by

$$\nabla f(\mathbf{x}) = [(\nabla f_1(x_1))^\top, \ldots, (\nabla f_N(x_N))^\top]^\top,$$

where $\nabla f_i(x_i)$ is the gradient of $f_i$ at $x_i$. In the online/streaming setting, we denote the stochastic gradient estimates of agents as

$$G(\mathbf{x}, \boldsymbol{\xi}) = [(G_i(x_1, \xi_1))^\top, \ldots, (G_N(x_N, \xi_N))^\top]^\top,$$

where $\boldsymbol{\xi} = [\xi_1^\top, \ldots, \xi_N^\top]$. Lastly, we define the following proximal operator of $r_i$

$$\text{prox}_{r_i}^\alpha(x) = \arg\min_u \frac{\alpha}{2} ||x - u||^2 + r_i(u),$$

where $\alpha$ is a parameter.

**Synopsis:** In Section II, the proposed SPPDM and PPDM algorithms are presented and their connections with existing methods are discussed. Based on an inexact stochastic primal-dual framework, it is shown how the SPPDM and PPDM algorithms are devised. Section III presents the theoretical results of the convergence conditions and convergence rate of the SPPDM and PPDM algorithms. The performance of the SPPDM and PPDM algorithms are illustrated in Section IV. Lastly, the conclusion is given in Section V.

## II. ALGORITHM DEVELOPMENT

### A. Network Model and Consensus Formulation

Let us denote the multi-agent network as a graph $\mathcal{G}$, which contains a node set $V := \{1, \ldots, N\}$ and an edge set $\mathcal{E}$

TABLE I
COMPARISONS OF DIFFERENT ALGORITHMS

| Algorithm | objective function | gradient | stepsize | momentum | computational | communication |
|---|---|---|---|---|---|---|
| D-PSGD [20] | $f(\mathbf{x})$ | stochastic | decreasing | ✗ | $\mathcal{O}(\frac{N}{\epsilon^2})$ | $\mathcal{O}(\frac{1}{\epsilon^2})$ |
| D$^2$ [21] | $f(\mathbf{x})$ | stochastic | decreasing | ✗ | $\mathcal{O}(\frac{N}{\epsilon^2})$ | $\mathcal{O}(\frac{1}{\epsilon^2})$ |
| GNSD [22] | $f(\mathbf{x})$ | stochastic | decreasing | ✗ | $\mathcal{O}(\frac{N}{\epsilon^2})$ | $\mathcal{O}(\frac{1}{\epsilon^2})$ |
| PR-SGD-M [27] | $f(\mathbf{x})$ | stochastic | decreasing | ✓ | $\mathcal{O}(\frac{N}{\epsilon^2})$ | $\mathcal{O}(\frac{1}{\epsilon^2})$ |
| PSGD [23] | $f(\mathbf{x}) + r(\mathbf{x})$ | stochastic | decreasing | ✗ | ✗ | ✗ |
| STOC-ADMM [28] | $f(\mathbf{x}) + r(\mathbf{x})$ | stochastic | fixed | ✗ | $\mathcal{O}(\frac{N}{\epsilon^2})$ | $\mathcal{O}(\frac{1}{\epsilon})$ |
| Prox-PDA [15] | $f(\mathbf{x})$ | full | fixed | ✗ | $\mathcal{O}(\frac{mN}{\epsilon})$ | $\mathcal{O}(\frac{1}{\epsilon})$ |
| Prox-DGD [11] | $f(\mathbf{x}) + r(\mathbf{x})$ | full | decreasing | ✗ | ✗ | ✗ |
| Prox-ADMM [17] | $f(\mathbf{x}) + r(\mathbf{x})$ | full | fixed | ✗ | $\mathcal{O}(\frac{mN}{\epsilon})$ | $\mathcal{O}(\frac{1}{\epsilon})$ |
| Proposed | $f(\mathbf{x}) + r(\mathbf{x})$ | full | fixed | ✓ | $\mathcal{O}(\frac{mN}{\epsilon})$ | $\mathcal{O}(\frac{1}{\epsilon})$ |
| | | stochastic | fixed | ✓ | $\mathcal{O}(\frac{N}{\epsilon^2})$ | $\mathcal{O}(\frac{1}{\epsilon})$ |

with cardinality $|\mathcal{E}|$. For each agent $i$, it has neighboring agents in the subset $\mathcal{N}_i := \{j \in V | (i,j) \in \mathcal{E}\}$ with size $d_i \geq 1$. It is assumed that each agent $i$ can communicate with its neighborhood $\mathcal{N}_i$. We also assume that the graph $\mathcal{G}$ is undirected and is connected in the sense that for any of two agents in the network there is a path connecting them through the edge links. Thus, problem (1) can be equivalently written as

$$\min_{\substack{x_i \\ i=1,\ldots,N}} \sum_{i=1}^{N} \Big( f_i(x_i) + r_i(x_i) \Big) \tag{4a}$$

$$\text{s.t. } x_i = x_j, \ \forall (i,j) \in \mathcal{E}. \tag{4b}$$

Let us introduce the incidence matrix $\tilde{\mathbf{A}} \in \mathbb{R}^{|\mathcal{E}| \times n}$ which has $\tilde{\mathbf{A}}(\ell, i) = 1$ and $\tilde{\mathbf{A}}(\ell, j) = -1$ if $(i,j) \in \mathcal{E}$ with $j > i$, and zero otherwise, for $\ell = 1, \ldots, |\mathcal{E}|$. Define the extended incidence matrix as $\mathbf{A} := \tilde{\mathbf{A}} \otimes \mathbf{I}_n$. Then (4) is equivalent to

$$\min_{\mathbf{x}} f(\mathbf{x}) + r(\mathbf{x}) \tag{5a}$$

$$\text{s.t. } \mathbf{A}\mathbf{x} = \mathbf{0}. \tag{5b}$$

### B. Proposed SPPDM and PPDM Algorithm

In this section, we present the proposed SPPDM algorithm for solving (5) under the online/streaming setting in (3). The algorithm steps are outlined in Algorithm 1. Before showing how the algorithm is developed in Section II-C, let us make a few comments about SPPDM.

In Algorithm 1, $\alpha, \beta, \gamma, c, \kappa, \eta_k$ are some positive constant parameters that depend on the problem instance (such as the Lipschitz constants of $\{\nabla f_i\}$ and the graph Laplacian matrix). Equations (7)-(10) are the updates performed by each agent $i$ within the $k$th communication round, for $k = 1, 2, \ldots$, and $i = 1, \ldots, N$. Specifically, step (7) is the introduced Nesterov momentum term $s_i^k$ for accelerating the algorithm convergence, where $\eta_k$ is the extrapolation coefficient at iteration $k$. Step (8) shows how the neighboring variables $\{x_j\}_{j \in \mathcal{N}_i}$ are used for local gradient update. Note here that in SPPDM the agent uses the sample average $\frac{1}{|\mathcal{I}|} \sum_{j=1}^{|\mathcal{I}|} G_i(s_i^k, \xi_{ij}^k)$ to approximate $\nabla f_i(s_i^k)$, where $\xi_{ij}^k \sim \mathcal{B}_i$, $j = 1, \ldots, |\mathcal{I}|$, denotes the samples drawn by agent $i$ in the $k$th iteration. Besides,

---

**Algorithm 1** Proposed SPPDM Algorithm

**Given** parameters $\alpha, \beta, \gamma, c, \kappa, \eta_k$ and initial values of $x_i^0$, $i = 1, \ldots, N$. Let

$$\psi_i = \gamma + 2cd_i + \kappa \tag{6}$$

and set $s_i^0 = x_i^0$, $i = 1, \ldots, N$. Do

$$x_i^{\frac{1}{2}} = (\gamma + cd_i + \kappa)\frac{x_i^0}{\psi_i} + \frac{c}{\psi_i}\sum_{j \in \mathcal{N}_i} x_j^0 - \frac{1}{\psi_i}\nabla f_i(x_i^0),$$

$$x_i^1 = \text{prox}_{r_i}^{\alpha_i}\left(x_i^{\frac{1}{2}}\right), \ i = 1, \ldots, N.$$

**for** communication round $k = 1, 2, \ldots$ **do**
  **for** agent $i = 1, 2, \ldots, N$ (in parallel) **do**

$$s_i^k = x_i^k + \eta_k(x_i^k - x_i^{k-1}), \tag{7}$$

$$x_i^{k+\frac{1}{2}} = x_i^{k-1+\frac{1}{2}} + \frac{d_i}{\psi_i}((c-\alpha)x_i^k - cx_i^{k-1}), \tag{8}$$

$$+ \frac{1}{\psi_i}\sum_{j \in \mathcal{N}_i}((c+\alpha)x_j^k - cx_j^{k-1})$$

$$+ \frac{1}{\psi_i}\Big(\gamma(s_i^k - s_i^{k-1}) + \kappa(z_i^k - z_i^{k-1})\Big)$$

$$- \frac{1}{\psi_i|\mathcal{I}|}\sum_{j=1}^{|\mathcal{I}|}(G_i(s_i^k, \xi_{ij}^k) - G_i(s_i^{k-1}, \xi_{ij}^{k-1})),$$

$$x_i^{k+1} = \text{prox}_{r_i}^{\psi_i}\left(x_i^{k+\frac{1}{2}}\right), \tag{9}$$

$$z_i^{k+1} = z_i^k + \beta(x_i^{k+1} - z_i^k). \tag{10}$$

  **end for**
**end for**

---

in (8), both approximate gradients at $s_i^k$ and $s_i^{k-1}$ are used. Step (9) performs the proximal gradient update with respect to the regularization term $r_i(x)$. In step (8), the variable $\{z_i^k\}$ is a "proximal" variable introduced for overcoming the non-convexity of $f_i$ (see (19)), and is updated as in step (10).

By stacking the variables for all $i = 1, \ldots, N$, one can write (7)-(10) in a vector form. Specifically, step (8) for $i =$

$1, \ldots, N$, can be expressed compactly as

$$\mathbf{x}^{k+\frac{1}{2}} = \mathbf{x}^{k-1+\frac{1}{2}} + \mathbf{U}\mathbf{x}^k - \tilde{\mathbf{U}}\mathbf{x}^{k-1}$$
$$+ \gamma\boldsymbol{\Psi}^{-1}(\mathbf{s}^k - \mathbf{s}^{k-1}) + \kappa\boldsymbol{\Psi}^{-1}(\mathbf{z}^k - \mathbf{z}^{k-1})$$
$$- \boldsymbol{\Psi}^{-1}(\bar{G}(\mathbf{s}^k, \boldsymbol{\xi}^k) - \bar{G}(\mathbf{s}^{k-1}, \boldsymbol{\xi}^{k-1})), \tag{11}$$

where $\mathbf{U}$ and $\tilde{\mathbf{U}}$ are two matrices satisfying

$$[\mathbf{U}]_{ij} = \begin{cases} \frac{d_i}{\psi_i}(c - \alpha), & i = j, \\ \frac{c+\alpha}{\psi_i}, & i \neq j \text{ and } (i,j) \in \mathcal{E}, \\ 0, & \text{otherwise.} \end{cases} \tag{12}$$

$$[\tilde{\mathbf{U}}]_{ij} = \begin{cases} \frac{d_i c}{\psi_i}, & i = j, \\ \frac{c}{\psi_i}, & i \neq j \text{ and } (i,j) \in \mathcal{E}, \\ 0 & \text{otherwise.} \end{cases} \tag{13}$$

for all $i, j = 1, \ldots, N$, $\boldsymbol{\Psi}$ is a diagonal matrix with its $i$th element being $\psi_i := \gamma + 2cd_i + \kappa$ for $i = 1, \ldots, N$, and

$$\bar{G}(\mathbf{s}^k, \boldsymbol{\xi}^k) := \frac{1}{|\mathcal{I}|} \sum_{j=1}^{|\mathcal{I}|} G(\mathbf{s}^k, \boldsymbol{\xi}_j^k). \tag{14}$$

When the full gradients $\nabla f_i$ are available under the offline/batch setting, the approximate gradient $G_i$ in (8) and (11) can be replaced by $\nabla f_i$. Then, the SPPDM algorithm reduces to the PPDM algorithm.

**Remark 1.** We show that the PPDM algorithm can have a close connection with the PG-EXTRA algorithm in [8]. Specifically, let us set $\eta_k = 0$ (no Nesterov momentum) and $\beta = 1$ (no proximal variable). Then, we have $s_i^k = z_i^k = x_i^k$ for all $k, i$, and the momentum and proximal variable update in (7) and (10) can be removed. As a result, (11) reduces to

$$\mathbf{x}^{k+\frac{1}{2}} = \mathbf{x}^{k-1+\frac{1}{2}} + \mathbf{W}\mathbf{x}^k - \tilde{\mathbf{W}}\mathbf{x}^{k-1}$$
$$- \boldsymbol{\Psi}^{-1}(\nabla f(\mathbf{x}^k) - \nabla f(\mathbf{x}^{k-1})), \tag{15}$$

where $\mathbf{W} = \mathbf{U} + (\gamma + \kappa)\boldsymbol{\Psi}^{-1}$ and $\tilde{\mathbf{W}} = \tilde{\mathbf{U}} + (\gamma + \kappa)\boldsymbol{\Psi}^{-1}$. One can see that (15) and (9) have an identical form as the PG-EXTRA algorithm in [8, Eqn. (3a)-(3b)]. Therefore, the proposed PPDM algorithm can be regarded as an accelerated version of the PG-EXTRA with extra capability to handle non-convex problems. One should note that, unlike (12) and (13), the PG-EXTRA allows a more flexible choice of the mixing matrix $\mathbf{W}$, and thus it is also closely related to the GT based methods [5].

**Remark 2.** The PPDM algorithm also has a close connection with the distributed Nesterov gradient (D-NG) algorithm in [24]. Specifically, let us set $\alpha = c$ and $\beta = 1$ (no proximal variable) and remove the non-smooth regularization term $r(\mathbf{x})$. Then, we have $z_i^k = x_i^k$ for all $k, i$, and the proximal gradient update (9) and the proximal variable update (10) can be removed. Under the setting, as shown in Appendix A, one can write (11) of the PPDM algorithm as

$$\mathbf{x}^{k+1} = \tilde{\mathbf{W}}\mathbf{s}^k - \boldsymbol{\Psi}^{-1}\nabla f(\mathbf{s}^k) + \mathbf{C}^k, \tag{16}$$

where $\mathbf{C}^k = (\tilde{\mathbf{U}}(\mathbf{x}^k - \mathbf{s}^k) + \kappa(\mathbf{x}^k - \mathbf{s}^k)\boldsymbol{\Psi}^{-1}) - \sum_{t=0}^k (\mathbf{I} - \tilde{\mathbf{W}})\mathbf{x}^t$ can regarded as a cumulative correction term. Note that the D-NG algorithm in [24, Eqn. (2)-(3)] is

$$\mathbf{s}^k = \mathbf{x}^k + \eta_k(\mathbf{x}^k - \mathbf{x}^{k-1}), \tag{17}$$
$$\mathbf{x}^{k+1} = \tilde{\mathbf{W}}\mathbf{s}^k - \boldsymbol{\Psi}^{-1}\nabla f(\mathbf{s}^k). \tag{18}$$

One can see that (18) and (16) have a similar form except for the correction term. Note that the convergence of the D-NG algorithm is proved in [24] only for convex problems with a diminishing step size. Therefore, the proposed PPDM algorithm is an enhanced counterpart of the D-NG algorithm with the ability to handle non-convex and non-smooth problems.

### C. Algorithm Development

In this subsection, let us elaborate how the SPPDM algorithm is devised. Our proposed algorithm is inspired by the proximal AL framework in [17]. First, we introduce a proximal term $\mathbf{z}$ to (5) as

$$\min_{\mathbf{x}, \mathbf{z}} \ f(\mathbf{x}) + r(\mathbf{x}) + \frac{\kappa}{2}\|\mathbf{x} - \mathbf{z}\|^2 \tag{19a}$$
$$\text{s.t. } \mathbf{A}\mathbf{x} = 0, \tag{19b}$$

where $\kappa > 0$ is a parameter. Obviously, (19) is equivalent to (5). The purpose of adding the proximal term $\frac{\kappa}{2}\|\mathbf{x} - \mathbf{z}\|^2$ is to make the objective function in (19a) strongly convex with respect to $\mathbf{x}$ when $\kappa > 0$ is large enough. Such strong convexity will be exploited for building the algorithm convergence.

Second, let us consider the AL function of (19) as follows

$$L_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda}) = f(\mathbf{x}) + r(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} \rangle$$
$$+ \frac{c}{2}\|\mathbf{A}\mathbf{x}\|^2 + \frac{\kappa}{2}\|\mathbf{x} - \mathbf{z}\|^2, \tag{20}$$

where $\boldsymbol{\lambda} \in \mathbb{R}^{|\mathcal{E}|}$ is the Lagrangian dual variable, and $c > 0$ is a positive penalty parameter. Then, the Lagrange dual problem of (19) can be expressed as

$$\max_{\boldsymbol{\lambda}} \min_{\mathbf{x}, \mathbf{z}} \ L_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda}). \tag{21}$$

We apply the following inexact stochastic primal-dual updates with momentum for problem (21): for $k = 0, 1, 2, \ldots,$

$$\boldsymbol{\lambda}^{k+1} = \boldsymbol{\lambda}^k + \alpha\mathbf{A}\mathbf{x}^k, \tag{22}$$
$$\mathbf{s}^k = \mathbf{x}^k + \eta_k(\mathbf{x}^k - \mathbf{x}^{k-1}), \tag{23}$$
$$\mathbf{x}^{k+1} = \arg\min_{\mathbf{x}} g(\mathbf{x}, \mathbf{x}^k, \mathbf{s}^k, \mathbf{z}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1}), \tag{24}$$
$$\mathbf{z}^{k+1} = \mathbf{z}^k + \beta(\mathbf{x}^{k+1} - \mathbf{z}^k). \tag{25}$$

Specifically, (22) is the dual ascent step with $\alpha > 0$ being the dual step size. In (23), the momentum variable $\mathbf{s}^k$ is introduced for the primal variable $\mathbf{x}$.

To update $\mathbf{x}$, we consider the inexact step as in (24) where $g(\mathbf{x}, \mathbf{x}^k, \mathbf{s}^k, \mathbf{z}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1})$ is a surrogate function given by

$$
\begin{aligned}
&g(\mathbf{x}, \mathbf{x}^k, \mathbf{s}^k, \mathbf{z}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1}) \\
&= \underbrace{f(\mathbf{s}^k) + \langle \bar{G}(\mathbf{s}^k, \boldsymbol{\xi}^k), \mathbf{x} - \mathbf{s}^k \rangle + \frac{\gamma}{2}\|\mathbf{x} - \mathbf{s}^k\|^2}_{(a)} \\
&\quad + \underbrace{r(\mathbf{x}) + \langle \boldsymbol{\lambda}^{k+1}, \mathbf{A}\mathbf{x} \rangle}_{} \\
&\quad + \underbrace{\frac{c}{2}\|\mathbf{A}\mathbf{x}\|^2 + \frac{c}{2}\|\mathbf{x} - \mathbf{x}^k\|_{\mathbf{B}^\top \mathbf{B}}^2 + \frac{\kappa}{2}\|\mathbf{x} - \mathbf{z}^k\|^2}_{(b)}. \quad (26)
\end{aligned}
$$

In (26), the term (a) is a quadratic approximation of $f$ at $\mathbf{s}^k$ using the stochastic gradient $\bar{G}$, where $\gamma > 0$ is a parameter. In term (b) of (26), $\mathbf{B}$ is the signless incidence matrix of the graph $\mathcal{G}$, i.e., $\mathbf{B} = |\mathbf{A}|$, which satisfies $\mathbf{A}^\top \mathbf{A} + \mathbf{B}^\top \mathbf{B} = 2\mathbf{D}$, where $\mathbf{D} = \text{diag}\{d_1, \ldots, d_N\}$ is the degree matrix of $\mathcal{G}$. As shown in [15], the introduction of $\frac{c}{2}\|\mathbf{x} - \mathbf{x}^k\|_{\mathbf{B}^\top \mathbf{B}}^2$ can "diagonalize" $\frac{c}{2}\|\mathbf{A}\mathbf{x}\|^2$ and lead to distributed implementation of (24). In particular, one can show that (24) with (26) can be expressed as

$$
\begin{aligned}
\mathbf{x}^{k+1} = \text{prox}_r^{\boldsymbol{\Psi}} \Big( \boldsymbol{\Psi}^{-1} \big( &\gamma \mathbf{s}^k + c\mathbf{B}^\top \mathbf{B}\mathbf{x}^k + \kappa \mathbf{z}^k \\
&- \bar{G}(\mathbf{s}^k, \boldsymbol{\xi}^k) - \mathbf{A}^\top \boldsymbol{\lambda}^{k+1} \big) \Big). \quad (27)
\end{aligned}
$$

As seen, due to the graphical structure of $\mathbf{B}^\top \mathbf{B}$, each $x_i^{k+1}$ in (27) can be obtained in a distributed fashion using only $x_j^k$, $j \in \mathcal{N}_i$ from its neighbors. Lastly, we update $\mathbf{z}$ by applying the gradient descent to $L_c(\mathbf{x}^{k+1}, \mathbf{z}; \boldsymbol{\lambda}^{k+1})$ with step size $\beta$, which then yields (25).

To show how (7)-(10) are obtained, let $\mathbf{p}^k = \mathbf{A}^\top \boldsymbol{\lambda}^k$ and define

$$
\begin{aligned}
\mathbf{x}^{k+\frac{1}{2}} = \boldsymbol{\Psi}^{-1} \big( &\gamma \mathbf{s}^k + c\mathbf{B}^\top \mathbf{B}\mathbf{x}^k + \kappa \mathbf{z}^k \\
&- \bar{G}(\mathbf{s}^k, \boldsymbol{\xi}^k) - \mathbf{p}^{k+1} \big). \quad (28)
\end{aligned}
$$

Then, (22) can be replaced by

$$
\mathbf{p}^{k+1} = \mathbf{p}^k + \alpha \mathbf{A}^\top \mathbf{A}\mathbf{x}^k, \quad (29)
$$

and (27) can be written as

$$
\mathbf{x}^{k+1} = \text{prox}_r^{\boldsymbol{\Psi}} \big( \mathbf{x}^{k+\frac{1}{2}} \big). \quad (30)
$$

Moreover, by subtracting $\mathbf{x}^{k-1+\frac{1}{2}}$ from $\mathbf{x}^{k+\frac{1}{2}}$, one obtains

$$
\begin{aligned}
\mathbf{x}^{k+\frac{1}{2}} = \mathbf{x}^{k-1+\frac{1}{2}} &+ \gamma \boldsymbol{\Psi}^{-1}(\mathbf{s}^k - \mathbf{s}^{k-1}) + \kappa \boldsymbol{\Psi}^{-1}(\mathbf{z}^k - \mathbf{z}^{k-1}) \\
&+ c\boldsymbol{\Psi}^{-1}\mathbf{B}^\top \mathbf{B}(\mathbf{x}^k - \mathbf{x}^{k-1}) - \boldsymbol{\Psi}^{-1}(\mathbf{p}^{k+1} - \mathbf{p}^k) \\
&- \boldsymbol{\Psi}^{-1}(\bar{G}(\mathbf{s}^k, \boldsymbol{\xi}^k) - \bar{G}(\mathbf{s}^{k-1}, \boldsymbol{\xi}^{k-1})). \quad (31)
\end{aligned}
$$

After substituting (29) into (31), we obtain

$$
\begin{aligned}
\mathbf{x}^{k+\frac{1}{2}} = \mathbf{x}^{k-1+\frac{1}{2}} &+ \mathbf{U}\mathbf{x}^k - \tilde{\mathbf{U}}\mathbf{x}^{k-1} \\
&+ \gamma \boldsymbol{\Psi}^{-1}(\mathbf{s}^k - \mathbf{s}^{k-1}) + \kappa \boldsymbol{\Psi}^{-1}(\mathbf{z}^k - \mathbf{z}^{k-1}) \\
&- \boldsymbol{\Psi}^{-1}(\bar{G}(\mathbf{s}^k, \boldsymbol{\xi}^k) - \bar{G}(\mathbf{s}^{k-1}, \boldsymbol{\xi}^{k-1})), \quad (32)
\end{aligned}
$$

which is exactly (11) since $\mathbf{U} = c\boldsymbol{\Psi}^{-1}\mathbf{B}^\top \mathbf{B} - \alpha \boldsymbol{\Psi}^{-1}\mathbf{A}^\top \mathbf{A}$ and $\tilde{\mathbf{U}} = c\boldsymbol{\Psi}^{-1}\mathbf{B}^\top \mathbf{B}$ by (12) and (13), respectively.

In summary, (22) and (24) can be equivalently written as (32) and (30), and therefore we obtain (23), (32), (30) and

(25) as the algorithm updates, which correspond to (7)-(10) in Algorithm 1

Before ending the section, we remark that it is possible to employ the existing stochastic primal-dual methods such as [28] for solving the non-smooth and non-convex problem (5). However, these methods require strict conditions on $\mathbf{A}$. For example, the stochastic ADMM method in [28] requires $\mathbf{A}$ to have full rank, which cannot happen for the distributed optimization problem (5) since the graph incidence matrix $\mathbf{A}$ for a connected graph must be rank deficient.

## III. CONVERGENCE ANALYSIS

In this section, we present the main theoretical results of the proposed SPPDM and PPDM algorithms by establishing their convergence conditions and convergence rate.

### A. Assumptions

We first make some proper assumptions on problem (5).

**Assumption 1.** *(i) The function $f(\mathbf{x})$ is a continuously differentiable function with Lipschitz continuous gradients, i.e., for constant $L > 0$,*

$$
\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad (33)
$$

*for all $\mathbf{x}, \mathbf{y}$. Moreover, assume that there exists a constant $\mu \geq -L$ (possibly negative) such that*

$$
f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{\mu}{2}\|\mathbf{x} - \mathbf{y}\|^2, \quad (34)
$$

*for all $\mathbf{x}, \mathbf{y}$.*
*(ii) The objective function $f(\mathbf{x}) + r(\mathbf{x})$ is bounded from below in the feasible set $\{\mathbf{x}|\mathbf{A}\mathbf{x} = 0\}$, i.e.,*

$$
f(\mathbf{x}) + r(\mathbf{x}) > \underline{f} > -\infty,
$$

*for some constant $\underline{f}$.*

**Assumption 2.** *The epigraph of each $r_i(x_i)$, i.e., $\{(x_i, y_i) \mid r_i(x_i) \leq y_i\}$, is a polyhedral set and has a compact form as*

$$
S_{x,i}x_i + S_{y,i}y_i \geq \zeta_i, \quad (35)
$$

*where $S_{x,i} \in \mathbb{R}^{q_i \times n}$, $S_{y,i} \in \mathbb{R}^{q_i}$ and $\zeta_i \in \mathbb{R}^{q_i}$ are some constant matrix and vectors.*

By (35), problem (5) can be written as

$$
\min_{\mathbf{x}, \mathbf{y}} \quad f(\mathbf{x}) + \mathbf{1}^\top \mathbf{y} \quad (36a)
$$

$$
\text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{0} \quad (36b)
$$

$$
\mathbf{S}_x \mathbf{x} + \mathbf{S}_y \mathbf{y} \geq \boldsymbol{\zeta}, \quad (36c)
$$

Here, $\mathbf{y} = [y_1, \ldots, y_N]^\top$, $\mathbf{S}_x = \text{diag}\{S_{x,1}, \ldots, S_{x,N}\}$, $\mathbf{S}_y = \text{diag}\{S_{y,1}, \ldots, S_{y,N}\}$, and $\boldsymbol{\zeta} = [\zeta_1^\top, \ldots, \zeta_N^\top]^\top$.

Let $\boldsymbol{\mu} = [\mu_1, \ldots, \mu_q]^\top \in \mathbb{R}^q$, $q = \sum_{i=1}^N q_i$, be the dual variable associated with (36c). Then, the Karush-Kuhn-Tucker (KKT) conditions of (36) are given by

$$
\nabla f(\mathbf{x}) + \mathbf{A}^\top \boldsymbol{\lambda} - \mathbf{S}_x^\top \boldsymbol{\mu} = 0, \quad \mathbf{S}_y^\top \boldsymbol{\mu} = \mathbf{1}, \quad (37a)
$$

$$
\mathbf{A}\mathbf{x} = 0, \quad \mathbf{S}_x \mathbf{x} + \mathbf{S}_y \mathbf{y} - \boldsymbol{\zeta} \geq 0, \quad \boldsymbol{\mu} \geq 0, \quad (37b)
$$

$$
\mu_j [\mathbf{S}_x \mathbf{x} + \mathbf{S}_y \mathbf{y} - \boldsymbol{\zeta}]_j = 0, \quad j = 1, \ldots, q. \quad (37c)
$$

For online/streaming learning, we also make the following standard assumptions that the gradient estimates are unbiased and have a bounded variance.

**Assumption 3.** *The stochastic gradient estimate $G_i(x, \xi)$ satisfies*

$$\mathbb{E}[G_i(x, \xi)] = \nabla f_i(x) \tag{38}$$

$$\mathbb{E}[\|G_i(x, \xi) - \nabla f_i(x)\|^2] \leq \sigma^2, \tag{39}$$

*for all $x$, where $\sigma > 0$ is a constant, and the expectation $\mathbb{E}$ is with respect to the random sample $\xi \sim \mathcal{B}_i$.*

It is easy to check that the gradient estimate of the mini-batch samples satisfies

$$\mathbb{E}\left[\left\|\frac{1}{|\mathcal{I}|}\sum_{j=1}^{|\mathcal{I}|} G_i(x, \xi_j) - \nabla f_i(x)\right\|^2\right] \leq \sigma^2/|\mathcal{I}|. \tag{40}$$

### B. Convergence Analysis of SPPDM

We define the following term

$$Q(\mathbf{x}, \boldsymbol{\lambda}) = \|\mathbf{x} - \text{prox}_r^1(\mathbf{x} - \nabla f(\mathbf{x}) - \mathbf{A}^\top \boldsymbol{\lambda})\|^2 + \|\mathbf{A}\mathbf{x}\|^2 \tag{41}$$

as the optimally gap for a primal-dual solution $(\mathbf{x}, \boldsymbol{\lambda})$ of problem (5). Obviously, one can shown that when $Q(\mathbf{x}^\star, \boldsymbol{\lambda}^\star) = 0$, $(\mathbf{x}^\star, \boldsymbol{\lambda}^\star)$ is a KKT solution of (5) which satisfies the conditions in (37) together with some $\mathbf{y}^\star$ and $\boldsymbol{\mu}^\star$. We define that $(\mathbf{x}^\star, \boldsymbol{\lambda}^\star)$ is an $\epsilon$-stationary solution of (5) if $Q(\mathbf{x}^\star, \boldsymbol{\lambda}^\star) < \epsilon$.

The convergence result is stated in the following theorem.

**Theorem 1.** *Assume that Assumptions 1-3 hold true, and let parameters satisfy*

$$\kappa > -\mu, \ \gamma > 3L, \tag{42}$$

$$\eta_k \leq \sqrt{\frac{\kappa + 2c + \gamma - 3L}{2(\gamma - \mu + 3L)}} := \bar{\eta}, \tag{43}$$

*moreover, let $0 < \alpha \leq c$ and $\beta > 0$ be both sufficiently small (see (89) and (90)). Then, for a sequence $\{\mathbf{x}^k, \mathbf{z}^k, \boldsymbol{\lambda}^k\}$ generated by Algorithm 1, it holds that*

$$\min_{k=0,\dots,K-1} \mathbb{E}[Q(\mathbf{x}^k, \boldsymbol{\lambda}^{k+1})] \leq C_0\left(\frac{\phi^0 - \underline{f}}{K} + \frac{C_1 N \sigma^2}{|\mathcal{I}|}\right), \tag{44}$$

*where $C_0$ and $C_1$ are some positive constants depending on the problem parameters (see (106) and (92)). In addition, $\phi^0$ is a constant defined in (73).*

To prove Theorem 1, the key is to define a novel stochastic potential function $\mathbb{E}[\phi^{k+1}]$ in (73) and analyze the conditions for which $\mathbb{E}[\phi^{k+1}]$ descends monotonically with the iteration number $k$ (Lemma 6). To achieve the goal, several approximation error bounds for the primal variable $\mathbf{x}^k$ (Lemma 2) and the dual variable $\boldsymbol{\lambda}^k$ (Lemma 4) are derived. Interested readers may refer to Appendix B for the details.

By Theorem 1, we immediately obtain the following corollary.

**Corollary 1.** *Let*

$$|\mathcal{I}| \geq \frac{2NC_0C_1\sigma^2}{\epsilon} \text{ and } K \geq \frac{2C_0(\phi^0 - \underline{f})}{\epsilon}. \tag{45}$$

*Then,*

$$\min_{k=0,\dots,K-1} \mathbb{E}[Q(\mathbf{x}^k, \boldsymbol{\lambda}^{k+1})] \leq \epsilon, \tag{46}$$

*that is, an $\epsilon$-stationary solution of problem (5) can be obtained in an expected sense.*

**Remark 3.** Given a mini-batch size $|\mathcal{I}| = \Omega(1/\epsilon)$, Corollary 1 implies that the proposed SPPDM algorithm has the convergence rate of $\mathcal{O}(1/\epsilon)$ to obtain an $\epsilon$-stationary solution. As a result, the corresponding communication complexity of the SPPDM algorithm is $\mathcal{O}(|\mathcal{E}|/\epsilon)$ while the computational complexity is $\mathcal{O}(N|\mathcal{I}|/\epsilon) = \mathcal{O}(N/\epsilon^2)$. As shown in Table I, the communication complexity $\mathcal{O}(1/\epsilon)$ of the SPPDM algorithm is smaller than $\mathcal{O}(1/\epsilon^2)$ of D-PSGD [20], $D^2$ [21], GNSD [22] and R-SGD-M [27]. The STOC-ADMM [28] has the same computation and communication complexity orders as the SPPDM algorithm, but it is not applicable to (5).

### C. Convergence Analysis of PPDM

When the full gradient $\nabla f(\mathbf{x}^k)$ is available for the PPDM algorithm, one can deduce a similar convergence result.

**Theorem 2.** *Assume Assumptions 1-2 and the same conditions in (42), (43), (90) and (89) hold true.*

- *Every limit point of the sequence $\{\mathbf{x}^k, \mathbf{z}^k, \boldsymbol{\lambda}^k\}$ generated by the PPDM algorithm is a KKT solution of (5).*
- *Given $K \geq \frac{C_0(\phi^0 - \underline{f})}{\epsilon}$, we have*

$$\min_{k=0,\dots,K-1} Q(\mathbf{x}^k, \boldsymbol{\lambda}^{k+1}) \leq C_0\left(\frac{\phi^0 - \underline{f}}{K}\right) \leq \epsilon.$$

The proof is presented in Appendix C..

To our knowledge, Theorem 1 and Theorem 2 are the first results that show the $\mathcal{O}(1/\epsilon)$ communication complexity of the distributed primal-dual method with momentum for non-convex and non-smooth problems. Numerical results in the next section will demonstrate that the SPPDM and PPDM algorithms can exhibit favorable convergence behavior than the existing methods.

## IV. NUMERICAL RESULTS

In this section, we examine the numerical performance of the proposed SPPDM/PPDM algorithm and present comparison results with the existing methods.

### A. Distributed Non-Convex Truncated Losses

We consider a linear regression model $y_j = h_j^\top x_* + \nu_j$, $j = 1, \dots, M$, where $M$ is the number of data samples. Here $y_j$ is the observed data sample and $h_j \in \mathbb{R}^n$ is the input data; $x_* \in \mathbb{R}^n$ is the ground truth; $\nu_j$ is the additive random noise.

Let $H := [h_1, \dots, h_M]^\top = [H_1^\top, \dots, H_N^\top]^\top \in \mathbb{R}^{M \times n}$, where each $H_i \in \mathbb{R}^{m \times n}$ corresponds to the data matrix owned by agent $i$ which has $m = M/N$ data points. The entries of $H$

are generated independently following the standard Gaussian distribution. The ground truth $x_*$ is a $S$-sparse vector to be estimated, whose non-zero entries are generated from the uniform distribution $U[-1, 1]$. The noise $\nu_j$ follows the Gaussian distribution $\mathcal{N}(0, 4)$. Then the data samples $y_j$, $j = 1, \ldots, M$, are generated by the above linear model.

Consider the following distributed regression problem with a nonconvex truncated loss [29]

$$\min_{x \in [-1, 1]} \sum_{i=1}^{N} \left( f_i(x) + \varsigma_i \|x\|_1 \right), \quad (47)$$

where

$$f_i(x) = \frac{\rho}{2N_i} \sum_{j=1}^{N_i} \log \left( 1 + \frac{\|y_j - a_j^\top x\|^2}{\rho} \right),$$

and $\rho$ is a parameter to determine the truncation level. We set $m = 150$, $n = 256$, $S = 16$, and $\rho = 3$. Moreover, we consider a circle graph with $N = 20$ agents.

For the online setting, we compare the SPPDM algorithm (Algorithm 1) with PSGD [23] and STOC-ADMM [28]. For the offline setting, we compare the PPDM algorithm with Prox-DGD [11], PG-EXTRA [8], Prox-ADMM [17] and STOC-ADMM [28]. Note that theoretically PG-EXTRA and STOC-ADMM are not guaranteed to converge for the nonconvex problem (5). We implement these two methods simply for comparison purpose.

For the PG-EXTRA, we choose the stepsize $\ell = 0.05$ according to the sufficient condition suggested in [8]. According to their convergence conditions, the diminishing step size $\ell = \frac{1}{3\sqrt{k+100}}$ is used for the PSGD and Prox-DGD. The primal and dual stepsize for the STOC-ADMM is chosen according to the convergence condition suggested in [28].

For PSGD, Prox-DGD and PG-EXTRA, the mixing matrix follows the metropolis weight

$$[\mathbf{W}]_{ij} \triangleq \begin{cases} \frac{1}{\max\{d_i, d_j\}+1}, & \text{for } (i,j) \in \mathcal{E}, \\ 0, & \text{for } (i,j) \notin \mathcal{E} \text{ and } i \neq j. \\ 1 - \sum_{j \neq i} w_{ij}, & \text{for } i = j \end{cases} \quad (48)$$

If not specified, the parameters of the SPPDM/PPDM and the Prox-ADMM are given as $\alpha = 2$, $\kappa = 1$, $c = 2$, $\gamma = 3$, $\beta = 0.9$[1] For the proposed SPPDM, we consider two cases about $\eta_k$, one is $\eta_k = 0$ without momentum, and the other is based on the Nesterov's extrapolation technique, i.e.,

$$\eta_k = \frac{\theta_{k-1} - 1}{\theta_k}, \quad \theta_{k+1} = \frac{1 + \sqrt{1 + 4\theta_k^2}}{2},$$

with $\theta_{-1} = \theta_0 = 1$. When $\eta_k = 0$, we denote SPPDM as SPPD.

---

[1]By analysis, the Hessian matrix for the function $f_i(x)$ is $\frac{1}{N_i} \sum_{j=1}^{N_i} \frac{\rho h_j h_j^T (\rho - \|h_j^T x - y_j\|^2)}{(\rho + \|h_j^T x - y_j\|^2)^2}$. It shows that the maximum eigenvalue of this Hessian matrix is smaller than 1 ($L < 1$) with the given parameter. Thus, the parameters of SPPDM/PPDM satisfy the conditions stated in Theorem 1.



Fig. 1. Comparison of proposed SPPDM with the PSGD and STOC-ADMM in terms of stationarity and consensus error; the batch size is $|\mathcal{I}| = 100$.

Define $\bar{x} = \frac{1}{N} \sum_{i=1}^{N} x_i$. The stationarity error and consensus error are defined below

$$\text{stationarity error} = \|\bar{x} - \text{prox}_r[\bar{x} - \nabla f(\bar{x})]\|^2,$$

$$\text{consensus error} = \frac{1}{N} \sum_{i=1}^{N} \|x_i - \bar{x}\|^2.$$

We run 10 independents trials for each algorithm with randomly generated data and random initial values. The convergence curves obtained by averaging over all 10 trials are plotted in Figs. 1-3.

In Fig. 1, we observe that the SPPDM, SPPM and STOC-ADMM all perform better than the PSGD in terms of stationarity error and consensus error. The reason is that these methods all use constant step sizes rather than the diminishing step size. In addition, the proposed SPPDM has better performance than SPPD and STOC-ADMM, due to the Nesterov momentum.

The impact of the mini-batch size $|\mathcal{I}|$, and parameters $\gamma$ and $c$ are analyzed in Fig. 2. One can see that the larger mini-batch size we use, the smaller error we can achieve, which corroborates Corollary 1. With the same mini-batch size, the larger values of $\gamma$ and $c$ correspond to smaller primal and dual step sizes. Thus, the SPPDM with larger values of $\gamma$ and $c$ has slower convergence; whereas as seen from the figures, larger values of $\gamma$ and $c$ can lead to smaller stationarity error and consensus error.

For the offline setting, the comparison results of the proposed PPDM with the existing methods are shown in Fig. 3. It can be observed that the proposed PPDM enjoys the fastest convergence. Compared with the Prox-ADMM, it is clear to see the advantage of the PPDM with momentum for speeding up the algorithm convergence.

### B. Distributed Neural Network

In this simulation, our task is to classify handwritten digits from the MNIST dataset. The local loss function $f_i(\theta_i)$ in each node is the cross-entropy function. In this example, we do not consider nonsmooth term and inequality constraint set. Thus, many existing methods, D-PSGD [20], $D^2$ [21] and PR-SGD-M [27] can be applied to train a classification DNN.

Assume the neural network contains one hidden layer with 500 neurons. The $6 \times 10^4$ training samples are divided into 10 subsets and assigned to the $N = 10$ agents in two ways.

Fig. 2. Comparison of the proposed SPPDM with different parameters in terms of stationarity and consensus error.



Fig. 3. Comparison of proposed PPDM with the existing methods in terms of stationarity and consensus error.



Fig. 4. Comparison of proposed SPPDM/SPPD algorithms with different methods under the IID case.



Fig. 5. Comparison of proposed SPPDM/SPPD algorithms with different methods under the Non-IID case.

The first is the *IID* case, where the samples are sufficiently shuffled, and then partitioned into 10 subsets with equal size ($m = 6000$). The second is the *Non-IID* case, where we first sort the samples according to their labels, divide it into 20 shards of size 3000, and assign each of 10 agents 2 shards. Thus most agents have samples of two digits only.

The communication graph is also a circle. We compare the SPPDM with the D-PSGD [20], $D^2$ [21] and PR-SGD-M [27]. The same mixing matrix in (48) is used for the three methods. Moreover, a fixed step size of $\ell = 0.05$ is used to ensure the convergence of these three methods in the simulation. For the proposed SPPDM, we set parameter $c = 1$, $\gamma = 3$, $\alpha = 0.001$, $\kappa = 0.1$, $\beta = 0.9$, and $\eta_k = 0.8$. The batch size is $|\mathcal{I}| = 128$.

We calculate the loss value and the classification accuracy based on the average model $\bar{\theta} = 1/N \sum_{i=1}^{N} \theta_i$. Fig. 4 and Fig. 5 show the training loss and the classification accuracy for the IID case and Non-IID case by averaging over all 5 trials, respectively. From Fig. 4, we see that $D^2$ and D-PSGD have a similar performance; meanwhile, the proposed SPPD performs better than these two methods. Besides, one can see that SPPDM and PR-SGD-M enjoy fast decreasing of the loss function and increasing of the classification accuracy, respectively. The reason is that both SPPDM and PR-SGD-M use the momentum technique. We should point out that the communication overhead of PR-SGD-M is twice of the SPPDM since the PR-SGD-M requires the agents to exchange not only the variable $x_i$ but also the momentum variables. Lastly, comparing the SPPDM with SPPD, it shows again that the momentum techniques can accelerate the algorithm convergence.

From Fig. 5 for the non-IID case, we can observe that

the $D^2$ performs better than D-PSGD and SPPD. In fact, by comparing the curves in Fig. 4 with those in Fig. 5, one can see that the convergence curve of $D^2$ remains almost the same due to the use of the variance reduction technique whereas D-PSGD and SPPD deteriorate under the non-IID setting. As seen, the convergence of SPPDM is also slowed, but it still performs best among the methods under test.

## V. CONCLUSION

In this paper, we have proposed a distributed stochastic proximal primal-dual algorithm with momentum for minimizing a non-convex and non-smooth function (5) over a connected multi-agent network. We have shown (in Remark 1 and Remark 2) that the proposed algorithm has a close connection with some of the existing algorithms that are for convex and smooth problems, and therefore can be regarded as an enhanced counterpart of these existing algorithms. Theoretically, under Assumptions 1-3, we have built the convergence conditions of the proposed algorithms in Theorem 1 and Theorem 2. In particular, we have shown that the proposed SPPDM can achieve an $\epsilon$-stationary solution with $\mathcal{O}(1/\epsilon^2)$ computational complexity and $\mathcal{O}(1/\epsilon)$ communication complexity, where the latter is better than many of the existing methods which have $\mathcal{O}(1/\epsilon^2)$ communication complexity (see Table 1). Experimental results have demonstrated that the proposed algorithms with momentum can effectively speed up the convergence. For distributed learning under non-IID data distribution (Fig. 5), we have also shown the proposed SPPDM performs better than the existing methods.

As future research directions, one may further relax Assumption 3 to accommodate a larger class of regularization functions. Besides, it will be also interesting to investigate the trade-off between the communication complexity (as measured

by the number of bits exchanged) and the convergence. Moreover, we will analytically investigate how data distribution affect the algorithm convergence and improve the robustness of the algorithms against to unbalanced and non-IID data distribution in the future.

## APPENDIX A
### DERIVATION OF (16)

When we remove the non-smooth regularization term $r(\mathbf{x})$, the proximal gradient update can be removed. Assume $\beta = 1$, then (27) can be written as

$$\mathbf{x}^{k+1} = \mathbf{\Psi}^{-1}\big(\gamma \mathbf{s}^k + c\mathbf{B}^\top \mathbf{B}\mathbf{x}^k + \kappa \mathbf{x}^k - \nabla f(\mathbf{s}^k) - \mathbf{A}^\top \boldsymbol{\lambda}^{k+1}\big)$$
$$= \mathbf{\Psi}^{-1}\big(\tilde{\mathbf{W}}\mathbf{s}^k - \mathbf{\Psi}^{-1}\nabla f(\mathbf{s}^k) + \mathbf{C}^k\big),$$

where

$$\tilde{\mathbf{W}} = c\mathbf{\Psi}^{-1}\mathbf{B}^\top \mathbf{B} + (\gamma + \kappa)\mathbf{\Psi}^{-1}, \tag{49}$$
$$\mathbf{C}^k = \mathbf{\Psi}^{-1}(c\mathbf{B}^\top \mathbf{B}(\mathbf{x}^k - \mathbf{s}^k) + \kappa(\mathbf{x}^k - \mathbf{s}^k)) - \mathbf{\Psi}^{-1}\mathbf{A}^\top \boldsymbol{\lambda}^{k+1}. \tag{50}$$

Using the definition of $\tilde{\mathbf{U}} = c\mathbf{\Psi}^{-1}\mathbf{B}^\top \mathbf{B}$ in (13), we rewrite (49) as

$$\tilde{\mathbf{W}} = \tilde{\mathbf{U}} + \mathbf{\Psi}^{-1}(\gamma + \kappa)\mathbf{I}.$$

According the definition of $\mathbf{\Psi} = (\gamma + \kappa)\mathbf{I} + c(\mathbf{A}^\top \mathbf{A} + \mathbf{B}^T \mathbf{B})$ and $c = \alpha$, we have

$$\mathbf{I} - \tilde{\mathbf{W}} = \mathbf{\Psi}^{-1}\mathbf{\Psi} - \tilde{\mathbf{W}} = \alpha \mathbf{\Psi}^{-1}\mathbf{A}^\top \mathbf{A}. \tag{51}$$

Based on

$$\mathbf{A}^\top \boldsymbol{\lambda}^{k+1} = \mathbf{p}^{k+1} = \mathbf{p}^k + \alpha \mathbf{A}^\top \mathbf{A}\mathbf{x}^k, \tag{52}$$

if $\mathbf{p}^0 = 0$, and applying (51), we obtain

$$\mathbf{\Psi}^{-1}\mathbf{A}^\top \boldsymbol{\lambda}^{k+1} = \sum_{t=0}^{k} \alpha \mathbf{\Psi}^{-1}\mathbf{A}^\top \mathbf{A}\mathbf{x}^t = \sum_{t=0}^{k} (\mathbf{I} - \tilde{\mathbf{W}})\mathbf{x}^t.$$

By substituting the above equality into (50), we obtain (16).

## APPENDIX B
### PROOF OF THEOREM 1

Let us recapitulate the augmented Lagrange function in (20) below

$$L_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda}) = f(\mathbf{x}) + r(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} \rangle$$
$$+ \frac{c}{2}\|\mathbf{A}\mathbf{x}\|^2 + \frac{\kappa}{2}\|\mathbf{x} - \mathbf{z}\|^2. \tag{53}$$

We introduce some auxiliary functions as follows

$$d(\mathbf{z}; \boldsymbol{\lambda}) = \min_{\mathbf{x}} L_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda}) \tag{54}$$
$$\mathbf{x}(\mathbf{z}; \boldsymbol{\lambda}) = \arg\min_{\mathbf{x}} L_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda}) \tag{55}$$
$$P(\mathbf{z}) = \min_{\mathbf{A}\mathbf{x}=0} f(\mathbf{x}) + r(\mathbf{x}) + \frac{\kappa}{2}\|\mathbf{x} - \mathbf{z}\|^2 \tag{56}$$
$$\mathbf{x}(\mathbf{z}) = \arg\min_{\mathbf{A}\mathbf{x}=0} f(\mathbf{x}) + r(\mathbf{x}) + \frac{\kappa}{2}\|\mathbf{x} - \mathbf{z}\|^2. \tag{57}$$

Besides, we define the full gradient iterate $\hat{\mathbf{x}}^{k+1}$ and $\hat{\mathbf{z}}^{k+1}$,

$$\hat{\mathbf{x}}^{k+1} := \arg\min_{\mathbf{x}} g(\mathbf{x}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1}) \tag{58}$$
$$\hat{\mathbf{z}}^{k+1} := \mathbf{z}^k + \beta(\hat{\mathbf{x}}^{k+1} - \mathbf{z}^k), \tag{59}$$

where $\mathbf{w}^k = [\mathbf{x}^k, \mathbf{s}^k, \mathbf{z}^k]$ and

$$g(\mathbf{x}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1})$$
$$= f(\mathbf{s}^k) + \langle \nabla f(\mathbf{s}^k), \mathbf{x} - \mathbf{s}^k \rangle + \frac{\gamma}{2}\|\mathbf{x} - \mathbf{s}^k\|^2 + r(\mathbf{x})$$
$$+ \langle \boldsymbol{\lambda}^{k+1}, \mathbf{A}\mathbf{x} \rangle + \frac{c}{2}\|\mathbf{A}\mathbf{x}\|^2 + \frac{c}{2}\|\mathbf{x} - \mathbf{x}^k\|_{\mathbf{B}^\top \mathbf{B}}^2 + \frac{\kappa}{2}\|\mathbf{x} - \mathbf{z}^k\|^2. \tag{60}$$

We also define

$$g(\mathbf{x}, \mathbf{w}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1}) := g(\mathbf{x}, \mathbf{x}^k, \mathbf{s}^k, \mathbf{z}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1})$$

for (26) at our disposal.

### A. Some Error Bounds

Firstly, we show the upper bound between $\mathbf{x}^{k+1}$ and $\hat{\mathbf{x}}^{k+1}$.

**Lemma 1.** *Suppose Assumption 3 holds, we have*

$$\mathbb{E}[\|\mathbf{x}^{k+1} - \hat{\mathbf{x}}^{k+1}\|^2] \leq \frac{N\sigma^2}{(\gamma + 2c + \kappa)^2|\mathcal{I}|}. \tag{61}$$

*Proof.* According to (38)-(39) in Assumption 3, we know

$$\mathbb{E}\Big[\big\|\frac{1}{|\mathcal{I}|}\sum_{j=1}^{|\mathcal{I}|} G(\mathbf{s}^k, \boldsymbol{\xi}_j^k) - \nabla f(\mathbf{x})\big\|^2\Big] \leq \frac{N}{|\mathcal{I}|}\sigma^2. \tag{62}$$

In addition, like (27), the proximity form of (58) is

$$\hat{\mathbf{x}}^{k+1} = \text{prox}_r^{\mathbf{\Psi}}\Big(\mathbf{\Psi}^{-1}\big(\gamma \mathbf{s}^k + c\mathbf{B}^\top \mathbf{B}\mathbf{x}^k + \kappa \mathbf{z}^k$$
$$- \nabla f(\mathbf{s}^k) - \mathbf{A}^\top \boldsymbol{\lambda}^{k+1}\big)\Big). \tag{63}$$

Using (27), (62)-(63) and applying the nonexpansive property of the proximal operator (see for example [30, p. 340]) we then obtain (61). □

**Lemma 2.** *Suppose $\kappa > -\mu$. There exists some positive constants $\sigma_1$, $\sigma_2$ such that the following primal error bound holds*

$$\|\mathbf{x}^k - \mathbf{x}(\mathbf{z}^k; \boldsymbol{\lambda}^{k+1})\| \leq \sigma_1 \|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\| + \sigma_2 \|\mathbf{x}^k - \mathbf{s}^k\|. \tag{64}$$

*Proof.* Based on $\kappa > -\mu$, we know that $L_c$ in (53) is strongly convex in $\mathbf{x}$ with modulus $\kappa + \mu$ and Lipschitz constant $\kappa + L + c\sigma_A^2$, where $\sigma_A$ is the spectral norm of the matrix. Thus, we can apply [31, Theorem 3.1] to upper bound the distance between $\mathbf{x}^k$ and the optimal solution $\mathbf{x}(\mathbf{z}^k; \boldsymbol{\lambda}^{k+1})$

$$\|\mathbf{x}^k - \mathbf{x}(\mathbf{z}^k; \boldsymbol{\lambda}^{k+1})\| \leq \varrho\|\tilde{\nabla}_{\mathbf{x}} L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1})\|, \tag{65}$$

where $\varrho = \frac{\kappa + L + c\sigma_A^2 + 1}{\kappa + \mu}$ and

$$\tilde{\nabla}_{\mathbf{x}} L_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda}) = \mathbf{x} - \text{prox}_r^{\mathbf{\Psi}}(\mathbf{x} - \nabla_{\mathbf{x}}(L_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda}) - r(\mathbf{x})))$$

is known as the proximal gradient.

We can bound $\|\tilde{\nabla}_{\mathbf{x}} L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1})\|$ as follows

$$\|\tilde{\nabla}_{\mathbf{x}} L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1})\|$$
$$= \|\mathbf{x}^k - \mathrm{prox}_r^{\boldsymbol{\Psi}}(\mathbf{x}^k - \nabla_{\mathbf{x}}(L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1}) - r(\mathbf{x}^k)))\|$$
$$\leq \|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|$$
$$\quad + \|\hat{\mathbf{x}}^{k+1} - \mathrm{prox}_r^{\boldsymbol{\Psi}}(\mathbf{x}^k - \nabla_{\mathbf{x}}(L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1}) - r(\mathbf{x}^k)))\|$$
$$= \|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|$$
$$\quad + \left\| \mathrm{prox}_r^{\boldsymbol{\Psi}}(\hat{\mathbf{x}}^{k+1} - \nabla_{\mathbf{x}}(g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1}) - r(\hat{\mathbf{x}}^{k+1}))) \right.$$
$$\quad \left. - \mathrm{prox}_r^{\boldsymbol{\Psi}}(\mathbf{x}^k - \nabla_{\mathbf{x}}(L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1}) - r(\mathbf{x}^k))) \right\|$$
$$\leq (2 + 2cd_{\max} + \gamma + \kappa)\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\| + (\gamma + L)\|\mathbf{x}^k - \mathbf{s}^k\|,$$

where $d_{\max} = \max\{d_1, \ldots, d_N\}$; the second equality is obtained by using the optimality condition of $\hat{\mathbf{x}}^{k+1}$ in (58), and the second inequality is based on the nonexpansive property of the proximal operator. Denote

$$\sigma_1 = \varrho(2 + 2cd_{\max} + \gamma + \kappa), \tag{66}$$
$$\sigma_2 = \gamma + L. \tag{67}$$

The proof is complete. $\qquad\square$

**Lemma 3.** *(Lemma 3.2 in [17]) Suppose $\kappa > -\mu$, and Assumption 1 holds. There exists some positive constants $\sigma_3, \sigma_4$ such that the following error bounds hold*

$$\|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2\| \geq \sigma_3 \|\mathbf{x}(\mathbf{z}; \boldsymbol{\lambda}_1) - \mathbf{x}(\mathbf{z}; \boldsymbol{\lambda}_2)\| \tag{68}$$
$$\|\mathbf{z}_1 - \mathbf{z}_2\| \geq \sigma_4 \|\mathbf{x}(\mathbf{z}_1; \boldsymbol{\lambda}) - \mathbf{x}(\mathbf{z}_2; \boldsymbol{\lambda})\|, \tag{69}$$

*where*

$$\sigma_3 = (\kappa + \mu)/\sigma_A \tag{70}$$
$$\sigma_4 = (\kappa + \mu)/\kappa. \tag{71}$$

**Lemma 4.** *Suppose that Assumptions 1-2 hold and $\kappa > \mu$. Then, there exist some positive scalars $\sigma_5$, $\Delta$ such that the following dual error bound holds*

$$\|\mathbf{x}(\mathbf{z}, \boldsymbol{\lambda}) - \mathbf{x}(\mathbf{z})\| \leq \sigma_5 \|\mathbf{A}\mathbf{x}(\mathbf{z}; \boldsymbol{\lambda})\|, \ for\ any\ \mathbf{z},\ \boldsymbol{\lambda}. \tag{72}$$

*where $\sigma_5$ depends only on the constants $L, \kappa, \sigma_A, \mu$ and the matrices $\mathbf{A}, \mathbf{S}_x, \mathbf{S}_y$.*

*Proof.* The lemma is an extension of [18, Lemma 3.2], where the non-smooth term $r(\mathbf{x})$ of (5) is limited to an indicator function of a polyhedral set. Due to limited space, the proof details are relegated to the supplementary document [32]. $\quad\square$

### B. Decent Lemmas

In order to show the convergence of Algorithm 1, we consider a new potential function,

$$\mathbb{E}[\phi^{k+1}] \triangleq \mathbb{E}[L_c(\mathbf{x}^{k+1}, \mathbf{z}^{k+1}; \boldsymbol{\lambda}^{k+1}) + \tau\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2]$$
$$\quad + \mathbb{E}[2P(\mathbf{z}^{k+1}) - 2d(\mathbf{z}^{k+1}; \boldsymbol{\lambda}^{k+1})], \tag{73}$$

for some $\tau > 0$. By the weak duality, we have

$$L_c(\mathbf{x}, \mathbf{z}; \boldsymbol{\lambda}) \geq d(\mathbf{z}; \boldsymbol{\lambda}), P(\mathbf{z}) \geq d(\mathbf{z}; \boldsymbol{\lambda}). \tag{74}$$

Thus, we have $\mathbb{E}[\phi^k] \geq \mathbb{E}[P(\mathbf{z}^k)]$. According to the definition of $P(\mathbf{z}^k)$ in (56) and Assumption 1 (ii), we obtain $P(\mathbf{z}^k) \geq \underline{f}$. As a result, $\mathbb{E}[\phi^k]$ is bounded below by $\underline{f}$.

**Lemma 5.** *For a sequence $\{\mathbf{x}^k, \mathbf{z}^k, \boldsymbol{\lambda}^k\}$ generated by Algorithm 1, if $\kappa > -\mu$, $\gamma > 3L$, $0 < \beta < 1$ and*

$$0 \leq \eta_k \leq \sqrt{\frac{\kappa + 2c + \gamma - 3L}{2(\gamma - \mu + 3L)}} := \bar{\eta}, \tag{75}$$

*there exit some positive constants $\tau$, $\hat{\sigma}_1$ and $\hat{\sigma}_2$ such that*

$$\hat{\sigma}_1 \triangleq \frac{\kappa + 2c + \gamma - 3L}{2} - 2\tau \geq 0 \tag{76}$$
$$\hat{\sigma}_2 \triangleq \frac{\mu - \gamma - 3L}{2}\bar{\eta}^2 + \tau \geq 0, \tag{77}$$

*then*

$$\mathbb{E}[L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^k) + \tau\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2]$$
$$\quad - \mathbb{E}[L_c(\mathbf{x}^{k+1}, \mathbf{z}^{k+1}; \boldsymbol{\lambda}^{k+1}) - \tau\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2]$$
$$\geq -\alpha\mathbb{E}[\|\mathbf{A}\mathbf{x}^k\|^2] + \frac{\kappa}{2\beta}\mathbb{E}[\|\mathbf{z}^k - \mathbf{z}^{k+1}\|^2]$$
$$\quad + \hat{\sigma}_1\mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2] + \hat{\sigma}_2\mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2]$$
$$\quad - \left( \frac{1}{2\gamma} + \frac{3L + 4\tau}{2(\gamma + 2c + \kappa)^2} \right) \frac{N\sigma^2}{|\mathcal{I}|}. \tag{78}$$

*Proof.* Firstly, according to (20) and (22), we have

$$\mathbb{E}[L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^k) - L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1})] = -\alpha\mathbb{E}[\|\mathbf{A}\mathbf{x}^k\|^2]. \tag{79}$$

Secondly, we have

$$\mathbb{E}[L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1}) - L_c(\mathbf{x}^{k+1}, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1})]$$
$$= \mathbb{E}[L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1}) - g(\mathbf{x}^k, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1})]$$
$$\quad + \mathbb{E}[g(\mathbf{x}^k, \mathbf{w}^k\boldsymbol{\lambda}^{k+1}) - g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1})]$$
$$\quad + \mathbb{E}[g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1}) - g(\mathbf{x}^{k+1}, \mathbf{w}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1})]$$
$$\quad + \mathbb{E}[g(\mathbf{x}^{k+1}, \mathbf{w}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1}) - L_c(\mathbf{x}^{k+1}, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1})]. \tag{80}$$

Next, we bound each of the terms in the right hand side of (80). Based on the definition of function $g$ in (60), we have

$$\mathbb{E}[L_c(\mathbf{x}^k, \mathbf{z}^k; \boldsymbol{\lambda}^{k+1}) - g(\mathbf{x}^k, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1})]$$
$$= \mathbb{E}[f(\mathbf{x}^k) - f(\mathbf{s}^k) - \langle \nabla f(\mathbf{s}^k), \mathbf{x}^k - \mathbf{s}^k \rangle - \frac{\gamma}{2}\|\mathbf{x}^k - \mathbf{s}^k\|^2]$$
$$\geq \frac{\mu - \gamma}{2}\mathbb{E}[\|\mathbf{x}^k - \mathbf{s}^k\|^2], \tag{81}$$

where the inequality comes from (34) in Assumption 1. Using the strongly convexity of the objective function $g$ (with modulus $\kappa + 2c + \gamma$) and the definition of $\hat{\mathbf{x}}^{k+1}$ in (58), we can obtain

$$\mathbb{E}[g(\mathbf{x}^k, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1}) - g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1})]$$
$$\geq \frac{\kappa + 2c + \gamma}{2}\mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2]. \tag{82}$$

In addition, we have

$$\mathbb{E}[g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1}) - g(\mathbf{x}^{k+1}, \mathbf{w}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1})]$$
$$= \mathbb{E}[g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1}) - g(\mathbf{x}^{k+1}, \mathbf{w}^k, \boldsymbol{\xi}^k; \boldsymbol{\lambda}^{k+1})]$$
$$\geq 0, \tag{83}$$

where the first equality dues to (38) in Assumption 3, and the above inequality dues to $\mathbf{x}^{k+1}$ is the optimal solution in (24). Lastly, we can bound

$$
\mathbb{E}[g(\mathbf{x}^{k+1},\mathbf{w}^k,\boldsymbol{\xi}^k;\boldsymbol{\lambda}^{k+1}) - L_c(\mathbf{x}^{k+1},\mathbf{z}^k;\boldsymbol{\lambda}^{k+1})]
$$
$$
= \mathbb{E}[f(\mathbf{s}^k)] + \frac{1}{|\mathcal{I}|}\sum_{j=1}^{|\mathcal{I}|}\mathbb{E}[\langle G(\mathbf{s}^k,\boldsymbol{\xi}_j^k),\mathbf{x}^{k+1}-\mathbf{s}^k\rangle]
$$
$$
+ \mathbb{E}\left[\frac{\gamma}{2}\|\mathbf{x}^{k+1}-\mathbf{s}^k\|^2 + \frac{c}{2}\|\mathbf{x}^{k+1}-\mathbf{x}^k\|_{\mathbf{B}^\top\mathbf{B}}^2 - f(\mathbf{x}^{k+1})\right]
$$
$$
\geq \frac{1}{|\mathcal{I}|}\sum_{j=1}^{|\mathcal{I}|}\mathbb{E}[\langle G(\mathbf{s}^k,\boldsymbol{\xi}_j^k)-\nabla f(\mathbf{s}^k),\mathbf{x}^{k+1}-\mathbf{s}^k\rangle]
$$
$$
+ \frac{\gamma-L}{2}\mathbb{E}[\|\mathbf{x}^{k+1}-\mathbf{s}^k\|^2] + \frac{c}{2}\mathbb{E}[\|\mathbf{x}^{k+1}-\mathbf{x}^k\|_{\mathbf{B}^\top\mathbf{B}}^2]
$$
$$
\geq -\frac{N\sigma^2}{2\gamma|\mathcal{I}|} - \frac{L}{2}\mathbb{E}[\|\mathbf{x}^{k+1}-\mathbf{s}^k\|^2], \tag{84}
$$

where the first inequality is obtained by applying the descent lemma [33, Lemma1.2.3]

$$
f(\mathbf{x}^{k+1}) \leq f(\mathbf{s}^k) + \langle\nabla f(\mathbf{s}^k),\mathbf{x}^{k+1}-\mathbf{s}^k\rangle + \frac{L}{2}\|\mathbf{x}^{k+1}-\mathbf{s}^k\|^2
$$

owing to gradient Lipschitz continuity in (33); the second inequality holds by using the Young's inequality $a^\top b \geq -\frac{\|a\|^2}{2\gamma}-\frac{\gamma\|b\|^2}{2}$ and (39) in Assumption 3. Using the convexity of the operator $\|\cdot\|^2$, we have

$$
\|\mathbf{x}^{k+1}-\mathbf{s}^k\|^2 \leq 3\|\mathbf{x}^{k+1}-\hat{\mathbf{x}}^{k+1}\|^2 + 3\|\hat{\mathbf{x}}^{k+1}-\mathbf{x}^k\|^2
$$
$$
+ 3\|\mathbf{x}^k-\mathbf{s}^k\|^2. \tag{85}
$$

Substituting (61) and (85) into (84) gives rise to

$$
\mathbb{E}[g(\mathbf{x}^{k+1},\mathbf{w}^k,\boldsymbol{\xi}^k;\boldsymbol{\lambda}^{k+1}) - L_c(\mathbf{x}^{k+1},\mathbf{z}^k;\boldsymbol{\lambda}^{k+1})]
$$
$$
\geq -\left(\frac{1}{2\gamma}+\frac{3L}{2(\gamma+2c+\kappa)^2}\right)\frac{N\sigma^2}{|\mathcal{I}|}
$$
$$
- \frac{3L}{2}\mathbb{E}[\|\hat{\mathbf{x}}^{k+1}-\mathbf{x}^k\|^2] - \frac{3L}{2}\mathbb{E}[\|\mathbf{s}^k-\mathbf{x}^k\|^2]. \tag{86}
$$

By further substituting (81)-(83) and (86) into (80), we obtain

$$
\mathbb{E}[L_c(\mathbf{x}^k,\mathbf{z}^k;\boldsymbol{\lambda}^{k+1}) + \tau\|\mathbf{x}^k-\mathbf{x}^{k-1}\|^2
$$
$$
- L_c(\mathbf{x}^{k+1},\mathbf{z}^k;\boldsymbol{\lambda}^{k+1}) - \tau\|\mathbf{x}^{k+1}-\mathbf{x}^k\|^2]
$$
$$
\geq \frac{\mu-\gamma}{2}\mathbb{E}[\|\mathbf{x}^k-\mathbf{s}^k\|^2] + \frac{\kappa+2c+\gamma}{2}\mathbb{E}[\|\mathbf{x}^k-\hat{\mathbf{x}}^{k+1}\|^2]
$$
$$
- \frac{3L}{2}\mathbb{E}[\|\hat{\mathbf{x}}^{k+1}-\mathbf{x}^k\|^2] - \frac{3L}{2}\mathbb{E}[\|\mathbf{s}^k-\mathbf{x}^k\|^2]
$$
$$
- \left(\frac{1}{2\gamma}+\frac{3L}{2(\gamma+2c+\kappa)^2}\right)\frac{\sigma^2}{|\mathcal{I}|} + \tau\mathbb{E}[\|\mathbf{x}^k-\mathbf{x}^{k-1}\|^2]
$$
$$
- 2\tau\mathbb{E}[\|\mathbf{x}^{k+1}-\hat{\mathbf{x}}^{k+1}\|^2] - 2\tau\mathbb{E}[\|\hat{\mathbf{x}}^{k+1}-\mathbf{x}^k\|^2]
$$
$$
= \hat{\sigma}_1\mathbb{E}[\|\mathbf{x}^k-\hat{\mathbf{x}}^{k+1}\|^2] + \hat{\sigma}_2\mathbb{E}[\|\mathbf{x}^k-\mathbf{x}^{k-1}\|^2]
$$
$$
- \left(\frac{1}{2\gamma}+\frac{3L+4\tau}{2(\gamma+2c+\kappa)^2}\right)\frac{N\sigma^2}{|\mathcal{I}|}, \tag{87}
$$

where $\hat{\sigma}_1$ and $\hat{\sigma}_2$ are defined in (76) and (77), respectively, and the equality is obtained by applying (23).

Thirdly, according to the definition of the $\mathbf{z}$ update in (25), we have

$$
\mathbb{E}[L_c(\mathbf{x}^{k+1},\mathbf{z}^k;\boldsymbol{\lambda}^{k+1}) - L_c(\mathbf{x}^{k+1},\mathbf{z}^{k+1};\boldsymbol{\lambda}^{k+1})]
$$
$$
\geq \frac{\kappa}{2\beta}(2-\beta)\mathbb{E}[\|\mathbf{z}^k-\mathbf{z}^{k+1}\|^2]
$$
$$
\geq \frac{\kappa}{2\beta}\mathbb{E}[\|\mathbf{z}^k-\mathbf{z}^{k+1}\|^2], \tag{88}
$$

where the last inequality is due to $0<\beta<1$. By combining (79), (87) and (88), we obtain (78). Besides, (76) and (77) implies (75). □

**Lemma 6.** *Under Assumptions 1-3, if $\kappa > -\mu$, $\gamma > 3L$, $\eta_k$ is a constant satisfies the condition (43), and*

$$
0 < \alpha \leq \min\left\{\frac{\hat{\sigma}_1}{4\sigma_A\sigma_1^2}, \frac{\hat{\sigma}_2}{4\sigma_A^2\sigma_2^2\eta_k^2}, c\right\}, \tag{89}
$$
$$
0 < \beta < \min\left\{\frac{\alpha}{12\kappa\sigma_5^2}, \frac{\sigma_4}{36}, 1\right\}, \tag{90}
$$

*where $\sigma_1$, $\sigma_2$, $\sigma_4$ and $\sigma_5$ are constants denoted in (66) and (67), (71) and (72), respectively. Then we have*

$$
\mathbb{E}[\phi^k - \phi^{k+1}]
$$
$$
\geq \frac{\kappa(1-\beta)\beta}{4}\mathbb{E}[\|\hat{\mathbf{x}}^{k+1}-\mathbf{z}^k\|^2] + \frac{\alpha}{2}\mathbb{E}[\|\mathbf{A}\mathbf{x}(\mathbf{z}^k,\boldsymbol{\lambda}^{k+1})\|^2]
$$
$$
+ \frac{\hat{\sigma}_1}{2}\mathbb{E}[\|\mathbf{x}^k-\hat{\mathbf{x}}^{k+1}\|^2] + \frac{\hat{\sigma}_2}{2}\mathbb{E}[\|\mathbf{x}^k-\mathbf{x}^{k-1}\|^2] - \frac{C_1N\sigma^2}{|\mathcal{I}|}, \tag{91}
$$

*where $\hat{\mathbf{x}}^{k+1}$ and $\hat{\mathbf{z}}^{k+1}$ are defined in (58) and (59), and*

$$
C_1 = \left(\frac{1}{2\gamma}+\frac{6L+8\tau+\kappa(1-\beta)}{4(\gamma+2c+\kappa)^2}\right). \tag{92}
$$

*Proof.* From the definition of $d(\mathbf{z};\boldsymbol{\lambda})$ in (54), we have

$$
\mathbb{E}[d(\mathbf{z}^k;\boldsymbol{\lambda}^{k+1}) - d(\mathbf{z}^k;\boldsymbol{\lambda}^k)]
$$
$$
= \mathbb{E}[L_c(\mathbf{x}(\mathbf{z}^k;\boldsymbol{\lambda}^{k+1}),\mathbf{z}^k;\boldsymbol{\lambda}^{k+1}) - L_c(\mathbf{x}(\mathbf{z}^k;\boldsymbol{\lambda}^k),\mathbf{z}^k;\boldsymbol{\lambda}^k)]
$$
$$
\geq \mathbb{E}[L_c(\mathbf{x}(\mathbf{z}^k;\boldsymbol{\lambda}^{k+1}),\mathbf{z}^k;\boldsymbol{\lambda}^{k+1}) - L_c(\mathbf{x}(\mathbf{z}^k;\boldsymbol{\lambda}^{k+1}),\mathbf{z}^k;\boldsymbol{\lambda}^k)]
$$
$$
= \alpha\mathbb{E}[\langle\mathbf{A}\mathbf{x}^k,\mathbf{A}\mathbf{x}(\mathbf{z}^k,\boldsymbol{\lambda}^{k+1})\rangle],
$$

where the inequality is due to $\mathbf{x}(\mathbf{z}^k;\boldsymbol{\lambda}^k) = \arg\min_{\mathbf{x}} L_c(\mathbf{x},\mathbf{z}^k;\boldsymbol{\lambda}^k)$ and the second equality comes from the iterates in (22). Using a similar technique, we have

$$
\mathbb{E}[d(\mathbf{z}^{k+1};\boldsymbol{\lambda}^{k+1}) - d(\mathbf{z}^k;\boldsymbol{\lambda}^{k+1})]
$$
$$
\geq \frac{\kappa}{2}\mathbb{E}[(\mathbf{z}^{k+1}-\mathbf{z}^k)^\top(\mathbf{z}^{k+1}+\mathbf{z}^k-2\mathbf{x}(\mathbf{z}^{k+1},\boldsymbol{\lambda}^{k+1}))].
$$

Combing the above two inequalities, we know

$$
\mathbb{E}[d(\mathbf{z}^{k+1};\boldsymbol{\lambda}^{k+1}) - d(\mathbf{z}^k;\boldsymbol{\lambda}^k)]
$$
$$
\geq \alpha\mathbb{E}[\langle\mathbf{A}\mathbf{x}^k,\mathbf{A}\mathbf{x}(\mathbf{z}^k,\boldsymbol{\lambda}^{k+1})\rangle] \tag{93}
$$
$$
+ \frac{\kappa}{2}\mathbb{E}[(\mathbf{z}^{k+1}-\mathbf{z}^k)^\top(\mathbf{z}^{k+1}+\mathbf{z}^k-2\mathbf{x}(\mathbf{z}^{k+1},\boldsymbol{\lambda}^{k+1}))].
$$

Based on Danskin's theorem [34, Proposition B.25] in convex analysis and $P(\mathbf{z})$ defined in (56) with $\kappa > -\mu$, we can have

$$
\nabla P(\mathbf{z}^k) = \kappa(\mathbf{z}^k-\mathbf{x}(\mathbf{z}^k)).
$$

Thus, it shows

$$\|\nabla P(\mathbf{z}^k) - \nabla P(\mathbf{z}^{k+1})\|$$
$$\leq \kappa\|\mathbf{z}^k - \mathbf{z}^{k+1}\| + \kappa\|\mathbf{x}(\mathbf{z}^{k+1}) - \mathbf{x}(\mathbf{z}^k)\|$$
$$\leq \kappa\tilde{\sigma}_4\|\mathbf{z}^{k+1} - \mathbf{z}^k\|,$$

where $\tilde{\sigma}_4 = 1 + \sigma_4^{-1}$ and the final inequality is due to Lemma 3. The above inequality shows the gradient of $P(\mathbf{z}^k)$ is Lipschitz continuous, which therefore it satisfies the descent lemma

$$\mathbb{E}[P(\mathbf{z}^{k+1}) - P(\mathbf{z}^k)] \qquad (94)$$
$$\leq \mathbb{E}[\kappa(\mathbf{z}^{k+1} - \mathbf{z}^k)^\top(\mathbf{z}^k - \mathbf{x}(\mathbf{z}^k))] + \frac{\kappa\tilde{\sigma}_4}{2}\mathbb{E}[\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2].$$

By combining (93), (94) and (78), we obtain

$$\mathbb{E}[\phi^k - \phi^{k+1}]$$
$$\geq -\alpha\mathbb{E}[\|\mathbf{A}\mathbf{x}^k\|^2] + \frac{\kappa}{2\beta}\mathbb{E}[\|\mathbf{z}^k - \mathbf{z}^{k+1}\|^2]$$
$$+ \hat{\sigma}_1\mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2] - 2\mathbb{E}[\kappa(\mathbf{z}^{k+1} - \mathbf{z}^k)^\top(\mathbf{z}^k - \mathbf{x}(\mathbf{z}^k))]$$
$$+ \hat{\sigma}_2\mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2] - \left(\frac{1}{2\gamma} + \frac{3L + 4\tau}{2(\gamma + 2c + \kappa)^2}\right)\frac{N\sigma^2}{|\mathcal{I}|}$$
$$- \kappa\tilde{\sigma}_4\mathbb{E}[\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2] + 2\alpha\mathbb{E}[\langle\mathbf{A}\mathbf{x}^k, \mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\rangle]$$
$$+ \kappa\mathbb{E}[(\mathbf{z}^{k+1} - \mathbf{z}^k)^\top(\mathbf{z}^{k+1} + \mathbf{z}^k - 2\mathbf{x}(\mathbf{z}^{k+1}, \boldsymbol{\lambda}^{k+1}))]$$
$$= \alpha\mathbb{E}[\|\mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2] - \alpha\mathbb{E}[\|\mathbf{A}(\mathbf{x}^k - \mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1}))\|^2]$$
$$+ \hat{\sigma}_1\mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2] + (\frac{\kappa}{2\beta} + \kappa - \kappa\tilde{\sigma}_4)\mathbb{E}[\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2]$$
$$+ 2\kappa\mathbb{E}[(\mathbf{z}^{k+1} - \mathbf{z}^k)^\top(\mathbf{x}(\mathbf{z}^k) - \mathbf{x}(\mathbf{z}^{k+1}; \boldsymbol{\lambda}^{k+1}))]$$
$$+ \hat{\sigma}_2\mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2] - \left(\frac{1}{2\gamma} + \frac{3L + 4\tau}{2(\gamma + 2c + \kappa)^2}\right)\frac{N\sigma^2}{|\mathcal{I}|}, \qquad (95)$$

where the equality comes from completing the square

$$\mathbb{E}[\|\mathbf{A}(\mathbf{x}^k - \mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1}))\|^2]$$
$$= \mathbb{E}[\|\mathbf{A}\mathbf{x}^k\|^2 - 2\langle\mathbf{A}\mathbf{x}^k, \mathbf{A}vx(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1}) + \|\mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2\rangle].$$

We further bound the right-hand-side terms of (95). By using the Young's inequality, we have

$$2(\mathbf{z}^{k+1} - \mathbf{z}^k)^\top(\mathbf{x}(\mathbf{z}^k) - \mathbf{x}(\mathbf{z}^k; \boldsymbol{\lambda}^{k+1}))$$
$$\geq -\frac{\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2}{6\beta} - 6\beta\|\mathbf{x}(\mathbf{z}^k) - \mathbf{x}(\mathbf{z}^k; \boldsymbol{\lambda}^{k+1})\|^2$$
$$\geq -\frac{\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2}{6\beta} - 6\beta\sigma_5^2\|\mathbf{A}\mathbf{x}(\mathbf{z}; \boldsymbol{\lambda})\|^2, \qquad (96)$$

where the lase inequality dues to (72) in Lemma 4. Besides, using the error bound (69) in Lemma 3, we have

$$(\mathbf{z}^{k+1} - \mathbf{z}^k)^\top(\mathbf{x}(\mathbf{z}^k; \boldsymbol{\lambda}^{k+1}) - \mathbf{x}(\mathbf{z}^{k+1}; \boldsymbol{\lambda}^{k+1}))$$
$$\geq -\|\mathbf{z}^{k+1} - \mathbf{z}^k\|\|\mathbf{x}(\mathbf{z}^k; \boldsymbol{\lambda}^{k+1}) - \mathbf{x}(\mathbf{z}^{k+1}; \boldsymbol{\lambda}^{k+1})\|$$
$$\geq -\frac{1}{\sigma_4}\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2. \qquad (97)$$

Also, based on the error bound (64) in Lemma 2, we obtain

$$\|\mathbf{A}(\mathbf{x}^k - \mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1}))\|^2$$
$$\leq 2\sigma_A^2\sigma_1^2\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2 + 2\sigma_A^2\sigma_2^2\|\mathbf{x}^k - \mathbf{s}^k\|^2. \qquad (98)$$

By substituting (96), (97) and (98) into (95), we therefore obtain

$$\mathbb{E}[\phi^k - \phi^{k+1}]$$
$$\geq (\alpha - 6\kappa\beta\sigma_5^2)\mathbb{E}[\|\mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2]$$
$$+ (\frac{\kappa}{2\beta} + \kappa - \kappa\tilde{\sigma}_5 - \frac{\kappa}{6\beta} - \frac{2\kappa}{\sigma_4})\mathbb{E}[\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2]$$
$$+ (\hat{\sigma}_1 - 2\alpha\sigma_A^2\sigma_1^2)\mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2]$$
$$+ (\hat{\sigma}_2 - 2\alpha\sigma_A^2\sigma_2^2\eta_k^2)\mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2]$$
$$- \left(\frac{1}{2\gamma} + \frac{3L + 4\tau}{2(\gamma + 2c + \kappa)^2}\right)\frac{N\sigma^2}{|\mathcal{I}|}. \qquad (99)$$

From (90), we know $\beta < \frac{\sigma_4}{36}$. By recalling $\tilde{\sigma}_4 = 1 + \sigma_4^{-1}$, we have

$$\frac{\kappa}{2\beta} + \kappa - \kappa\tilde{\sigma}_4 - \frac{\kappa}{6\beta} - \frac{2\kappa}{\sigma_4} \geq \frac{\kappa}{4\beta}.$$

As $\beta < \frac{\alpha}{12\kappa\sigma_5^2}$ by (90), we have

$$\alpha - 6\kappa\beta\sigma_5^2 \geq \frac{\alpha}{2}.$$

Similarly, based on (89), we have

$$\hat{\sigma}_1 - 2\alpha\sigma_A^2\sigma_1^2 \geq \frac{\hat{\sigma}_1}{2}, \quad \hat{\sigma}_2 - 2\alpha\sigma_A^2\sigma_2^2\eta_k^2 \geq \frac{\hat{\sigma}_2}{2}.$$

Thus, it follows from (99) that

$$\mathbb{E}[\phi^k - \phi^{k+1}]$$
$$\geq \frac{\kappa}{4\beta}\mathbb{E}[\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2] + \frac{\alpha}{2}\mathbb{E}[\|\mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2]$$
$$+ \frac{\hat{\sigma}_1}{2}\mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2] + \frac{\hat{\sigma}_2}{2}\mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2]$$
$$- \left(\frac{1}{2\gamma} + \frac{3L + 4\tau}{2(\gamma + 2c + \kappa)^2}\right)\frac{N\sigma^2}{|\mathcal{I}|}. \qquad (100)$$

Note that by using the definition of $\hat{\mathbf{z}}^{k+1}$ in (59) and by (25), we have

$$\hat{\mathbf{z}}^{k+1} = \mathbf{z}^{k+1} + \beta(\hat{\mathbf{x}}^{k+1} - \mathbf{x}^{k+1}). \qquad (101)$$

Thus, we can bound $\mathbb{E}[\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2]$ as

$$\mathbb{E}[\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2]$$
$$\geq (1 - \frac{1}{\beta})\mathbb{E}[\|\mathbf{z}^{k+1} - \hat{\mathbf{z}}^{k+1}\|^2] + (1 - \beta)\mathbb{E}[\|\hat{\mathbf{z}}^{k+1} - \mathbf{z}^k\|^2]$$
$$= \beta(\beta - 1)\mathbb{E}[\|\mathbf{x}^{k+1} - \hat{\mathbf{x}}^{k+1}\|^2] + (1 - \beta)\mathbb{E}[\|\hat{\mathbf{z}}^{k+1} - \mathbf{z}^k\|^2]$$
$$\geq \frac{\beta(\beta - 1)}{(\gamma + 2c + \kappa)^2}\frac{N\sigma^2}{|\mathcal{I}|} + (1 - \beta)\beta^2\mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathbf{z}^k\|^2],$$

where the last inequality comes from (61) and (59). By substituting the above inequality (100), we obtain (91). $\qquad\square$

## C. Proof of Theorem 1

We are ready to prove Theorem 1. By summing (91) for $k = 0, 1, \ldots, K-1$, we obtain

$$\mathbb{E}[\phi^0 - \phi^K]$$

$$\geq \frac{\kappa(1-\beta)\beta}{4} \sum_{k=0}^{K-1} \mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathbf{z}^k\|^2] - K\frac{C_1 N\sigma^2}{|\mathcal{I}|}$$

$$+ \frac{\alpha}{2} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2] + \frac{\hat{\sigma}_1}{2} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2]$$

$$+ \frac{\hat{\sigma}_2}{2} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2]. \tag{102}$$

Recall the definition of $Q(\mathbf{x}, \boldsymbol{\lambda})$ in (41)

$$Q(\mathbf{x}, \boldsymbol{\lambda}) = \|\mathbf{x} - \mathrm{prox}_r^1(\mathbf{x} - \nabla f(\mathbf{x}) - \mathbf{A}^\top\boldsymbol{\lambda})\|^2 + \|\mathbf{A}\mathbf{x}\|^2. \tag{103}$$

To obtain the desired result, we first consider

$$\mathbb{E}[\|\mathbf{x}^k - \mathrm{prox}_r^1(\mathbf{x}^k - \nabla_{\mathbf{x}} f(\mathbf{x}^k) - \mathbf{A}^\top\boldsymbol{\lambda}^{k+1})\|^2]$$

$$\leq 2\mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2] \tag{104}$$

$$+ 2\mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathrm{prox}_r^1(\mathbf{x}^k - \nabla_{\mathbf{x}} f(\mathbf{x}^k) - \mathbf{A}^\top\boldsymbol{\lambda}^{k+1})\|^2]$$

where the inequality dues to $\|a+b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$. Notice

$$\mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathrm{prox}_r^1(\mathbf{x}^k - \nabla_{\mathbf{x}} f(\mathbf{x}^k) - \mathbf{A}^\top\boldsymbol{\lambda}^{k+1})\|^2]$$

$$= \mathbb{E}[\|\mathrm{prox}_r^1(\hat{\mathbf{x}}^{k+1} - \nabla_{\mathbf{x}} g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1}))$$

$$- \mathrm{prox}_r^1(\mathbf{x}^k - \nabla_{\mathbf{x}} f(\mathbf{x}^k) - \mathbf{A}^\top\boldsymbol{\lambda}^{k+1})\|^2]$$

$$\leq \mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k - \nabla_{\mathbf{x}} g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1})$$

$$+ \nabla_{\mathbf{x}} f(\mathbf{x}^k) + \mathbf{A}^\top\boldsymbol{\lambda}^{k+1}\|^2]$$

$$\leq 2\mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k\|^2]$$

$$+ 2\mathbb{E}[\|\nabla_{\mathbf{x}} g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1}) - \nabla_{\mathbf{x}} f(\mathbf{x}^k) - \mathbf{A}^\top\boldsymbol{\lambda}^{k+1}\|^2],$$

$$= 2\mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k\|^2] + 2\mathbb{E}[\|\nabla_{\mathbf{x}} f(\mathbf{s}^k) - \nabla_{\mathbf{x}} f(\mathbf{x}^k)$$

$$+ \gamma(\hat{\mathbf{x}}^{k+1} - \mathbf{s}^k) + c\mathbf{D}(\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k) + c\mathbf{A}^T\mathbf{A}\mathbf{x}^k$$

$$+ \kappa(\hat{\mathbf{x}}^{k+1} - \mathbf{z}^k)\|^2]$$

$$\leq 2\mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k\|^2] + 10\mathbb{E}[\|\nabla_{\mathbf{x}} f(\mathbf{s}^k) - \nabla_{\mathbf{x}} f(\mathbf{x}^k)\|^2$$

$$+ \|\gamma(\hat{\mathbf{x}}^{k+1} - \mathbf{s}^k)\|^2 + \|c\mathbf{D}(\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k)\|^2 + \|c\mathbf{A}^T\mathbf{A}\mathbf{x}^k\|^2$$

$$+ \|\kappa(\hat{\mathbf{x}}^{k+1} - \mathbf{z}^k)\|^2]$$

$$\leq (2 + 10c^2 d_{max}^2 + 20\gamma^2)\mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathbf{x}^k\|^2] + 10c^2\sigma_A^2\mathbb{E}[\|\mathbf{A}\mathbf{x}^k\|^2]$$

$$+ (10L^2 + 20\gamma^2)\mathbb{E}[\|\mathbf{x}^k - \mathbf{s}^k\|^2] + 10\kappa^2\mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathbf{z}^k\|^2],$$

where $d_{\max}$ is the largest value of matrix $\mathbf{D}$, the first equality is due to the optimal condition for (58), i.e., $\hat{\mathbf{x}}^{k+1} = \mathrm{prox}_r^1(\hat{\mathbf{x}}^{k+1} - \nabla_{\mathbf{x}} g(\hat{\mathbf{x}}^{k+1}, \mathbf{w}^k; \boldsymbol{\lambda}^{k+1}))$; the first inequality is owing to the nonexpansive property of the proximal operator; the second inequality dues to $\|a+b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$; the second equality is obtained by the definition of function $g$ in (60); the last inequality dues to the $L$-smooth in (33).

Next, we show the upper bound for $\|\mathbf{A}\mathbf{x}^k\|$ as

$$\mathbb{E}[\|\mathbf{A}\mathbf{x}^k\|^2]$$

$$\leq 2\mathbb{E}[\|\mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2] + 2\sigma_A^2\mathbb{E}[\|\mathbf{x}^k - \mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2]$$

$$\leq 2\mathbb{E}[\|\mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2] + 4\sigma_A^2\sigma_1^2\mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2]$$

$$+ 4\sigma_A^2\sigma_2^2\mathbb{E}[\|\mathbf{x}^k - \mathbf{s}^k\|^2], \tag{105}$$

where the last inequality comes from Lemma 2. Now, we consider the upper bound of (103). Using the above inequalities (104)-(105), we can obtain

$$\min_{k=0,\ldots,K-1} \mathbb{E}[Q(\mathbf{x}^k, \boldsymbol{\lambda}^{k+1})]$$

$$\leq \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{x}^k - \mathrm{prox}_r^1(\mathbf{x}^k - \nabla_{\mathbf{x}} f(\mathbf{x}^k) - \mathbf{A}^\top\boldsymbol{\lambda}^{k+1})\|^2]$$

$$+ \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{A}\mathbf{x}^k\|^2]$$

$$\leq \frac{K_1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2] + \frac{K_2}{K} \sum_{k=0}^{K-1} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2]$$

$$+ \frac{K_3}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\hat{\mathbf{x}}^{k+1} - \mathbf{z}^k\|^2] + \frac{K_4}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2].$$

where

$$K_1 = 6 + 40\gamma + 20c^2 d_{\max}^2 + 4(20c^2\sigma_A^2 + 1)\sigma_A^2\sigma_1^2,$$

$$K_2 = (20L^2 + 20\gamma^2)\bar{\eta}^2 + 4(20c^2\sigma_A^2 + 1)\sigma_A^2\sigma_2^2\bar{\eta}^2,$$

$$K_3 = 20\kappa^2, \quad K_4 = 2(20c^2\sigma_A^2 + 1).$$

Further applying (102), we have

$$\min_{k=0,\ldots,K-1} \mathbb{E}[Q(\mathbf{x}^k, \boldsymbol{\lambda}^{k+1})] \leq C_0 \left( \frac{\mathbb{E}[\phi^0 - \phi^K]}{K} + \frac{C_1 N\sigma^2}{|\mathcal{I}|} \right)$$

$$\leq C_0 \left( \frac{\phi^0 - \underline{f}}{K} + \frac{C_1 N\sigma^2}{|\mathcal{I}|} \right),$$

where $\underline{f}$ is the lower bound of $\phi$ and $C_0$ is defined as follows,

$$C_0 \triangleq \frac{2K_1}{\hat{\sigma}_1} + \frac{K_2}{\hat{\sigma}_2} + \frac{4K_3\beta}{\kappa(1-\beta)} + \frac{2K_4}{\alpha}, \tag{106}$$

## APPENDIX C
## PROOF OF THEOREM 2

*Proof.* If we know the full gradient $\nabla f(\mathbf{x}^k)$, i.e., $G(\mathbf{x}^k, \boldsymbol{\xi}^k) = \nabla f(\mathbf{x}^k)$, then $\sigma^2 = 0$. Substituting it into (91), we have

$$\phi^k - \phi^{k+1}$$

$$\geq \frac{\kappa(1-\beta)\beta}{4} \|\hat{\mathbf{x}}^{k+1} - \mathbf{z}^k\|^2 + \frac{\alpha}{2} \|\mathbf{A}\mathbf{x}(\mathbf{z}^k, \boldsymbol{\lambda}^{k+1})\|^2$$

$$+ \frac{\hat{\sigma}_1}{2} \|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|^2 + \frac{\hat{\sigma}_2}{2} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 \geq 0.$$

Thus $\phi^k$ is monotonically decreasing and it has lower bound $\underline{f}$. This implies that

$$\max\{\|\mathbf{x}^k - \hat{\mathbf{x}}^{k+1}\|, \|\mathbf{z}^k - \hat{\mathbf{x}}^{k+1}\|, \|\mathbf{A}\mathbf{x}(\mathbf{z}^k; \boldsymbol{\lambda}^{k+1})\|\} \to 0.$$

Thus, according to [18, Theorem 2.4], every limit point generated by PPDM algorithm is a KKT point of problem (5). In addition, substituting $\sigma^2 = 0$ into (44) and picking $K \geq \frac{C_0(\phi^0 - \underline{f})}{\epsilon}$, we have

$$\min_{k=0,\ldots,K-1} Q(\mathbf{x}^k, \boldsymbol{\lambda}^{k+1}) \leq C_0 \left( \frac{\phi^0 - \underline{f}}{K} \right) \leq \epsilon.$$

Therefore, the proof is completed. $\qquad\square$

## References

[1] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.

[2] P. Mach and Z. Becvar, "Mobile edge computing: A survey on architecture and computation offloading," *IEEE Communications Surveys & Tutorials*, vol. 19, no. 3, pp. 1628–1656, 2017.

[3] S. Vlaski and A. H. Sayed, "Distributed learning in non-convex environments–Part I: Agreement at a linear rate," *arXiv preprint arXiv:1907.01848*, 2019.

[4] T.-H. Chang, M. Hong, and X. Wang, "Multi-agent distributed optimization via inexact consensus ADMM," *IEEE Transactions on Signal Processing*, vol. 63, no. 2, pp. 482–497, 2015.

[5] T.-H. Chang, M. Hong, H.-T. Wai, X. Zhang, and S. Lu, "Distributed learning in the non-convex world: From batch to streaming data, and beyond," *IEEE Signal Processing Magazine*, vol. 37, pp. 26–38, 2020.

[6] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.

[7] W. Shi, Q. Ling, G. Wu, and W. Yin, "EXTRA: An exact first-order algorithm for decentralized consensus optimization," *SIAM Journal on Optimization*, vol. 25, no. 2, pp. 944–966, 2015.

[8] W. Shi, Q. Ling, G. Wu, and W. Yin, "A proximal gradient algorithm for decentralized composite optimization," *IEEE Transactions on Signal Processing*, vol. 63, no. 22, pp. 6013–6023, 2015.

[9] 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*, vol. 62, no. 7, pp. 1750–1761, 2014.

[10] S. A. Alghunaim, E. K. Ryu, K. Yuan, and A. H. Sayed, "Decentralized proximal gradient algorithms with linear convergence rates," *arXiv preprint arXiv:1909.06479*, 2019.

[11] J. Zeng and W. Yin, "On nonconvex decentralized gradient descent," *IEEE Transactions on Signal Processing*, vol. 66, no. 11, pp. 2834–2848, 2018.

[12] P. Di Lorenzo and G. Scutari, "NEXT: In-network nonconvex optimization," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 2, pp. 120–136, 2016.

[13] G. Scutari and Y. Sun, "Parallel and distributed successive convex approximation methods for big-data optimization," in *Multi-agent Optimization*, pp. 141–308, Springer, 2018.

[14] G. Scutari and Y. Sun, "Distributed nonconvex constrained optimization over time-varying digraphs," *Mathematical Programming*, vol. 176, no. 1-2, pp. 497–544, 2019.

[15] M. Hong, D. Hajinezhad, and M.-M. Zhao, "Prox-PDA: The proximal primal-dual algorithm for fast distributed nonconvex optimization and learning over networks," in *Proceedings of the 34th International Conference on Machine Learning*, pp. 1529–1538, JMLR, 2017.

[16] H. Sun and M. Hong, "Distributed non-convex first-order optimization and information processing: Lower complexity bounds and rate optimal algorithms," *IEEE Transactions on Signal Processing*, vol. 67, no. 22, pp. 5912–5928, 2019.

[17] J. Zhang and Z.-Q. Luo, "A proximal alternating direction method of multiplier for linearly constrained nonconvex minimization," *Accepted for publication in SIAM Journal on Optimization*, 2020.

[18] J. Zhang and Z. Luo, "A global dual error bound and its application to the analysis of linearly constrained nonconvex optimization," *arXiv preprint arXiv:2006.16440*, 2020.

[19] M. Hong and T.-H. Chang, "Stochastic proximal gradient consensus over random networks," *IEEE Transactions on Signal Processing*, vol. 65, no. 11, pp. 2933–2948, 2017.

[20] X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu, "Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent," in *Advances in Neural Information Processing Systems*, pp. 5330–5340, 2017.

[21] H. Tang, X. Lian, M. Yan, C. Zhang, and J. Liu, "$D^2$: Decentralized training over decentralized data," *arXiv preprint arXiv:1803.07068*, 2018.

[22] S. Lu, X. Zhang, H. Sun, and M. Hong, "GNSD: A gradient-tracking based nonconvex stochastic algorithm for decentralized optimization," in *2019 IEEE Data Science Workshop, DSW 2019*, pp. 315–321, Institute of Electrical and Electronics Engineers Inc., 2019.

[23] P. Bianchi and J. Jakubowicz, "Convergence of a multi-agent projected stochastic gradient algorithm for non-convex optimization," *IEEE Transactions on Automatic Control*, vol. 58, no. 2, pp. 391–405, 2012.

[24] D. Jakovetić, J. Xavier, and J. M. Moura, "Fast distributed gradient methods," *IEEE Transactions on Automatic Control*, vol. 59, no. 5, pp. 1131–1146, 2014.

[25] H. Li, C. Fang, W. Yin, and Z. Lin, "A sharp convergence rate analysis for distributed accelerated gradient methods," *arXiv preprint arXiv:1810.01053*, 2018.

[26] Y. Yan, T. Yang, Z. Li, Q. Lin, and Y. Yang, "A unified analysis of stochastic momentum methods for deep learning," *arXiv preprint arXiv:1808.10396*, 2018.

[27] H. Yu, R. Jin, and S. Yang, "On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization," in *Proceedings of the 36th International Conference on Machine Learning*, 2019.

[28] F. Huang and S. Chen, "Mini-batch stochastic ADMMs for nonconvex nonsmooth optimization," *arXiv preprint arXiv:1802.03284*, 2018.

[29] Y. Xu, S. Zhu, S. Yang, C. Zhang, R. Jin, and T. Yang, "Learning with non-convex truncated losses by SGD," *arXiv preprint arXiv:1805.07880*, 2018.

[30] R. T. Rockafellar, *Convex Analysis*. Princeton University Press, Princeton, NJ, 1970.

[31] J.-S. Pang, "A posteriori error bounds for the linearly-constrained variational inequality problem," *Mathematics of Operations Research*, vol. 12, no. 3, pp. 474–484, 1987.

[32] Z. Wang, J. Zhang, T.-H. Chang, J. Li, and Z.-Q. Luo, "Supplementary material for distributed consensus optimization with momentum for nonconvex nonsmooth problems," *https://www.researchgate.net/publication/343418255*, 2020.

[33] Y. Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*. Kluwer Academic, Dordrecht, 2004.

[34] D. P. Bertsekas, *Nonlinear programming*. Athena Scientific, 1999.