

Accelerated Distributed Nesterov Gradient Descent for Smooth and Strongly Convex Functions

Guannan Qu, Na Li

Abstract—This paper considers the distributed optimization problem over a network, where the objective is to optimize a global function formed by a sum of local functions, using only local computation and communication. We develop an Accelerated Distributed Nesterov Gradient Descent (Acc-DNGD) method for strongly-convex and smooth functions. We show that it achieves a linear convergence rate and analyze how the convergence rate depends on the condition number and the underlying graph structure.

I. INTRODUCTION

Given a set of agents $\mathcal{N} = \{1, 2, \dots, n\}$, each of which has a local convex cost function $f_i(x) : \mathbb{R}^N \rightarrow \mathbb{R}$, the objective of distributed optimization is to find x that minimizes the average of all the functions,

$$\min_{x \in \mathbb{R}^N} f(x) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(x)$$

using local communication and local computation. The local communication is defined through an undirected communication graph $\mathcal{G} = (V, E)$, where the nodes $V = \mathcal{N}$ and edges $E \subset V \times V$. This problem has recently received much attention and has found various applications in multi-agent control, distributed state estimation over sensor networks, large scale computation in machine/statistical learning, etc [1]–[3].

There exist many studies on developing distributed algorithms for this problem, e.g., [4]–[15], most of which are distributed gradient descent algorithms based on a consensus scheme [16]. These methods have achieved sublinear convergence rates for convex functions. When the convex functions are nonsmooth, the sublinear convergence rate matches the centralized gradient method. More recent work [17]–[23], have improved these results to achieve linear convergence rates for strongly convex and smooth functions, which matches the centralized gradient method as well.

It is known that among all centralized gradient based algorithms, centralized Nesterov Gradient Descent (CNGD) [16] achieves the optimal convergence speed for smooth and convex functions in terms of first-order oracle complexity. Specifically, for L -smooth and convex problems, the convergence rate is $O(\frac{1}{t^2})$; for L -smooth and μ -strongly convex problems, the convergence rate is $O((1 - \sqrt{\mu\eta})^t)$ for step size $\eta \in (0, \frac{1}{L}]$. The nice convergence rates lead to the question of this paper: how to decentralize the Nesterov

Gradient methods to achieve similar convergence rates? [24] has developed Distributed Nesterov Gradient (D-NG) method and shown that for convex and L -smooth problems, it has a convergence rate of $O(\frac{\log t}{t})$.¹

In this paper, we focus on the μ -strongly convex and L -smooth case. We propose an Accelerated Distributed Nesterov Gradient Descent (Acc-DNGD) method. Different from [24], our algorithm features an additional gradient estimation scheme. The gradient estimation scheme firstly appears in [25] and independently in [17], and it uses a consensus scheme to estimate the global gradient. It has been generalized to distributed optimization with time-varying graphs [19], directed graphs [22], and with uncoordinated step sizes [20]. The gradient estimation scheme takes advantage of the smoothness of the functions. As a result, our Acc-DNGD method achieves similar convergence rate $O((1 - \sqrt{\mu\eta})^t)$ as CNGD but for a more restricted range of step size $\eta \in (0, \frac{1}{L}(\frac{1-\sigma}{64})^3(\frac{\mu}{L})^{3/2})$ where σ is a parameter related to the size and topology of the communication graph.

The rest of the paper is organized as follows. Section II formally defines the problem and presents our algorithm and results. Section III proves the convergence of our algorithm. Lastly, Section IV provides numerical simulations and Section V concludes the paper.

Notations. In this paper, n is the number of agents, and N is the dimension of the domain of the f_i 's. $i, j \in \{1, 2, \dots, n\}$ are indexes for agents, while $t, k, \ell \in \mathbb{N}$ are indexes for iteration steps. We use x^* and f^* to denote the minimizer and the minimal value of f , respectively. If f has multiple minimizers, x^* can be any of them. $\|\cdot\|$ denotes 2-norm for vectors, and Frobenius norm for matrices. $\langle \cdot, \cdot \rangle$ denotes inner product for vectors and matrices. $\rho(\cdot)$ denotes spectral radius for square matrices, and $\mathbf{1}$ denotes a n -dimensional all one column vector. All vectors, when having dimension N (the dimension of the domain of the f_i 's), will all be regarded as row vectors. As a special case, all gradients, $\nabla f_i(x)$ and $\nabla f(x)$ are interpreted as N -dimensional row vectors. ' \leq ', when applied to vectors of the same dimension, denotes element wise 'less than or equal to'.

II. PROBLEM AND ALGORITHM

A. Problem Formulation

Consider n agents, $\mathcal{N} = \{1, 2, \dots, n\}$, each of which has a convex function $f_i : \mathbb{R}^N \rightarrow \mathbb{R}$. The objective of distributed

¹The convergence rate for strongly convex and smooth functions is not studied in [24]. From our simulation results in Section IV, the convergence rate for strongly convex and smooth functions is sublinear.

Guannan Qu and Na Li are affiliated with John A. Paulson School of Engineering and Applied Sciences at Harvard University. Email: gqu@g.harvard.edu, nali@seas.harvard.edu. This work is supported under NSF ECCS 1608509 and NSF CAREER 1553407.

optimization is to find x to minimize the average of all the functions, i.e.

$$\min_{x \in \mathbb{R}^N} f(x) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(x) \quad (1)$$

using local communication and local computation. The local communication is defined through a *connected undirected* communication graph $\mathcal{G} = (V, E)$, where the nodes $V = \mathcal{N}$ and edges $E \subset V \times V$. Agent i and j can send information to each other if and only if $(i, j) \in E$. The local computation means that each agent can only make his decision based on the local function f_i and the information obtained from his neighbors.

Throughout the paper, we assume that the set of minimizers of f is non-empty and compact. We denote x^* as one of the minimizers and f^* as the minimal value. We will study the case where each f_i is μ -strongly convex (Assumption 1) and L -smooth (Assumption 2).

Assumption 1. $\forall i \in \mathcal{N}$, f_i is μ -strongly convex, i.e. $\forall x, y \in \mathbb{R}^N$, we have

$$f_i(y) \geq f_i(x) + \langle \nabla f_i(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2.$$

As a result, f is also μ -strongly convex.

Assumption 2. $\forall i \in \mathcal{N}$, f_i is L -smooth, that is, f_i is differentiable and the gradient is L -Lipschitz continuous, i.e., $\forall x, y \in \mathbb{R}^N$,

$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|.$$

As a result, f is L -smooth.

B. Centralized Nesterov Gradient Descent (CNGD)

We here briefly introduce a version of centralized Nesterov Gradient Descent (CNGD) algorithm that is derived from [16]. It is designated for μ -strongly convex and L -smooth functions. Given η , let $\alpha = \sqrt{\mu\eta}$, CNGD keeps updating three variables $x(t), v(t), y(t) \in \mathbb{R}^N$, starting from an initial point $x(0) = v(0) = y(0) \in \mathbb{R}^N$, and the update equation is given by

$$x(t+1) = x(t) - \eta \nabla f(y(t)) \quad (2a)$$

$$v(t+1) = (1 - \alpha)v(t) + \alpha y(t) - \frac{\alpha}{\mu} \nabla f(y(t)) \quad (2b)$$

$$y(t+1) = \frac{x(t+1) + \alpha v(t+1)}{1 + \alpha}. \quad (2c)$$

The following theorem (adapted from [16]) gives the convergence rate of CNGD.

Theorem 1. *Under Assumption 1 and Assumption 2, when $0 < \eta \leq \frac{1}{L}$, in CNGD we have $f(x(t)) - f^* = O((1 - \sqrt{\mu\eta})^t)$.*

We note here that the version of CNGD shown in (2) has a different form compared to a more common version, where only two variables need to be updated (see, e.g., [26]). In fact, the two versions are equivalent, and their relationship can be found in [16].

C. Our Algorithm: Accelerated Distributed Nesterov Gradient Descent (Acc-DNGD)

We design our algorithm based on a consensus matrix $W = [w_{ij}] \in \mathbb{R}^{n \times n}$. Here w_{ij} stands for how much agent i weighs his neighbor j 's information. W satisfies the following properties:

- (a) $\forall (i, j) \in E$, $w_{ij} > 0$. $\forall i$, $w_{ii} > 0$. $w_{ij} = 0$ elsewhere.
- (b) W is doubly stochastic, i.e. $\sum_{i'} w_{i'j} = \sum_{j'} w_{ij'} = 1$ for all $i, j \in \mathcal{N}$.

As a result, $\exists \sigma \in (0, 1)$ which depends on the spectrum of W , such that for any $\omega \in \mathbb{R}^{n \times 1}$, we have the ‘averaging property’, $\|W\omega - \mathbf{1}\bar{\omega}\| \leq \sigma\|\omega - \mathbf{1}\bar{\omega}\|$ where $\bar{\omega} = \frac{1}{n}\mathbf{1}^T\omega$ (the average of the entries in ω) [27]. The selection of the consensus matrix to satisfy these properties has been intensely studied, see [27], [28]. We will use the ‘averaging’ property of matrix W frequently in the rest of the paper.

In our algorithm Acc-DNGD, each agent keeps a copy of the three variables in CNGD, $x_i(t), v_i(t), y_i(t)$ and in addition $s_i(t)$ which serves as a gradient estimator. The initial condition is $x_i(0) = v_i(0) = y_i(0) = 0$ and $s_i(0) = \nabla f(0)$,² and it updates as follows:

$$x_i(t+1) = \sum_j w_{ij} y_j(t) - \eta s_i(t) \quad (3a)$$

$$v_i(t+1) = (1 - \alpha) \sum_j w_{ij} v_j(t) + \alpha \sum_j w_{ij} y_j(t) - \frac{\alpha}{\mu} s_i(t) \quad (3b)$$

$$y_i(t+1) = \frac{x_i(t+1) + \alpha v_i(t+1)}{1 + \alpha} \quad (3c)$$

$$s_i(t+1) = \sum_j w_{ij} s_j(t) + \nabla f_i(y_i(t+1)) - \nabla f_i(y_i(t)) \quad (3d)$$

where $[w_{ij}]_{n \times n}$ are the consensus weights and $\eta > 0$ is a fixed step size and $\alpha = \sqrt{\mu\eta}$. Because $w_{ij} = 0$ when $(i, j) \notin E$, each node i only needs to send $x_i(t), v_i(t), y_i(t)$ and $s_i(t)$ to its neighbors. Therefore, the algorithm can be operated in a fully distributed fashion with only local communication. The additional term $s_i(t)$ allows each agent to obtain an estimate on the global gradient $\frac{1}{n} \sum_i f_i(y_i(t))$. Compared with distributed algorithms without this estimation term, it helps improve the convergence speed. As a result, we call this method as Accelerated Distributed Nesterov Gradient Descent (Acc-DNGD) method.

D. Convergence of the Algorithm

To state the convergence results, we need to define the following average sequence,

$$\bar{x}(t) = \frac{1}{n} \sum_{i=1}^n x_i(t) \in \mathbb{R}^{1 \times N}.$$

²We note that the initial condition $s_i(0) = \nabla f(0)$ requires the agents to conduct an initial run of consensus. We impose this initial condition for technical reasons, while we expect the result of this paper to hold for a relaxed initial condition, $s_i(0) = \nabla f_i(0)$ which does not need initial coordination.

We summarize our convergence results here.

Theorem 2. *Under the smooth and strongly convex assumptions (Assumption 1 and Assumption 2), when $0 < \eta \leq \frac{1}{L}(\frac{1-\sigma}{64})^3(\frac{\mu}{L})^{3/2}$, we have $f(\bar{x}(t)) - f^* = O((1 - \sqrt{\mu\eta})^t)$.*

III. CONVERGENCE ANALYSIS

In this section, we will provide the proof of Theorem 2. We will first provide a proof overview in Section III-A and then defer the detailed proof to the rest of the section.

A. Proof Overview

We first introduce matrix notations $x(t), v(t), y(t), s(t), \nabla(t) \in \mathbb{R}^{n \times N}$ to simplify the mathematical expressions,³

$$x(t) = [x_1(t)^T, x_2(t)^T, \dots, x_n(t)^T]^T$$

$$v(t) = [v_1(t)^T, v_2(t)^T, \dots, v_n(t)^T]^T$$

$$y(t) = [y_1(t)^T, y_2(t)^T, \dots, y_n(t)^T]^T$$

$$s(t) = [s_1(t)^T, s_2(t)^T, \dots, s_n(t)^T]^T$$

$$\nabla(t) = [\nabla f_1(y_1(t))^T, \nabla f_2(y_2(t))^T, \dots, \nabla f_n(y_n(t))^T]^T.$$

Now our algorithm in (3) can be written as

$$x(t+1) = Wy(t) - \eta s(t) \quad (4a)$$

$$v(t+1) = (1 - \alpha)Wv(t) + \alpha Wy(t) - \frac{\alpha}{\mu} s(t) \quad (4b)$$

$$y(t+1) = \frac{x(t+1) + \alpha v(t+1)}{1 + \alpha} \quad (4c)$$

$$s(t+1) = Ws(t) + \nabla(t+1) - \nabla(t). \quad (4d)$$

Apart from the average sequence $\bar{x}(t) = \frac{1}{n} \sum_{i=1}^n x_i(t)$ that we have defined, we also define several other average sequences, $\bar{v}(t) = \frac{1}{n} \sum_{i=1}^n v_i(t)$, $\bar{y}(t) = \frac{1}{n} \sum_{i=1}^n y_i(t)$, $\bar{s}(t) = \frac{1}{n} \sum_{i=1}^n s_i(t)$, and $g(t) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(y_i(t))$.

Overview of the Proof. In our proof, we firstly derive the update formula for the average sequences (Lemma 3). Then, we show that the update rule for the average sequences is in fact centralized Nesterov Gradient Descent (CNGD) with inexact gradients [29], and the inexactness is characterized by ‘‘consensus error’’ $\|y(t) - \mathbf{1}\bar{y}(t)\|$ (Lemma 4). The consensus error is bounded in Lemma 5. Lastly, we prove the convergence of the algorithm by applying the proof of CNGD (see e.g. [16]) to the average sequences in spite of the consensus error (Lemma 6).

Lemma 3. *The following equalities hold.*

$$\bar{x}(t+1) = \bar{y}(t) - \eta g(t) \quad (5a)$$

$$\bar{v}(t+1) = (1 - \alpha)\bar{v}(t) + \alpha\bar{y}(t) - \frac{\alpha}{\mu}g(t) \quad (5b)$$

$$\bar{y}(t+1) = \frac{\bar{x}(t+1) + \alpha\bar{v}(t+1)}{1 + \alpha} \quad (5c)$$

$$\bar{s}(t+1) = \bar{s}(t) + g(t+1) - g(t) = g(t+1) \quad (5d)$$

³Without causing any confusion with notations in (2), in this section we abuse the use of notation $x(t), v(t), y(t)$.

Proof: We omit the proof since these can be easily derived using the fact that W is doubly stochastic. For (5d) we also need to use the fact that $\bar{s}(0) = g(0)$. \square

From (5a)-(5c) we see that the sequences $\bar{x}(t)$, $\bar{v}(t)$ and $\bar{y}(t)$ follow a update rule similar to the CNGD in (2). The only difference is that the $g(t)$ in (5a)-(5c) is not the exact gradient $\nabla f(\bar{y}(t))$ in CNGD. In the following Lemma, we show that $g(t)$ is an inexact gradient with error $O(\|y(t) - \mathbf{1}\bar{y}(t)\|^2)$.

Lemma 4. *$\forall t$, $g(t)$ is an inexact gradient of f at $\bar{y}(t)$ with error $O(\|y(t) - \mathbf{1}\bar{y}(t)\|^2)$ in the sense that $\exists \hat{f}(t) \in \mathbb{R}$ s.t. $\forall \omega \in \mathbb{R}^N$,*

$$f(\omega) \geq \hat{f}(t) + \langle g(t), \omega - \bar{y}(t) \rangle + \frac{\mu}{2} \|\omega - \bar{y}(t)\|^2 \quad (6)$$

$$f(\omega) \leq \hat{f}(t) + \langle g(t), \omega - \bar{y}(t) \rangle + L \|\omega - \bar{y}(t)\|^2 + L \frac{1}{n} \|y(t) - \mathbf{1}\bar{y}(t)\|^2. \quad (7)$$

Proof: Define

$$\hat{f}(t) = \frac{1}{n} \sum_{i=1}^n [f_i(y_i(t)) + \langle \nabla f_i(y_i(t)), \bar{y}(t) - y_i(t) \rangle].$$

Then, for any $\omega \in \mathbb{R}^N$, we have

$$\begin{aligned} f(\omega) &= \frac{1}{n} \sum_{i=1}^n f_i(\omega) \\ &\geq \frac{1}{n} \sum_{i=1}^n [f_i(y_i(t)) + \langle \nabla f_i(y_i(t)), \omega - y_i(t) \rangle + \frac{\mu}{2} \|\omega - y_i(t)\|^2] \\ &= \frac{1}{n} \sum_{i=1}^n [f_i(y_i(t)) + \langle \nabla f_i(y_i(t)), \bar{y}(t) - y_i(t) \rangle + \frac{1}{n} \sum_{i=1}^n \langle \nabla f_i(y_i(t)), \omega - \bar{y}(t) \rangle + \frac{1}{n} \sum_{i=1}^n \frac{\mu}{2} \|\omega - y_i(t)\|^2] \\ &\geq \hat{f}(t) + \langle g(t), \omega - \bar{y}(t) \rangle + \frac{\mu}{2} \|\omega - \bar{y}(t)\|^2 \end{aligned}$$

which shows (6). For (7), similarly,

$$\begin{aligned} f(\omega) &\leq \hat{f}(t) + \langle g(t), \omega - \bar{y}(t) \rangle + \frac{L}{2} \frac{1}{n} \sum_{i=1}^n \|(\omega - \bar{y}(t)) + (\bar{y}(t) - y_i(t))\|^2 \\ &\leq \hat{f}(t) + \langle g(t), \omega - \bar{y}(t) \rangle + L \|\omega - \bar{y}(t)\|^2 + L \frac{1}{n} \sum_{i=1}^n \|\bar{y}(t) - y_i(t)\|^2 \\ &= \hat{f}(t) + \langle g(t), \omega - \bar{y}(t) \rangle + L \|\omega - \bar{y}(t)\|^2 + L \frac{1}{n} \|y(t) - \mathbf{1}\bar{y}(t)\|^2 \end{aligned}$$

where in the second inequality we have used the elementary fact that $\|u + v\|^2 \leq 2\|u\|^2 + 2\|v\|^2$ for all $u, v \in \mathbb{R}^N$. \square

The consensus error in the previous lemma is bounded by the following lemma whose proof is deferred to Section III-B.

Lemma 5. When $0 < \eta \leq \frac{1}{L}(\frac{1-\sigma}{64})^3(\frac{\mu}{L})^{3/2}$, we have

$$\|y(t) - \mathbf{1}\bar{y}(t)\| \leq \frac{\sqrt{n}}{4} \sqrt{\frac{\mu}{L}} \|\bar{y}(t) - \bar{x}(t)\| + \frac{2\eta\sqrt{n}}{1-\sigma} \sqrt{\frac{\mu}{L}} \|g(t)\|.$$

Given the bounded consensus error, the average sequence satisfies the following convergence rate, whose proof is deferred to Section III-C.

Lemma 6. When $0 < \eta \leq \frac{(1-\sigma)^2}{18L}$ and when

$$\|y(t) - \mathbf{1}\bar{y}(t)\| \leq \frac{\sqrt{n}}{4} \sqrt{\frac{\mu}{L}} \|\bar{y}(t) - \bar{x}(t)\| + \frac{2\eta\sqrt{n}}{1-\sigma} \sqrt{\frac{\mu}{L}} \|g(t)\|$$

we have $f(\bar{x}(t)) - f^* = O((1 - \sqrt{\mu\eta})^t)$.

B. Proof of the Bounded Consensus Error (Lemma 5)

Proof of Lemma 5:

Overview of the proof. The proof is separated into three steps. In step 1, we first treat the algorithm (4) as a linear system and derive a linear system inequality (8). In step 2, we analyze the state transition matrix in (8) and prove several properties (e.g. bounded spectral radius). In step 3, we further analyze the linear system (8) and bound the state by the input, from which the conclusion of the lemma follows.

Step 1: A Linear System Inequality. Define state $z(t) = [\|v(t) - \mathbf{1}\bar{v}(t)\|, \|y(t) - \mathbf{1}\bar{y}(t)\|, \|s(t) - \mathbf{1}g(t)\|]^T \in \mathbb{R}^3$ and we will derive a linear system inequality that bounds the update of $z(t)$, given as follows

$$z(t+1) \leq G(\eta)z(t) + b(t). \quad (8)$$

Here $b(t) = [0, 0, \sqrt{n}a(t)]^T \in \mathbb{R}^3$ is the input to the system with

$$a(t) \triangleq \frac{1-\alpha}{1+\alpha} L \|\bar{y}(t) - \bar{x}(t)\| + \frac{2\lambda\eta L}{1+\alpha} \|g(t)\|$$

and $\lambda \triangleq \frac{4}{1-\sigma} > 1$. The state transition matrix $G(\eta) \in \mathbb{R}^{3 \times 3}$ is given by

$$G(\eta) = \begin{bmatrix} (1-\alpha)\sigma & \alpha\sigma & \frac{\eta}{1+\alpha} \\ \frac{1-\alpha}{1+\alpha}\alpha\sigma & \frac{1+\alpha^2}{1+\alpha}\sigma & \frac{2\eta}{1+\alpha} \\ \frac{1-\alpha}{1+\alpha}\alpha\sigma L & [\frac{1+\alpha^2}{1+\alpha}\sigma + 1]L & \sigma + \frac{2\eta L}{1+\alpha} \end{bmatrix}.$$

We now prove (8). By (4a) and (5a), we have

$$\begin{aligned} & \|x(t+1) - \mathbf{1}\bar{x}(t+1)\| \\ &= \|[Wy(t) - \mathbf{1}\bar{y}(t)] - \eta[s(t) - \mathbf{1}g(t)]\| \\ &\leq \sigma\|y(t) - \mathbf{1}\bar{y}(t)\| + \eta\|s(t) - \mathbf{1}g(t)\|. \end{aligned} \quad (9)$$

By (4b) and (5b), we have

$$\begin{aligned} & \|v(t+1) - \mathbf{1}\bar{v}(t+1)\| \\ &\leq \|(1-\alpha)[Wv(t) - \mathbf{1}\bar{v}(t)] + \alpha[Wy(t) - \mathbf{1}\bar{y}(t)] \\ &\quad - \frac{\alpha}{\mu}[s(t) - \mathbf{1}g(t)]\| \\ &\leq (1-\alpha)\sigma\|v(t) - \mathbf{1}\bar{v}(t)\| + \alpha\sigma\|y(t) - \mathbf{1}\bar{y}(t)\| \\ &\quad + \frac{\eta}{\alpha}\|s(t) - \mathbf{1}g(t)\| \end{aligned} \quad (10)$$

where in the last inequality we have used $\alpha = \sqrt{\mu\eta}$.

By (4c) and (5c), we have

$$\begin{aligned} & \|y(t+1) - \mathbf{1}\bar{y}(t+1)\| \\ &\leq \frac{1}{1+\alpha}\|x(t+1) - \mathbf{1}\bar{x}(t+1)\| \\ &\quad + \frac{\alpha}{1+\alpha}\|v(t+1) - \mathbf{1}\bar{v}(t+1)\| \\ &\leq \frac{1}{1+\alpha} \left[\sigma\|y(t) - \mathbf{1}\bar{y}(t)\| + \eta\|s(t) - \mathbf{1}g(t)\| \right] \\ &\quad + \frac{\alpha}{1+\alpha} \left[(1-\alpha)\sigma\|v(t) - \mathbf{1}\bar{v}(t)\| \right. \\ &\quad \left. + \alpha\sigma\|y(t) - \mathbf{1}\bar{y}(t)\| + \frac{\eta}{\alpha}\|s(t) - \mathbf{1}g(t)\| \right] \\ &\leq \frac{1-\alpha}{1+\alpha}\alpha\sigma\|v(t) - \mathbf{1}\bar{v}(t)\| + \frac{1+\alpha^2}{1+\alpha}\sigma\|y(t) - \mathbf{1}\bar{y}(t)\| \\ &\quad + \frac{2\eta}{1+\alpha}\|s(t) - \mathbf{1}g(t)\| \end{aligned} \quad (11)$$

where we have used (9) and (10) in the second inequality.

By (4d) and (5d), we have

$$\begin{aligned} & \|s(t+1) - \mathbf{1}g(t+1)\| \\ &= \|Ws(t) - \mathbf{1}g(t) \\ &\quad + [\nabla(t+1) - \nabla(t) - \mathbf{1}(g(t+1) - g(t))]\| \\ &\stackrel{(a)}{\leq} \sigma\|s(t) - \mathbf{1}g(t)\| + \|\nabla(t+1) - \nabla(t)\| \\ &\stackrel{(b)}{\leq} \sigma\|s(t) - \mathbf{1}g(t)\| + L\|y(t+1) - y(t)\| \end{aligned} \quad (12)$$

where in (a) we have used the fact that

$$\begin{aligned} & \left\| [\nabla(t+1) - \nabla(t)] - [\mathbf{1}g(t+1) - \mathbf{1}g(t)] \right\|^2 \\ &= \|\nabla(t+1) - \nabla(t)\|^2 - n\|g(t+1) - g(t)\|^2 \\ &\leq \|\nabla(t+1) - \nabla(t)\|^2 \end{aligned}$$

and in (b) we have used (24) (Proposition 10 in Appendix-A).

Now we expand $y(t+1) - y(t)$.

$$\begin{aligned} & \|y(t+1) - y(t)\| \\ &\leq \|y(t+1) - \mathbf{1}\bar{y}(t+1)\| + \|y(t) - \mathbf{1}\bar{y}(t)\| \\ &\quad + \|\mathbf{1}\bar{y}(t+1) - \mathbf{1}\bar{y}(t)\| \\ &\leq \frac{1-\alpha}{1+\alpha}\alpha\sigma\|v(t) - \mathbf{1}\bar{v}(t)\| + \left[\frac{1+\alpha^2}{1+\alpha}\sigma + 1 \right] \|y(t) - \mathbf{1}\bar{y}(t)\| \\ &\quad + \frac{2\eta}{1+\alpha}\|s(t) - \mathbf{1}g(t)\| + \|\mathbf{1}\bar{y}(t+1) - \mathbf{1}\bar{y}(t)\| \\ &\leq \frac{1-\alpha}{1+\alpha}\alpha\sigma\|v(t) - \mathbf{1}\bar{v}(t)\| + \left[\frac{1+\alpha^2}{1+\alpha}\sigma + 1 \right] \|y(t) - \mathbf{1}\bar{y}(t)\| \\ &\quad + \frac{2\eta}{1+\alpha}\|s(t) - \mathbf{1}g(t)\| + \sqrt{n}\frac{1-\alpha}{1+\alpha}\|\bar{y}(t) - \bar{x}(t)\| \\ &\quad + \sqrt{n}\frac{2\eta}{1+\alpha}\|g(t)\| \end{aligned}$$

where in the last inequality we have used (22) (Proposition 9 in Appendix-A). Therefore

$$\|s(t+1) - \mathbf{1}g(t+1)\|$$

$$\begin{aligned}
&\leq \sigma \|s(t) - \mathbf{1}g(t)\| + L \|y(t+1) - y(t)\| \\
&\leq \frac{1-\alpha}{1+\alpha} \alpha \sigma L \|v(t) - \mathbf{1}\bar{v}(t)\| + \left[\frac{1+\alpha^2}{1+\alpha} \sigma \right. \\
&\quad \left. + 1 \right] L \|y(t) - \mathbf{1}\bar{y}(t)\| + \left(\sigma + \frac{2\eta L}{1+\alpha} \right) \|s(t) - \mathbf{1}g(t)\| \\
&\quad + \sqrt{n} \frac{1-\alpha}{1+\alpha} L \|\bar{y}(t) - \bar{x}(t)\| + \sqrt{n} \frac{2\eta L}{1+\alpha} \|g(t)\|. \quad (13)
\end{aligned}$$

Combining (10) (11) (13) gives the linear system inequality (8). Note that we have intentionally loosened inequality (13) by multiplying a factor λ to the last term. The reason of doing so will be clear later.

Step 2: Spectral Properties of $G(\eta)$. We give the following two lemmas regarding the eigenvalues and eigenvectors of $G(\eta)$. Their proofs are provided in Appendix-B.

Lemma 7. *When $0 < \eta \leq \frac{1}{L}$, we have $\sigma < \rho(G(\eta)) < \sigma + 4(\eta L)^{1/3}$.*

Lemma 8. *When $0 < \eta \leq \frac{1}{L} \left(\frac{1-\sigma}{64}\right)^3 \left(\frac{\mu}{L}\right)^{3/2}$, $G(\eta)$ has a leading eigenvector $\chi = [\chi_1, \chi_2, \chi_3]^T$ with all positive entries, that satisfy $\chi_2 = 1$, $\chi_3 \geq \frac{16L}{1-\sigma} \sqrt{\frac{L}{\mu}}$, and $\chi_1 \leq D\chi_3$ for $D = \sqrt{2\left[\frac{1}{L^2} + \frac{1}{\sigma\mu L}\right]}$.*

In light of the above two lemmas, we define $\theta \triangleq \rho(G(\eta))$ and we have $\sigma < \theta < \sigma + 4(\eta L)^{1/3} < \frac{1+\sigma}{2}$. We also define $\kappa \triangleq \frac{2}{(1-\theta)\chi_3} \leq \frac{2}{1-\frac{1+\sigma}{2}} \frac{1}{\chi_3} \leq \frac{1}{4L} \sqrt{\frac{\mu}{L}}$.

Step 3: Bound $z(t)$ by $a(t)$. With the above preparations, now we prove, by induction, the following statement,

$$z(t) \leq \kappa \sqrt{na(t)} \chi \quad (14a)$$

$$a(t+1) \geq \frac{1+\theta}{2} a(t). \quad (14b)$$

For $t=0$, (14a) is automatically satisfied since $z(0)$ has all zero entries. For (14b), notice that $\bar{x}(0) = \bar{y}(0) = 0$, we have $a(0) = \frac{2\lambda\eta L}{1+\alpha} \|g(0)\|$, and

$$\begin{aligned}
a(1) &= \frac{1-\alpha}{1+\alpha} L \|\bar{y}(1) - \bar{x}(1)\| + \frac{2\lambda\eta L}{1+\alpha} \|g(1)\| \\
&\geq \frac{2\lambda\eta L}{1+\alpha} \|g(1)\| \\
&\geq \frac{2\lambda\eta L}{1+\alpha} \|g(0)\| - \frac{2\lambda\eta L}{1+\alpha} \|g(1) - g(0)\|.
\end{aligned}$$

Since $s(0) = \mathbf{1}g(0)$, it's easy to check that $x(1) = \mathbf{1}\bar{x}(1)$, $v(1) = \mathbf{1}\bar{v}(1)$ and $y(1) = \mathbf{1}\bar{y}(1)$. We then have

$$\begin{aligned}
\|g(1) - g(0)\| &= \|\nabla f(\bar{y}(1)) - \nabla f(\bar{y}(0))\| \\
&\leq L \|\bar{y}(1) - \bar{y}(0)\| \\
&= L \left\| \frac{1-\alpha}{1+\alpha} [\bar{y}(0) - \bar{x}(0)] - \frac{2\eta}{1+\alpha} g(0) \right\| \\
&= \frac{2\eta L}{1+\alpha} \|g(0)\|.
\end{aligned}$$

Hence

$$\begin{aligned}
a(1) &\geq \frac{2\lambda\eta L}{1+\alpha} \|g(0)\| - \frac{2\lambda\eta L}{1+\alpha} \frac{2\eta L}{1+\alpha} \|g(0)\| \\
&= \left(1 - \frac{2\eta L}{1+\alpha}\right) a(0).
\end{aligned}$$

Hence (14b) is true for $t=0$ as long as $\frac{2\eta L}{1+\alpha} < \frac{1-\theta}{2}$, which is true as long as $2\eta L < \frac{1-\frac{1+\sigma}{2}}{2} = \frac{1-\sigma}{4}$, and this is met by our step size selection.

Assume (14) holds for t . Then for $t+1$, we have

$$\begin{aligned}
z(t+1) &\stackrel{(a)}{\leq} G(\eta)z(t) + b(t) \\
&\stackrel{(b)}{\leq} G(\eta)\kappa\sqrt{na(t)}\chi + \sqrt{na(t)}\frac{1}{\chi_3}\chi \\
&\stackrel{(c)}{=} \theta\kappa\sqrt{na(t)}\chi + \sqrt{na(t)}\frac{1}{\chi_3}\chi \\
&= \kappa\sqrt{na(t)}\chi\left(\theta + \frac{1}{\chi_3\kappa}\right) \\
&\stackrel{(d)}{\leq} \kappa\sqrt{na(t+1)}\chi\left(\theta + \frac{1}{\chi_3\kappa}\right)\frac{2}{1+\theta} \\
&\stackrel{(e)}{=} \kappa\sqrt{na(t+1)}\chi \quad (15)
\end{aligned}$$

where (a) is due to (8), and (b) is due to induction assumption (14a), and (c) is because θ is an eigenvalue of $G(\eta)$ with eigenvector χ , and (d) is due to induction assumption (14b), and in (e), we have used the fact that by definition, κ satisfies $\left(\theta + \frac{1}{\chi_3\kappa}\right)\frac{2}{1+\theta} = 1$.

Now, (14a) is proven for $t+1$. For (14b), using (23) we have

$$\begin{aligned}
a(t+2) &= \frac{1-\alpha}{1+\alpha} L \|\bar{y}(t+2) - \bar{x}(t+2)\| + \frac{2\lambda\eta L}{1+\alpha} \|g(t+2)\| \\
&= \frac{1-\alpha}{1+\alpha} L \left\| \frac{1-\alpha}{1+\alpha} (\bar{y}(t+1) - \bar{x}(t+1)) \right. \\
&\quad \left. - \frac{1-\alpha}{1+\alpha} \eta g(t+1) \right\| + \frac{2\lambda\eta L}{1+\alpha} \|g(t+2)\| \\
&\geq \left[\frac{1-\alpha}{1+\alpha} \right]^2 L \|\bar{y}(t+1) - \bar{x}(t+1)\| \\
&\quad - \left[\frac{1-\alpha}{1+\alpha} \right]^2 \eta L \|g(t+1)\| \\
&\quad + \frac{2\lambda\eta L}{1+\alpha} \|g(t+1)\| - \frac{2\lambda\eta L}{1+\alpha} \|g(t+2) - g(t+1)\|.
\end{aligned}$$

Therefore, $a(t+1) = \frac{1-\alpha}{1+\alpha} L \|\bar{y}(t+1) - \bar{x}(t+1)\| + \frac{2\lambda\eta L}{1+\alpha} \|g(t+1)\|$ implies that

$$\begin{aligned}
a(t+1) - a(t+2) &\leq \left[\frac{1-\alpha}{1+\alpha} \right] \frac{2\alpha}{1+\alpha} L \|\bar{y}(t+1) - \bar{x}(t+1)\| \\
&\quad + \left[\frac{1-\alpha}{1+\alpha} \right]^2 \eta L \|g(t+1)\| + \frac{2\lambda\eta L}{1+\alpha} \|g(t+2) - g(t+1)\| \\
&\leq \max\left(\frac{2\alpha}{1+\alpha}, \frac{(1-\alpha)^2}{2\lambda(1+\alpha)}\right) a(t+1) \\
&\quad + \frac{2\lambda\eta L}{1+\alpha} \|g(t+2) - g(t+1)\|. \quad (16)
\end{aligned}$$

Here in the second inequality we have used the elementary fact that for four positive numbers a_1, a_2, a_3, a_4 , and nonnegative numbers u, v , we have $a_1u + a_2v = \frac{a_1}{a_3}a_3u + \frac{a_2}{a_4}a_4v \leq \max\left(\frac{a_1}{a_3}, \frac{a_2}{a_4}\right)(a_3u + a_4v)$.

Further, we expand $\|g(t+2) - g(t+1)\|$,

$$\|g(t+2) - g(t+1)\|$$

$$\begin{aligned}
&\leq \|g(t+2) - \nabla f(\bar{y}(t+2))\| + \|g(t+1) - \nabla f(\bar{y}(t+1))\| \\
&\quad + \|\nabla f(\bar{y}(t+2)) - \nabla f(\bar{y}(t+1))\| \\
&\stackrel{(a)}{\leq} \frac{L}{\sqrt{n}} \|y(t+2) - \mathbf{1}\bar{y}(t+2)\| + \frac{L}{\sqrt{n}} \|y(t+1) - \mathbf{1}\bar{y}(t+1)\| \\
&\quad + L\|\bar{y}(t+2) - \bar{y}(t+1)\| \\
&\stackrel{(b)}{\leq} \frac{L}{\sqrt{n}} \frac{1-\alpha}{1+\alpha} \alpha\sigma \|v(t+1) - \mathbf{1}\bar{v}(t+1)\| \\
&\quad + \frac{L}{\sqrt{n}} \left[\frac{1+\alpha^2}{1+\alpha} \sigma + 1 \right] \|y(t+1) - \mathbf{1}\bar{y}(t+1)\| \\
&\quad + \frac{L}{\sqrt{n}} \frac{2\eta}{1+\alpha} \|s(t+1) - \mathbf{1}g(t+1)\| + a(t+1) \\
&\stackrel{(c)}{\leq} L \frac{1-\alpha}{1+\alpha} \alpha\sigma\kappa\chi_1 a(t+1) + L \left[\frac{1+\alpha^2}{1+\alpha} \sigma + 1 \right] \kappa\chi_2 a(t+1) \\
&\quad + L \frac{2\eta}{1+\alpha} \kappa\chi_3 a(t+1) + a(t+1) \\
&\stackrel{(d)}{\leq} a(t+1) \left\{ L \frac{1-\alpha}{1+\alpha} \alpha\sigma \frac{4D}{1-\sigma} + \left[\frac{1+\alpha^2}{1+\alpha} \sigma + 1 \right] \right. \\
&\quad \left. + L \frac{2\eta}{1+\alpha} \frac{4}{1-\sigma} + 1 \right\} \\
&\leq a(t+1) \left\{ L \frac{4D\sigma}{1-\sigma} + \frac{8}{1-\sigma} + 3 \right\} \stackrel{(e)}{\leq} Ca(t+1)
\end{aligned}$$

where C is a constant defined as $C \triangleq 20 \frac{1}{1-\sigma} \sqrt{\frac{L}{\mu}}$. Here (a) is due to (25), (b) is due to (11) and (22), (c) is due to the fact that (14a) is true for $t+1$, (d) is due to the fact that $\chi_2\kappa = \kappa \leq \frac{1}{L}$, $\chi_3\kappa = \frac{2}{1-\theta} < \frac{2}{1-\frac{1+\sigma}{2}} = \frac{4}{1-\sigma}$ and $\chi_1 \leq D\chi_3$, and (e) is due to the fact that $4D\sigma L < 8\sqrt{\frac{L}{\mu}}$.

Combining the above expansion for $\|g(t+2) - g(t+1)\|$ with (16), we have

$$\begin{aligned}
&a(t+1) - a(t+2) \\
&\leq \max\left(\frac{2\alpha}{1+\alpha}, \frac{(1-\alpha)^2}{2\lambda(1+\alpha)}\right) a(t+1) + \frac{2\lambda\eta L}{1+\alpha} Ca(t+1) \\
&\leq \left[\max(2\sqrt{\mu\eta}, \frac{1-\sigma}{8}) + \frac{8}{1-\sigma} \eta LC \right] a(t+1).
\end{aligned}$$

Therefore, since $\sqrt{\mu\eta} < \sqrt{\eta L} < \frac{1-\sigma}{16}$, and $\frac{8}{1-\sigma} \eta LC < \frac{1-\sigma}{8}$, we have $a(t+1) - a(t+2) \leq \frac{1-\sigma}{4} a(t+1) \leq \frac{1-\theta}{2} a(t+1)$, where we have used $\theta < \frac{1+\sigma}{2}$. Hence $a(t+2) \geq \frac{1+\theta}{2} a(t+1)$. Now we have completed the induction. Therefore, (14) is true for all t , and in particular, we have $\|y(t) - \mathbf{1}\bar{y}(t)\| \leq \kappa\sqrt{n}a(t)$. Combining this with $\kappa \leq \frac{1}{4L} \sqrt{\frac{\mu}{L}}$ gives the lemma. \square

C. Proof of Convergence of the Average Sequence (Lemma 6)

Proof of Lemma 6: The proof essentially follows the same argument as the CNGD (see e.g. Section 2.2 in [16]). Firstly we define a series of functions $\Phi_t(\omega)$ recursively, with $\Phi_0(\omega) = f(\bar{x}(0)) + \frac{\mu}{2} \|\omega - \bar{v}(0)\|^2$ and

$$\begin{aligned}
\Phi_{t+1}(\omega) &= (1-\alpha)\Phi_t(\omega) \\
&\quad + \alpha(\hat{f}(t) + \langle g(t), \omega - \bar{y}(t) \rangle + \frac{\mu}{2} \|\omega - \bar{y}(t)\|^2). \quad (17)
\end{aligned}$$

It's clear that Φ_t is always a quadratic function. Since $\nabla^2\Phi_0(\omega) = \mu I$ and $\nabla^2\Phi_{t+1}(\omega) = (1-\alpha)\nabla^2\Phi_t(\omega) + \alpha\mu I$, we get $\nabla^2\Phi_t(\omega) = \mu I$ for all t . We now claim that $\Phi_t(\omega)$ can be written as

$$\Phi_t(\omega) = \phi_t^* + \frac{\mu}{2} \|\omega - \bar{v}(t)\|^2 \quad (18)$$

for some $\phi_t^* \in \mathbb{R}$. In other words, we are claiming that Φ_t achieves its minimum at $\bar{v}(t)$. We prove the claim using induction. Firstly, Φ_0 achieves its minimum at $\bar{v}(0)$. Next, by (17) and the induction assumption, we have

$$\nabla\Phi_{t+1}(\omega) = (1-\alpha)\mu(\omega - \bar{v}(t)) + \alpha(g(t) + \mu(\omega - \bar{y}(t))).$$

Set $\omega = \bar{v}(t+1)$, we can easily verify $\nabla\Phi_{t+1}(\bar{v}(t+1)) = 0$, and hence Φ_{t+1} achieves its optimum at $\bar{v}(t+1)$, and hence (18) is true.

Notice that $\phi_0^* = f(\bar{x}(0))$. We now derive a recursive formula for ϕ_t^* . Since by (17), $\Phi_{t+1}(\bar{y}(t)) = (1-\alpha)\Phi_t(\bar{y}(t)) + \alpha\hat{f}(t)$. Combining this with (18), we get

$$\begin{aligned}
\phi_{t+1}^* &= (1-\alpha)\phi_t^* + \frac{\mu}{2}(1-\alpha)\alpha\|\bar{v}(t) - \bar{y}(t)\|^2 + \alpha\hat{f}(t) \\
&\quad - \frac{1}{2}\eta\|g(t)\|^2 + (1-\alpha)\alpha\langle g(t), \bar{v}(t) - \bar{y}(t) \rangle. \quad (19)
\end{aligned}$$

We now claim that

$$\phi_t^* \geq f(\bar{x}(t)). \quad (20)$$

This is true for $t=0$, since $\phi_0^* = f(\bar{x}(0))$. Suppose it's true for t . Then for $t+1$, we get, by (19),

$$\begin{aligned}
&\phi_{t+1}^* \\
&\stackrel{(a)}{\geq} (1-\alpha)f(\bar{x}(t)) + \frac{\mu}{2}(1-\alpha)\alpha\|\bar{v}(t) - \bar{y}(t)\|^2 \\
&\quad + \alpha\hat{f}(t) - \frac{1}{2}\eta\|g(t)\|^2 + (1-\alpha)\alpha\langle g(t), \bar{v}(t) - \bar{y}(t) \rangle \\
&\stackrel{(b)}{\geq} (1-\alpha)\{ \hat{f}(t) + \langle g(t), \bar{x}(t) - \bar{y}(t) \rangle + \frac{\mu}{2}\|\bar{x}(t) - \bar{y}(t)\|^2 \} \\
&\quad + \alpha\hat{f}(t) - \frac{1}{2}\eta\|g(t)\|^2 + (1-\alpha)\alpha\langle g(t), \bar{v}(t) - \bar{y}(t) \rangle \\
&= \hat{f}(t) - \frac{1}{2}\eta\|g(t)\|^2 + (1-\alpha)\frac{\mu}{2}\|\bar{x}(t) - \bar{y}(t)\|^2 \\
&\quad + \langle g(t), (1-\alpha)\alpha(\bar{v}(t) - \bar{y}(t)) + (1-\alpha)(\bar{x}(t) - \bar{y}(t)) \rangle
\end{aligned}$$

where (a) is due to the induction assumption, and (b) is due to (6). But because $(1-\alpha)\alpha(\bar{v}(t) - \bar{y}(t)) + (1-\alpha)(\bar{x}(t) - \bar{y}(t)) = (1-\alpha)[\alpha\bar{v}(t) + \bar{x}(t) - (1+\alpha)\bar{y}(t)] = 0$, we have

$$\phi_{t+1}^* \geq \hat{f}(t) - \frac{1}{2}\eta\|g(t)\|^2 + (1-\alpha)\frac{\mu}{2}\|\bar{x}(t) - \bar{y}(t)\|^2.$$

Therefore, we only need to prove

$$f(\bar{x}(t+1)) \leq \hat{f}(t) - \frac{1}{2}\eta\|g(t)\|^2 + (1-\alpha)\frac{\mu}{2}\|\bar{x}(t) - \bar{y}(t)\|^2.$$

We now use the bound on $\|y(t) - \mathbf{1}\bar{y}(t)\|$ and get

$$\|y(t) - \mathbf{1}\bar{y}(t)\|^2 \leq n \frac{\mu}{8L} \|\bar{y}(t) - \bar{x}(t)\|^2 + \frac{8\eta^2 n}{(1-\sigma)^2} \frac{\mu}{L} \|g(t)\|^2. \quad (21)$$

Now, by (7),

$$\begin{aligned}
& f(\bar{x}(t+1)) \\
& \leq \hat{f}(t) + \langle g(t), \bar{x}(t+1) - \bar{y}(t) \rangle \\
& \quad + L \|\bar{x}(t+1) - \bar{y}(t)\|^2 + L \frac{1}{n} \|y(t) - \mathbf{1}\bar{y}(t)\|^2 \\
& \leq \hat{f}(t) - \eta \|g(t)\|^2 + L \eta^2 \|g(t)\|^2 + \frac{8\eta^2}{(1-\sigma)^2} \mu \|g(t)\|^2 \\
& \quad + \frac{\mu}{8} \|\bar{x}(t) - \bar{y}(t)\|^2 \\
& = \hat{f}(t) - \frac{1}{2} \eta \|g(t)\|^2 + \left(-\frac{1}{2} \eta + \left(L + \frac{8}{(1-\sigma)^2} \mu\right) \eta^2\right) \|g(t)\|^2 \\
& \quad + \frac{1}{4} \frac{\mu}{2} \|\bar{x}(t) - \bar{y}(t)\|^2
\end{aligned}$$

where in the second inequality we have used (21). The induction is complete if $-\frac{1}{2}\eta + (L + \frac{8}{(1-\sigma)^2}\mu)\eta^2 < 0$, and $\frac{1}{4} < 1 - \alpha = 1 - \sqrt{\mu\eta}$, which are satisfied by our step size selection.

Lastly we show that $\Phi_t(\omega) \leq f(\omega) + (1 - \alpha)^t(\Phi_0(\omega) - f(\omega))$. This is true for $t = 0$. Then, assuming it's true for t , then for $t + 1$, we have by (17) and (6),

$$\begin{aligned}
\Phi_{t+1}(\omega) &= (1 - \alpha)\Phi_t(\omega) + \alpha(\hat{f}(t) \\
& \quad + \langle g(t), \omega - \bar{y}(t) \rangle + \frac{\mu}{2} \|\omega - \bar{y}(t)\|^2) \\
& \leq (1 - \alpha)\Phi_t(\omega) + \alpha f(\omega) \\
& \leq f(\omega) + (1 - \alpha)^{t+1}(\Phi_0(\omega) - f(\omega)).
\end{aligned}$$

As a result, we have

$$f(\bar{x}(t)) \leq \phi_t^* \leq \Phi_t(x^*) \leq f(x^*) + (1 - \alpha)^t(\Phi_0(x^*) - f(x^*)).$$

This shows that $f(\bar{x}(t)) - f^* = O((1 - \alpha)^t)$. \square

IV. NUMERICAL EXPERIMENTS

We test our algorithm on a 100-node connected graph in which each pair of nodes are connected with probability 0.3. The objective functions are randomly generated quadratic functions with $\mu = 1$ and $L = 100$. We compare our proposed algorithm Acc-DNGD in (3) with the CNGD in (2) and the distributed Nesterov Gradient algorithm (D-NG) developed in [24]. Simulation results are shown in Figure 1.

V. CONCLUSION

In this paper we propose an Accelerated Distributed Nesterov Gradient Descent algorithm for distributed optimization of strongly convex and smooth functions. In the future we are interested in studying the case where the functions are smooth and convex but possibly not strongly convex. We expect the method will achieve the convergence rate of $O(\frac{1}{t^2})$ under certain conditions.

REFERENCES

- [1] B. Johansson, "On distributed optimization in networked systems," 2008.
- [2] J. A. Bazerque and G. B. Giannakis, "Distributed spectrum sensing for cognitive radio networks by exploiting sparsity," *IEEE Transactions on Signal Processing*, vol. 58, no. 3, pp. 1847–1862, 2010.

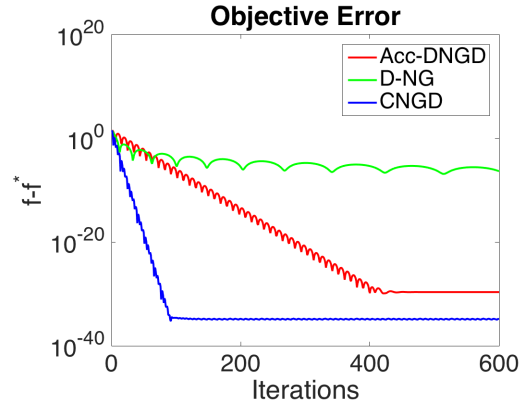


Fig. 1: Simulation results. Blue is CNGD (2) with step size $\eta = \frac{1}{L} = 0.01$. Green is the D-NG in [24] with $c = \frac{1}{L} = 0.01$. Red is the proposed algorithm (3) with step size $\eta = 0.002$.

- [3] P. A. Forero, A. Cano, and G. B. Giannakis, "Consensus-based distributed support vector machines," *Journal of Machine Learning Research*, vol. 11, no. May, pp. 1663–1707, 2010.
- [4] J. N. Tsitsiklis, D. P. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," in *1984 American Control Conference*, 1984, pp. 484–489.
- [5] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and distributed computation: numerical methods*. Prentice hall Englewood Cliffs, NJ, 1989, vol. 23.
- [6] A. Nedić and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *Automatic Control, IEEE Transactions on*, vol. 54, no. 1, pp. 48–61, 2009.
- [7] I. Lobel and A. Ozdaglar, "Convergence analysis of distributed subgradient methods over random networks," in *Communication, Control, and Computing, 2008 46th Annual Allerton Conference on*. IEEE, 2008, pp. 353–360.
- [8] J. C. Duchi, A. Agarwal, and M. J. Wainwright, "Dual averaging for distributed optimization: convergence analysis and network scaling," *Automatic control, IEEE Transactions on*, vol. 57, no. 3, pp. 592–606, 2012.
- [9] S. S. Ram, A. Nedić, and V. V. Veeravalli, "Distributed stochastic subgradient projection algorithms for convex optimization," *Journal of optimization theory and applications*, vol. 147, no. 3, pp. 516–545, 2010.
- [10] A. Nedic and A. Olshevsky, "Stochastic gradient-push for strongly convex functions on time-varying directed graphs," *arXiv preprint arXiv:1406.2075*, 2014.
- [11] —, "Distributed optimization over time-varying directed graphs," *Automatic Control, IEEE Transactions on*, vol. 60, no. 3, pp. 601–615, 2015.
- [12] I. Matei and J. S. Baras, "Performance evaluation of the consensus-based distributed subgradient method under random communication topologies," *Selected Topics in Signal Processing, IEEE Journal of*, vol. 5, no. 4, pp. 754–771, 2011.
- [13] A. Olshevsky, "Linear time average consensus on fixed graphs and implications for decentralized optimization and multi-agent control," *arXiv preprint arXiv:1411.4186*, 2014.
- [14] M. Zhu and S. Martínez, "On distributed convex optimization under inequality and equality constraints," *Automatic Control, IEEE Transactions on*, vol. 57, no. 1, pp. 151–164, 2012.
- [15] I. Lobel, A. Ozdaglar, and D. Feijer, "Distributed multi-agent optimization with state-dependent communication," *Mathematical Programming*, vol. 129, no. 2, pp. 255–284, 2011.
- [16] Y. Nesterov, *Introductory lectures on convex optimization: A basic course*. Springer Science & Business Media, 2013, vol. 87.
- [17] G. Qu and N. Li, "Harnessing smoothness to accelerate distributed optimization," *arXiv preprint arXiv:1605.07112*, 2016.
- [18] 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.

- [19] A. Nedich, A. Olshevsky, and W. Shi, “Achieving geometric convergence for distributed optimization over time-varying graphs,” *arXiv preprint arXiv:1607.03218*, 2016.
- [20] A. Nedic, A. Olshevsky, W. Shi, and C. A. Uribe, “Geometrically convergent distributed optimization with uncoordinated step-sizes,” *arXiv preprint arXiv:1609.05877*, 2016.
- [21] C. Xi and U. A. Khan, “On the linear convergence of distributed optimization over directed graphs,” *arXiv preprint arXiv:1510.02149*, 2015.
- [22] —, “Add-opt: Accelerated distributed directed optimization,” *arXiv preprint arXiv:1607.04757*, 2016.
- [23] J. Zeng and W. Yin, “Extrapush for convex smooth decentralized optimization over directed networks,” *arXiv preprint arXiv:1511.02942*, 2015.
- [24] D. Jakovetic, J. Xavier, and J. M. Moura, “Fast distributed gradient methods,” *Automatic Control, IEEE Transactions on*, vol. 59, no. 5, pp. 1131–1146, 2014.
- [25] J. Xu, S. Zhu, Y. C. Soh, and L. Xie, “Augmented distributed gradient methods for multi-agent optimization under uncoordinated constant stepsizes,” in *2015 54th IEEE Conference on Decision and Control (CDC)*. IEEE, 2015, pp. 2055–2060.
- [26] S. Bubeck *et al.*, “Convex optimization: Algorithms and complexity,” *Foundations and Trends® in Machine Learning*, vol. 8, no. 3-4, pp. 231–357, 2015.
- [27] A. Olshevsky and J. N. Tsitsiklis, “Convergence speed in distributed consensus and averaging,” *SIAM Journal on Control and Optimization*, vol. 48, no. 1, pp. 33–55, 2009.
- [28] R. Olfati-Saber, J. A. Fax, and R. M. Murray, “Consensus and cooperation in networked multi-agent systems,” *Proceedings of the IEEE*, vol. 95, no. 1, pp. 215–233, Jan 2007.
- [29] O. Devolder, F. Glineur, and Y. Nesterov, “First-order methods of smooth convex optimization with inexact oracle,” *Mathematical Programming*, vol. 146, no. 1-2, pp. 37–75, 2014.
- [30] C. D. Meyer, *Matrix analysis and applied linear algebra*. Siam, 2000, vol. 2.

APPENDIX

A. Frequently used propositions

In Section III, we frequently use the following propositions to prove the convergence of the algorithm.

Proposition 9. *The following equalities are true*

$$\bar{y}(t+1) - \bar{y}(t) = \frac{1-\alpha}{1+\alpha}[\bar{y}(t) - \bar{x}(t)] - \frac{2\eta}{1+\alpha}g(t) \quad (22)$$

$$\bar{y}(t+1) - \bar{x}(t+1) = \frac{1-\alpha}{1+\alpha}(\bar{y}(t) - \bar{x}(t)) - \frac{1-\alpha}{1+\alpha}\eta g(t) \quad (23)$$

Proof: We omit the proof since it can be easily derived from (5). \square

Proposition 10. *The following inequalities are true.*

$$\|\nabla(t+1) - \nabla(t)\| \leq L\|y(t+1) - y(t)\| \quad (24)$$

$$\|g(t) - \nabla f(\bar{y}(t))\| \leq \frac{L}{\sqrt{n}}\|y(t) - \mathbf{1}\bar{y}(t)\| \quad (25)$$

Proof: We omit the proof and refer the readers to the Lemma 6 of [17]. \square

B. Proof of Lemma 7 and Lemma 8

Proof of Lemma 7: We first calculate the characteristic polynomial of $G(\eta)$ as

$$p(\zeta) = (\zeta - \sigma)\left(\zeta - \frac{1-\alpha}{1+\alpha}\sigma\right)\left(\zeta - \sigma - \frac{2\eta L}{1+\alpha}\right) - k_1\left(\zeta - \sigma - \frac{2\eta L}{1+\alpha}\right) - k_2$$

where k_1 and k_2 are positive constants given by

$$k_1 = \eta L \left\{ \frac{2}{1+\alpha} \left[\frac{1+\alpha^2}{1+\alpha} \sigma + 1 \right] + \frac{1-\alpha}{1+\alpha} \sigma \right\} < 5\eta L$$

$$k_2 = \frac{\eta^2 L^2}{(1+\alpha)^2} \left\{ 4 \left[\frac{1+\alpha^2}{1+\alpha} \sigma + 1 \right] + 2(1-\alpha)\sigma \right\} + \frac{\eta L}{1+\alpha} \left\{ \left[\frac{1+\alpha^2}{1+\alpha} \sigma + 1 \right] (1+\alpha)\sigma + \sigma^2 \alpha \frac{1-\alpha^2}{1+\alpha} \right\} < 15\eta L.$$

Let $\zeta_0 = \sigma + 4(\eta L)^{1/3}$, then $\zeta_0 > \sigma + \frac{2\eta L}{1+\alpha}$. Since $p(\sigma + \frac{2\eta L}{1+\alpha}) < 0$, and $p(\zeta)$ is a strictly increasing function on $[\zeta_0, +\infty)$, and also since $G(\eta)$'s largest eigenvalue in magnitude must be a real number (Perron-Frobenius Theorem [30]), the lemma will follow from $p(\zeta_0) > 0$, which is shown below.

$$\begin{aligned} p(\zeta_0) &> (16(\eta L)^{2/3} - k_1)\left(\zeta_0 - \sigma - \frac{2\eta L}{1+\alpha}\right) - k_2 \\ &> [16(\eta L)^{2/3} - 5(\eta L)^{2/3}][4(\eta L)^{1/3} - 2(\eta L)^{1/3}] \\ &\quad - 15\eta L \\ &= 7\eta L > 0 \end{aligned}$$

\square

Proof of Lemma 8: Since η is positive, $G(\eta)$ is a positive matrix, by Perron-Frobenius Theorem [30] we have $G(\eta)$ has a unique largest (in magnitude) eigenvalue, and it's a real number $\theta = \rho(G(\eta))$, and θ has a eigenvector $\chi = [\chi_1, \chi_2, \chi_3]^T$ with all positive entries. Notice $\theta\chi = G(\eta)\chi$. Now we compute the third row,

$$\begin{aligned} \theta\chi_3 &= \frac{1-\alpha}{1+\alpha}\alpha\sigma L\chi_1 + \left[\frac{1+\alpha^2}{1+\alpha}\sigma + 1\right]L\chi_2 + \left[\sigma + \frac{2\eta L}{1+\alpha}\right]\chi_3 \\ &\geq L\chi_2 + \sigma\chi_3 \end{aligned}$$

Hence, by Lemma 7 and the step size selection $\frac{\chi_3}{\chi_2} \geq \frac{L}{\theta - \sigma} > \frac{L}{4(\eta L)^{1/3}} \geq \frac{16L}{1-\sigma} \sqrt{\frac{L}{\mu}}$. Now, still consider the third row,

$$\theta\chi_3 \geq \frac{1-\alpha}{1+\alpha}\alpha\sigma L\chi_1 + \sigma\chi_3 \Rightarrow \frac{\chi_3}{\chi_1} \geq \frac{1-\alpha}{1+\alpha}\alpha\sigma L \frac{1}{\theta - \sigma}$$

Now compute the first row,

$$\theta\chi_1 = (1-\alpha)\sigma\chi_1 + \alpha\sigma\chi_2 + \frac{\eta}{\alpha}\chi_3 \leq \sigma\chi_1 + \alpha\sigma \frac{\theta - \sigma}{L}\chi_3 + \frac{\alpha}{\mu}\chi_3$$

where we have used the fact that $\chi_2 \leq \frac{\theta - \sigma}{L}\chi_3$. This implies

$$\frac{\chi_3}{\chi_1} \geq \frac{\theta - \sigma}{\alpha[\sigma \frac{\theta - \sigma}{L} + \frac{1}{\mu}]}$$

Combine the above two lower bounds on $\frac{\chi_3}{\chi_1}$, and use the elementary fact that for two positive numbers p and q , $\max(p, q) \geq \sqrt{pq}$, we get

$$\begin{aligned} \frac{\chi_3}{\chi_1} &\geq \sqrt{\frac{1-\alpha}{1+\alpha}\alpha\sigma L \frac{1}{\theta - \sigma} \cdot \frac{\theta - \sigma}{\alpha[\sigma \frac{\theta - \sigma}{L} + \frac{1}{\mu}]}} \\ &= \sqrt{\frac{1-\alpha}{1+\alpha}\sigma L \frac{1}{\sigma \frac{\theta - \sigma}{L} + \frac{1}{\mu}}} \geq \frac{1}{D}. \end{aligned}$$

The lemma follows by letting $\chi_2 = 1$.