## FAST CONVERGENCE RATES OF DISTRIBUTED SUBGRADIENT METHODS WITH ADAPTIVE QUANTIZATION

THINH T. DOAN<sup>†,\*</sup>, SIVA THEJA MAGULURI<sup>\*</sup>, JUSTIN ROMBERG<sup>†</sup> <sup>†</sup> SCHOOL OF ELECTRICAL AND COMPUTER ENGINEERING <sup>\*</sup> H. MILTON STEWART SCHOOL OF INDUSTRIAL AND SYSTEMS ENGINEERING GEORGIA INSTITUTE OF TECHNOLOGY, GA, 30332, USA {THINHDOAN, SIVA.THEJA}@GATECH.EDU, JROM@ECE.GATECH.EDU.

**Abstract.** We study distributed optimization problems over a network when the communication between the nodes is constrained, and so information that is exchanged between the nodes must be quantized. Recent advances using the distributed gradient algorithm with a quantization scheme at a fixed resolution have established convergence, but at rates significantly slower than when the communications are unquantized.

In this paper, we introduce a novel quantization method, which we refer to as adaptive quantization, that allows us to match the convergence rates under perfect communications. Our approach adjusts the quantization scheme used by each node as the algorithm progresses: as we approach the solution, we become more certain about where the state variables are localized, and adapt the quantizer codebook accordingly.

We bound the convergence rates of the proposed method as a function of the communication bandwidth, the underlying network topology, and structural properties of the constituent objective functions. In particular, we show that if the objective functions are convex or strongly convex, then using adaptive quantization does not affect the rate of convergence of the distributed subgradient methods when the communications are quantized, except for a constant that depends on the resolution of the quantizer. To the best of our knowledge, the rates achieved in this paper are better than any existing work in the literature for distributed gradient methods under finite communication bandwidths. We also provide numerical simulations that compare convergence properties of the distributed gradient methods with and without quantization for solving distributed regression problems for both quadratic and absolute loss functions.

**1. Introduction.** We consider a distributed subgradient algorithm for solving optimization problems of the form

$$\underset{\mathbf{x}\in\mathcal{X}}{\text{minimize }} f(\mathbf{x}) \triangleq \sum_{i=1}^{n} f_i(\mathbf{x}).$$
(1.1)

Each functional  $f_i$  in the sum above is associated with a computational node, and these nodes have connected into a network specified by a graph. Each node has knowledge of only their function  $f_i$ , and so they must work together, communicating only with their neighbors on the graph, to find a minimizer of (1.1). The distributed subgradient (DSG) method, described in full in Section 2.2 below, is a popular approach for solving this type of problem. In DSG, each node keeps a local estimate of the decision variable **x** and iterates by communicating the local state to its neighbors on the graph, averaging the estimates it received from its neighbors, then taking a gradient step. The convergence properties of this algorithm are well understood, see [1] along with the other references in Section 1.2, and essentially match the rates of standard centralized subgradient methods with a constant factor that depends on the connectivity of the graph.

In this paper, we take a step towards understanding how imperfect communication affects the convergence of DSG. In particular, we derive convergence rates for a modified version of the DSG algorithm when the communications between the nodes are *quantized*, modeling scenarios in which the bandwidth available for communication is often limited. Unlike previous works [2, 3], we show if the quantization intervals are adjusted as the algorithm approaches a solution, then we can match the convergence rates of unquantized DSG within a constant that depends on the number of bits used for quantization. We call this novel coding method *adaptive quantization*, as it changes at every iteration based on the stepsize being used for the subgradient descent.

1.1. Main Contribution. We first propose a modified version of DSG, which directly takes into account the quantized error at every iteration. Second, we design a novel adaptive quantization method, where the nodes quantize their values based on the progress of their updates. The entire method is summarized in Algorithm 2.1 and described fully in Section 2 below.

Our analytical contribution is to show that the convergence rates of DSG are unaffected when the communications use adaptive quantization, except for a factor which captures the size of communication bandwidth.

Our analysis treats both the cases where the objective functions are convex and strongly convex. When the  $f_i$  are convex, we show that at iteration k, each node i has an estimate  $\mathbf{z}_i(k)$  that obeys

$$f(\mathbf{z}_i(k)) - f^* \lesssim \frac{\Delta^2}{(1 - \sigma_2)^2} \cdot \frac{\ln k}{\sqrt{k}},$$

where  $1 - \sigma_2$  is the spectral gap that quantifies the connectivity of the underlying network and  $\Delta$  is the resolution of the quantizer. When the  $f_i$  are strongly convex, we derive a rate on the convergence of the  $\mathbf{z}_i$  to the unique solution  $\mathbf{x}^*$ ,

$$\|\mathbf{z}_i(k) - \mathbf{x}^*\|^2 \lesssim \frac{\Delta^2}{(1 - \sigma_2)^2} \cdot \frac{\ln k}{k}.$$

These rates match those for the standard, unquantized verion of DSG [4].

The numerical results in Section 4 show that for stylized problems, both smooth and not smooth, quantizing to 8 bits is essentially the same as communicating real numbers, while using 5 or 6 bits results in only a modest increase in the number iterations required to converge to a specified tolerance.

1.2. Related Work. The DSG algorithms for solving problem (1.1) have a long history, probably first studied in [5] and recently received a wide attention; see for example, [1,6–11] and the recent survey paper [4]. Convergence results of DSG have been explicitly studied in this literature; however, they are mostly established under a critical assumption on the perfect communication between nodes. Such assumption is not often held in practice, therefore, there is a necessity to study the performance of DSG under imperfect communication. In particular, our focus in this paper is to study the convergence rates of this method when information exchanged between the nodes are quantized, modeling the practical applications with finite communication bandwidth.

Distributed algorithms with random (dithered) quantization have been considered in [12] for solving network consensus problems, a special case of the problem considered in this paper. On the other hand, different variants of distributed gradient methods under quantized communication have been studied in [2,3,13–17]. In [13,14] the authors only show the convergence to a neighborhood around the optimal of the problem, while an exact convergence has been studied in [15, 16]; however, a condition on the growing communication bandwidth is assumed in the latter work. To remove such strong condition, the authors in [2,3] show the asymptotic convergence of DSG methods under random quantization using only finite bandwidth. In particular, in [3] the authors provide a rate in expectation  $O(1/k^{(1-\gamma)/2})$ , for some  $\gamma \in (0, 1)$ , when the problem objectives are smooth and strongly convex. On the other hand, in [2] we study distributed subgradient methods for nonsmooth problems and analyze their convergence rates by utilizing techniques from stochastic approximation approach. Specifically, such algorithms asymptotically converge to the optimal value in expectation at a rate  $\mathcal{O}(\ln(k) / k^{1/4})$  and  $\mathcal{O}(\ln(k) / k^{1/3})$  for convex and strongly convex functions, respectively. The rates established in these two papers, however, are sub-optimal and much slower than the ones in this paper as stated in Section 1.1.

The adaptive quantization studied in this paper seems to share some similarity with the so-called "zoom in" and "zoom out" quantization to study the stability of linear systems [18]. Moreover, this "zoom in" and "zoom out" quantization has also been applied in distributed optimization with finite bandwidths, where an asymptotic convergence to a solution has been derived in [17]. However, there is a lack of understanding how fast the algorithm converges, which is one of the main focus of this paper.

We also want to note some related work [19–21] and the references therein, in which distributed stochastic gradient with quantization under master/worker models is considered. Such models consider a special star graph communication structure, while consensus-based gradient methods are designed for any network topology. In general, these two approaches are fundamentally different, therefore, the results studied in master/worker models cannot be extended to cover the problem considered in this paper. Moreover, the quantized communication constraint studied in this paper is one example of imperfect exchange of information between nodes. Another example of imperfect exchange is latency in the communications. Convergence rates of DSG optimization methods in the presence of communication delays have been studied in [22–25].

Finally, there are some related methods based on primal-dual approach for solving problem (1.1), such as, the accelerated primal-dual methods [26, 27], the *alternating direction method of multipliers* (ADMM) [28–32], and the *distributed dual methods* (mirror descent/dual averaging) [33–35]. Our focus in this paper will be on DSG algorithms, as they are both simple and have convergence guarantees that are as strong or stronger than those for dual methods.

**1.3. Notation.** We introduce here a set of definitions and notation that is used throughout this paper. We use boldface to denote vectors in  $\mathbb{R}^d$  to distinguish them from scalars in  $\mathbb{R}$ . Given a collection of vectors  $\mathbf{x}_1, \ldots, \mathbf{x}_n$  in  $\mathbb{R}^d$ , we denote by  $\mathbf{X}$  a matrix in  $\mathbb{R}^{n \times d}$ , whose *i*-th row is  $\mathbf{x}_i^T$ . We then denote by  $\|\mathbf{x}\|$  and  $\|\mathbf{X}\|$  the Euclidean norm and the Frobenius norm of the vector  $\mathbf{x}$  and the meatrix  $\mathbf{X}$ , respectively. We use  $\mathbf{1} \in \mathbb{R}^d$  to denote the the vector whose entries are all 1 and  $\mathbf{I} \in \mathbb{R}^{n \times n}$  to denote the identity matrix. Given a closed convex set  $\mathcal{X}$ , we denote by  $[\mathbf{x}]_{\mathcal{X}}$  the projection of  $\mathbf{x}$  to  $\mathcal{X}$ .

Given a nonsmooth convex function  $f : \mathbb{R}^d \to \mathbb{R}$ , we denote by  $\partial f(\mathbf{x})$  its subdifferential at x, which is defined as the set of subgradients of f at  $\mathbf{x}$ , i.e.,  $\partial f(\mathbf{x}) \triangleq \{\mathbf{g} \in \mathbb{R}^d | f(\mathbf{y}) \ge f(\mathbf{x}) + \mathbf{g}^T(\mathbf{y} - \mathbf{x}), \forall \mathbf{y} \in \mathbb{R}^d\}$ . Since f is convex,  $\partial f(\cdot)$  is nonempty. A function f is said to be L-Lipschitz continuous if

$$|f(\mathbf{x}) - f(\mathbf{y})| \le L \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$
(1.2)

Note that the *L*-Lipschitz continuity of f is equivalent to the subgradients of f being uniformly bounded by L [36]. A function f is said to be  $\mu$ -strongly convex if f satisfies

$$f(\mathbf{y}) - f(\mathbf{x}) - \mathbf{g}(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \ge \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2 \qquad \forall \mathbf{x}, \mathbf{y}.$$
 (1.3)

Note that since the set  $\mathcal{X}$  is compact, there exists a point  $\mathbf{x}^*$  which solves Problem (1.1). However,  $\mathbf{x}^*$  may not be unique. We will use  $\mathcal{X}^*$  to denote the set of optimal solutions to Problem (1.1). Given a solution  $\mathbf{x}^* \in \mathcal{X}^*$  we denote  $f^* = \sum_{i=1}^n f_i(\mathbf{x}^*)$ . Also, due to the compactness of  $\mathcal{X}$ , the subgradients of  $f_i$  are uniformly bounded in  $\mathcal{X}$ . We state this observation formally in the following proposition.

PROPOSITION 1.1. There exists a positive constant  $L_i$ , for all  $i \in \mathcal{V}$ , such that the 2-norm of subgradients  $\mathbf{g}_i(\cdot)$  of  $f_i$  are uniformly bounded by  $L_i$  in  $\mathcal{X}$ , i.e., the following condition holds

$$\|\mathbf{g}_i(\mathbf{x})\| \le L_i, \qquad \text{for all } \mathbf{x} \in \mathcal{X}. \tag{1.4}$$

2. Distributed Subgradient Methods. For solving problem (1.1), we are interested in DSG methods [37], where each node *i* maintains its own version of the decision variables  $\mathbf{x}_i \in \mathbb{R}^d$ ; the goal is to have all the  $\mathbf{x}_i$  converge to  $\mathbf{x}^*$ , a solution of problem (1.1). Each node is only allowed to interact with its neighbors that are directly connected to it through a connected and undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V} = \{1, \ldots, n\}$  and  $\mathcal{E} = (\mathcal{V} \times \mathcal{V})$  are the vertex and edge sets, respectively. Each node *i* then iteratively updates  $\mathbf{x}_i$  as

$$\mathbf{x}_{i}(k+1) = \left[\sum_{j \in \mathcal{N}_{i}} a_{ij}\mathbf{x}_{j}(k) - \alpha(k)\mathbf{g}_{i}(\mathbf{x}_{i}(k))\right]_{\mathcal{X}},$$
(2.1)

where  $\alpha(k)$  is some sequence of stepsizes,  $\mathbf{g}_i(\mathbf{x}_i(k)) \in \partial f_i(\mathbf{x}_i(k))$ , and  $\mathcal{N} := \{j \in \mathcal{V} \mid (i, j) \in \mathcal{E}\}$  is the set of node *i*'s neighbors. The  $a_{ij}$  above are positive weights that can be non-zero if there is an edge between nodes *i* and *j*, and otherwise can be assigned by node *i*. We will collect these weights into a  $n \times n$  matrix  $\mathbf{A}$ , and assume it meets the following conditions throughout.

ASSUMPTION 1. The matrix  $\mathbf{A}$ , whose (i, j)-th entries are  $a_{ij}$ , is doubly stochastic, i.e.,  $\sum_{i=1}^{n} a_{ij} = 1$  for all j and  $\sum_{j=1}^{n} a_{ij} = 1$  for all i. Moreover,  $\mathbf{A}$  is irreducible and aperiodic. Finally, the weights  $a_{ij} > 0$  if and only if  $(i, j) \in \mathcal{E}$  otherwise  $a_{ij} = 0$ .

This assumption also implies that **A** has a largest singular value of 1, and its other singular values are strictly less than 1; see for example, the Perron-Frobenius theorem [38]. We denote by  $\sigma_2 \in (0, 1)$  the second largest singular value of **A**, which is a key quantity in the analysis of the mixing time of a Markov chain with transition probabilities given by **A**.

**2.1. Adaptive Quantization.** Each iteration in (2.1) requires every node to communicate its estimate of the decision variables to its neighbors. Theorems 3.6 and 3.8 below study the convergence rate of a modified version of (2.1) when these communications are quantized using our proposed adaptive quantization method. We first present in this section some fundamentals of quantized communication.

To explain the main idea of our approach, we start with the uniform quantization method to quantize a single real number  $x \in [\ell, u]$ . In particular, we divide the interval into *B* bins whose end points are denoted by  $\tau_i$ ,  $\ell = \tau_1 \leq \tau_2 \leq \ldots \leq \tau_B = u$ . We assume that the points  $\tau_i$  are uniformly spaced with distance  $\Delta$ , i.e.,  $\Delta = \tau_{i+1} - \tau_i =$  $(u - \ell) / (B - 1)$  for all  $i = 0, \ldots, B - 1$ . Thus  $b = \lceil \log_2(B) \rceil$  bits can be used to index the  $\{\tau_i\}$ .

Next, given a value  $x \in [\ell, u]$  we denote by  $q = \mathcal{Q}(x)$  its quantized value where

$$\mathcal{Q}(x) \triangleq \min_{\ell} |\tau_{\ell} - x|.$$
(2.2)

If  $\tau_i$  and  $\tau_{i+1}$  achieve the minimal value in Eq. (2.3), then without loss of generality we set  $Q(x) = \tau_i$ . Also, by Eq. (2.2) the quantized error is given as

$$|x - \mathcal{Q}(x)| \le \Delta = \frac{u - \ell}{B - 1}.$$
(2.3)

When the quantized interval depends on time, i.e.,  $[\ell(k), u(k)]$ , we denote by  $\mathcal{Q}_k(x)$ the uniform quantization of x over the time-varying interval  $[\ell(k), u(k)]$  at time k. We can see from Eq. (2.3) that when B is fixed the quantized error depends directly on the size of the interval  $[\ell(k), u(k)]$ . The main idea of our approach is to refine this interval at every time step so that this interval is shrinking, which implies that the quantized error decays to zero. We will refer to this scheme as adaptive quantization, where a distributed implementation will be proposed in the next section.

Moreover, since the constraint set  $\mathcal{X}$  is compact, it is contained in a rectangular set

$$\mathcal{X} \subset \mathcal{R} \triangleq [\ell, \mathbf{u}] = [\ell^1, u^1] \times \ldots \times [\ell^d, u^d],$$

for some  $\{(\ell^i, u^i)\}$ . We quantize a  $\mathbf{x} \in \mathcal{R}$  by applying the procedure above independently on each component, reusing the notation  $\mathcal{Q}(\mathbf{x})$  to indicate that (2.2) is applied component wise. Thus, in this case the total number of bits required to quantize the whole vector is  $b \times d$ . Moreover, the constants in the convergence results presented below will be a function of the interval lengths  $u^i - \ell^i$  for each coordinate i of  $\mathcal{X}$ .

Finally, as will be seen in Eq. (2.4) in the next section, to implement the adaptive quantization method and update its estimate each node  $i \in \mathcal{V}$  has to apply the following two steps of its encoding/decoding scheme in each iteration  $k \geq 0$ .

- 1. Node *i* computes the quantized interval  $\mathcal{R}_i(k)$  to find its quantized value  $\mathbf{q}_i(k) = \mathcal{Q}_k(\mathbf{x}_i(k))$  by using the uniform quantization over this interval. As mentioned, this interval has to be computed in such a way that the quantized error  $\Delta_i(k) = \mathbf{x}_i(k) \mathbf{q}_i(k)$  decays to zero.
- 2. We note that for each iteration  $k \ge 0$  node *i* only receives a sequence of *b* bits  $\mathbf{q}_{j}^{b}(k)$ , e.g.  $\mathbf{q}_{j}^{b}(k) = 01010010$  a sequence of 8 bits, representing the value of  $\mathbf{q}_{j}(k)$  over  $\mathcal{R}_{j}(k)$ . Thus, node *i* has to decode  $\mathbf{q}_{j}^{b}(k)$  to recover  $\mathbf{q}_{j}(k)$ . This step can be done if node *i* knows  $\mathcal{R}_{j}(k)$  for  $j \in \mathcal{N}_{i}$ .

2.2. Distributed Subgradient Methods under Adaptive Quantization. Our focus in this paper is to study the impact of quantized communication between the nodes on the performance of DSG. In particular, at any iteration  $k \ge 0$  the nodes are only allowed to send and receive the quantized values of their estimates to their neighboring nodes. Due to the quantization, we first modify the update in Eq. (2.1) to take into account the quantized error. That is, each node i, for all  $i \in \mathcal{V}$ , now considers the following update

$$\mathbf{x}_{i}(k+1) = \left[ \underbrace{\mathbf{x}_{i}(k) - \mathbf{q}_{i}(k)}_{\text{"quantization error"}} + \sum_{j \in \mathcal{N}_{i}} a_{ij} \mathbf{q}_{j}(k) - \alpha(k) \mathbf{g}_{i}(\mathbf{x}_{i}(k)) \right]_{\mathcal{X}}, \quad (2.4)$$

where  $\mathbf{q}_i(k) = \mathcal{Q}_k(\mathbf{x}_i(k))$ , for all  $i \in \mathcal{V}$ , is the quantized value of  $\mathbf{x}_i(k)$  over some interval  $\mathcal{R}_i(k) \triangleq [\boldsymbol{\ell}_i(k), \mathbf{u}_i(k)]$  for each step  $k \geq 0$ . The update (2.4) has a simple interpretation as follows. At time  $k \geq 0$ , each node *i* first obtains the quantized value  $\mathbf{q}_i(k)$  of its value  $\mathbf{x}_i(k)$ . Each node *i* then formulates the weighted average of

its quantized value  $\mathbf{q}_i(k)$  and the quantized values  $\mathbf{q}_i(k)$  received from its neighbors  $j \in \mathcal{N}_i$ , with the goal of seeking a consensus on their estimates. In addition, each node also introduces its quantized error into its own update, with the goal of eliminating the bias due to the quantized error over the network. Each node then moves along the subgradients of its respective objective function to update its estimates, pushing the consensus point toward the optimal set  $\mathcal{X}^*$ . The distributed subgradient algorithm under random quantization is formally stated in Algorithm 2.1.

Algorithm 2.1 Distributed Subgradient Algorithm Under Adaptive Quantization

1. Initialize: Let  $L = \sum_{i \in \mathcal{V}} L_i, \gamma = \frac{96+48L}{1-\sigma_2}$ .

Each node  $i \in \mathcal{V}$  initializes

- a. Divide  $[\ell, \mathbf{u}]$  (that contains  $\mathcal{X}$ ) into  $B^d$  rectangular bins uniformly componentwise as described above, i.e.,  $\boldsymbol{\ell} = \boldsymbol{\tau}_1 \leq \ldots \leq \boldsymbol{\tau}_{B^d} = \mathbf{u}$  .
- b. Set  $\mathbf{x}_i(0) = \mathbf{q}_i(0) = \boldsymbol{\tau}_m$  for some arbitrarily index  $m \in [1, B^d]$ . Compute  $\mathbf{q}_{i}^{b}(0)$  using the quantization scheme  $\mathcal{Q}_{0}(\mathbf{x}_{i}(0))$  over the set  $[\boldsymbol{\ell}, \mathbf{u}]$ .
- c. A sequence of positive and nonincreasing step sizes  $\{\alpha(k)\}$ .
- 2. Iteration: For  $k = 0, 1, \ldots$ , node  $i \in \mathcal{V}$  implements
  - a. Send  $\mathbf{q}_i^b(k)$  to node  $j \in \mathcal{N}_i$
  - b. Receive  $\mathbf{q}_i^b(k)$  from node  $j \in \mathcal{N}_i$ 
    - If k = 0 then  $\mathcal{R}_i(0) = [\ell, \mathbf{u}]$ , otherwise

$$\mathcal{R}_{j}(k) \triangleq \left[\mathbf{q}_{j}(k-1) - \frac{\gamma}{2}\alpha(k-1)\mathbf{1}, \, \mathbf{q}_{j}(k-1) + \frac{\gamma}{2}\alpha(k-1)\mathbf{1}\right]$$

- Recover  $\mathbf{q}_{i}(k)$  from  $\mathbf{q}_{i}^{b}(k)$  by using uniform quantization over  $\mathcal{R}_{i}(k)$ c. Use  $\mathbf{q}_i(k)$  and update

$$\mathbf{x}_{i}(k+1) = \left[\mathbf{x}_{i}(k) - \mathbf{q}_{i}(k) + \sum_{j \in \mathcal{N}_{i}} a_{ij}\mathbf{q}_{j}(k) - \alpha(k)\mathbf{g}_{i}(\mathbf{x}_{i}(k))\right]_{\mathcal{X}}$$

d. Compute  $\mathbf{q}_i^b(k+1)$  and  $\mathbf{q}_i(k+1)$  by using  $\mathcal{Q}_{k+1}(\mathbf{x}_i(k+1))$  over the interval

$$\mathcal{R}_i(k+1) \triangleq \left[\mathbf{q}_i(k) - \frac{\gamma}{2}\alpha(k)\mathbf{1}, \, \mathbf{q}_i(k) + \frac{\gamma}{2}\alpha(k)\mathbf{1}\right]$$

e. Update the output

$$\mathbf{z}_{i}(k) = \frac{\sum_{t=0}^{k} \alpha(t) \mathbf{x}_{i}(t)}{\sum_{t=0}^{k} \alpha(t)} \cdot$$
(2.5)

Steps (a)-(d) in Algorithm 2.1 are to guarantee the conditions in the two steps (1) and (2) mentioned in the preceding subsection. In particular, step (d) shows how each node i defines its quantized interval  $\mathcal{R}_i(k)$  and computes  $\mathbf{q}_i(k)$  and  $\mathbf{q}_i^b(k)$ . In addition, it is clear from the definition of  $\mathcal{R}_i(k)$  that

$$\|\Delta_i(k)\| = \|\mathbf{x}_i(k) - \mathbf{q}_i(k)\| \lesssim \mathcal{O}(\alpha(k)),$$

which decays to zero. On the other hand, step (b) shows how node i can recover the 6

quantized value  $\mathbf{q}_j(k)$  from the encoded bits,  $\mathbf{q}_j^b(k)$  and  $\mathcal{R}_j(k)$ . Note that the interval  $\mathcal{R}_j(k)$  can be calculated locally at node *i* at any time *k* since node *i* knows the previous value  $\mathbf{q}_j(k-1)$  for  $j \in \mathcal{N}_i$ . Thus, Algorithm 2.1 is fully distributed, that is, the updates at the nodes are executed in parallel and based only on local interactions between nodes. Finally, we show in the next section that by executing Eq. (2.4) we have  $\mathbf{x}_i(k) \in \mathcal{R}_i(k)$  for all  $k \geq 0$ . This observation explains our motivation in defining  $\mathcal{R}_i(k)$  in step (d).

3. Main Results. This section establishes rates of convergence for Algorithm 2.1 in the cases where  $f_i$  are convex and strongly convex. These results show that under adaptive quantization, the convergence rates of the distributed subgradient algorithm are essentially unaffected by the finite communication bandwidths, except for a constant factor that captures the size of these bandwidths.

The main steps of our analysis are as follows. As we observed in the previous section, the adaptive quantization scheme forces the quantization error  $\Delta_i(k)$  to decay to zero at the same rate as the step sizes  $\alpha(k)$ . Our first step is to use this fact to show that the distance between the estimates  $\mathbf{x}_i(k)$  to the average  $\bar{\mathbf{x}}(k)$  converges to zero, implying the nodes eventually reach consensus. Next we will show that the update of (descent on)  $\bar{\mathbf{x}}(k)$  mirrors the update of standard centralized subgradient methods. This allows us to study the convergence rate of Algorithm 2.1 using the standard outline for the analysis of centralized subgradient methods.

When the  $f_i$  are convex, and we use stepsizes  $\alpha(k) = 1/\sqrt{k+1}$ , we show that the time-weighted average  $\mathbf{z}_i(k)$  in (2.5) obeys

$$f(\mathbf{z}_i(k)) - f^* \lesssim \frac{1}{(2^b - 1)^2 (1 - \sigma_2)} \cdot \frac{\ln k}{\sqrt{k}},$$

where  $1 - \sigma_2$  is the spectral gap that quantifies the connectivity of the underlying network, and  $1/(2^b-1)$  is the resolution of the quantizer using *b* communication bits. When the objective function is strongly convex, using stepsizes  $\alpha(k) = a/(k+1)$  for appropriately chosen constant *a*, we have a refined rate on the convergence of the decision variables themselves,

$$\|\mathbf{z}_i(k) - \mathbf{x}^*\|^2 \lesssim \frac{1}{(2^b - 1)^2 (1 - \sigma_2)} \cdot \frac{\ln k}{k}.$$

Aside from the constant  $1/(2^b - 1)$ , these results match the standard results for the distributed subgradient method with perfect communication, meaning that the quantization does not qualitatively affect the behavior of the algorithm.

**3.1. Preliminaries.** Given a vector  $\mathbf{v} \in \mathbb{R}^d$  we denote by  $\boldsymbol{\xi}(\mathbf{v})$  the error due to the projection of  $\mathbf{v}$  on to  $\mathcal{X}$ ,

$$\boldsymbol{\xi}(\mathbf{v}) = \mathbf{v} - \left[\mathbf{v}\right]_{\mathcal{X}},$$

and rewrite Eq. (2.4) as

$$\mathbf{v}_{i}(k) = \left(\sum_{j \in \mathcal{N}_{i}} a_{ij} \mathbf{x}_{j}(k)\right) + \mathbf{x}_{i}(k) - \mathbf{q}_{i}(k) + \sum_{j \in \mathcal{N}_{i}} a_{ij}(\mathbf{q}_{j}(k) - \mathbf{x}_{j}(k)) - \alpha(k)\mathbf{g}_{i}(\mathbf{x}_{i}(k)),$$
  
$$\mathbf{x}_{i}(k+1) = [\mathbf{v}_{i}(k)]_{\mathcal{X}} = \mathbf{v}_{i}(k) - \boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k)).$$
  
(3.1)

Stacking the  $\mathbf{x}_i^T$  as rows in the  $n \times d$  matrix  $\mathbf{X}$  (and doing likewise with the  $\mathbf{v}_i, \mathbf{q}_i, \boldsymbol{\xi}_i$ ), we can write the above in matrix form as

$$\mathbf{V}(k) = \mathbf{A}\mathbf{X}(k) + (\mathbf{I} - \mathbf{A})(\mathbf{X}(k) - \mathbf{Q}(k)) - \alpha(k)\mathbf{G}(\mathbf{X}(k)),$$
  
$$\mathbf{X}(k+1) = \mathbf{V}(k) - \mathbf{\Xi}(\mathbf{V}(k)),$$
  
(3.2)

where **A** is the adjacency matrix in Assumption 1. Let  $\bar{\mathbf{x}}(k)$  and  $\bar{\boldsymbol{\xi}}(k)$  be the averages of  $\mathbf{x}_i(k)$  and  $\boldsymbol{\xi}_i(\mathbf{v}_i(k))$  across all nodes at time k:

$$\bar{\mathbf{x}}(k) = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_{i}(k) = \frac{1}{n} \mathbf{X}^{T} \mathbf{1} \in \mathbb{R}^{d} \quad \text{and} \quad \bar{\boldsymbol{\xi}}(k) = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k)) = \frac{1}{n} \mathbf{\Xi}(\mathbf{V}(k))^{T} \mathbf{1} \in \mathbb{R}^{d}.$$

Since  $\mathbf{1}^T \mathbf{A} = \mathbf{1}^T$ , (3.2) gives

$$\bar{\mathbf{v}}(k) = \bar{\mathbf{x}}(k) - \frac{\alpha(k)}{n} \sum_{i=1}^{n} \mathbf{g}_i(\mathbf{x}_i(k)),$$

$$\bar{\mathbf{x}}(k+1) = \bar{\mathbf{v}}(k) - \bar{\boldsymbol{\xi}}(k).$$
(3.3)

Recall that  $\Delta_i(k) = \mathbf{x}_i(k) - \mathbf{q}_i(k) \in \mathbb{R}^d$  and let  $\Delta(k)$  be defined as

$$\Delta(k) = \begin{bmatrix} -\Delta_1(k)^T - \\ \cdots \\ -\Delta_n(k)^T - \end{bmatrix} \in \mathbb{R}^{n \times d}.$$

We now consider the following sequence of lemmas, which provides fundamental preliminaries for our main results given in the next sections. In the sequel, we will consider two choices of the step sizes  $\alpha(k)$ , that is,  $\alpha(k) = 1/\sqrt{k+1}$  or  $\alpha(k) = 1/(k+1)$ . These choices of step sizes also are used to establish our main results in the next section. For ease of exposition we delay all the proofs of the results in this section to Appendix A.

We first provide an upper bound for the projection error  $\boldsymbol{\xi}_i$  in the following lemma.

LEMMA 3.1. Suppose that Assumption 1 holds. Let the sequence  $\{\mathbf{x}_i(k)\}$ , for all  $i \in \mathcal{V}$ , be generated by Algorithm 2.1. Then for all  $i \in \mathcal{V}$  we have

$$\|\boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))\| \leq \sum_{j \in \mathcal{N}_{i}} a_{ij} \|\Delta_{i}(k) - \Delta_{j}(k)\| + L_{i}\alpha(k).$$
(3.4)

In addition, let  $L = \sum_{i \in \mathcal{V}} L_i$ . Then, we obtain

$$\sum_{i=1}^{n} \|\boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))\|^{2} \leq 8 \sum_{i=1}^{n} \|\Delta_{i}(k)\|^{2} + 2L^{2}\alpha^{2}(k).$$
(3.5)

Next we provide an upper bound for the consensus errors  $\|\mathbf{x}_i(k) - \bar{\mathbf{x}}(k)\|$  in the following lemma.

LEMMA 3.2. Suppose that Assumption 1 holds. Let the sequence  $\{\mathbf{x}_i(k)\}$ , for all  $i \in \mathcal{V}$ , be generated by Algorithm 2.1. In addition, let  $\{\alpha(k)\}$  be a nonnegative nonincreasing sequence of stepsizes. Then, we have

$$\|\mathbf{X}(k+1) - \mathbf{1}\bar{\mathbf{x}}(k+1)^T\| \le 6\sum_{\substack{t=0\\8}}^k \sigma_2^{k-t} \|\Delta(t)\| + 3L\sum_{t=0}^k \sigma_2^{k-t}\alpha(t).$$
(3.6)

As mentioned, the main motivation of the adaptive quantization is to eliminate the impact of quantized errors. In particular, we will show that the quantized errors produced by Algorithm 2.1 decrease to zero at the same rate with the step size  $\alpha(k)$ . To do that, we require the following technical condition.

ASSUMPTION 2. Let  $\gamma = 48(2 + L) / (1 - \sigma_2)$ . Then the number bits b of the communication bandwdith satisfies

$$\sqrt{nd\gamma} \le 2^b - 1. \tag{3.7}$$

The following lemma is to show that the quantized error  $\|\Delta_i(k)\| \leq \alpha(k)$ , for all  $i \in \mathcal{V}$ , for the case  $\sigma_2^k \leq \alpha(k)$  for all  $k \geq 0$ . When  $\sigma_2^k \geq \alpha(k)$  for some small k, e.g.,  $\sigma_2$  is closed to 1 but not equal, one can show from our analysis that  $\|\Delta_i(k)\| \leq \sigma_2^k$ , which is eventually converge to zero faster than  $\alpha(k)$ . Since this issue has been captured by the rate of consensus in Eq. (3.6), we skip it here for simplicity.

LEMMA 3.3. Suppose that Assumption 1 and 2 hold. Let the sequence  $\{\mathbf{x}_i(k)\}$ , for all  $i \in \mathcal{V}$ , be generated by Algorithm 2.1. Let  $\alpha(k)$  be either  $\alpha(k) = 1/\sqrt{k+1}$  or  $\alpha(k) = 1/(k+1)$ . Then we have all  $k \ge 0$ 

$$\mathbf{x}_{i}(k+1) \in \mathcal{R}_{i}(k+1) \triangleq \left[\mathbf{q}_{i}(k) - \frac{\gamma}{2}\alpha(k)\mathbf{1}, \, \mathbf{q}_{i}(k) + \frac{\gamma}{2}\alpha(k)\mathbf{1}\right].$$
(3.8)

In addition, we also have

$$\|\Delta_i(k)\| \le \frac{\sqrt{d\gamma}}{2^b - 1} \alpha(k). \tag{3.9}$$

The following lemma is a consequence of Lemmas 3.2 and 3.3.

LEMMA 3.4. Suppose that Assumptions 1 and 2 hold. Let the sequence  $\{\mathbf{x}_i(k)\}$ , for all  $i \in \mathcal{V}$ , be generated by Algorithm 2.1. Let  $\alpha(k)$  be either  $\alpha(k) = 1/\sqrt{k+1}$  or  $\alpha(k) = 1/(k+1)$ . Then, we have

$$\lim_{k \to \infty} \mathbf{x}_i(k) = \lim_{k \to \infty} \mathbf{x}_j(k), \qquad \forall i, j \in \mathcal{V}.$$
(3.10)

In addition, if  $\alpha(k)$  is also square-summable, i.e.,

$$\sum_{k=0}^{\infty} \alpha^2(k) < \infty, \tag{3.11}$$

then for all  $k \ge 0$  we have

$$\sum_{t=0}^{k} \alpha(t) \| \mathbf{X}(t) - \mathbf{1}\bar{\mathbf{x}}(t)^{T} \| \le \left( \frac{6\sqrt{nd\gamma} + 3L2^{b}}{(1 - \sigma_{2})(2^{b} - 1)} \right) \sum_{t=0}^{k} \alpha^{2}(t) < \infty.$$
(3.12)

If  $\alpha(k) = 1 / \sqrt{k+1}$  then we have for all  $k \ge 0$ ,

$$\sum_{t=0}^{k} \alpha(t) \|\mathbf{X}(t) - \mathbf{1}\bar{\mathbf{x}}(t)^{T}\| \le \left(\frac{6\sqrt{nd\gamma} + 3L2^{b}}{(1-\sigma_{2})(2^{b}-1)}\right) (\ln(k+1)+1).$$
(3.13)

Finally, we provide an upper bound for the optimal distance  $\|\bar{\mathbf{x}}(k) - \mathbf{x}^*\|^2$  in the following lemma.

LEMMA 3.5. Suppose that Assumptions 1 and 2 hold. Let the sequence  $\{\mathbf{x}_i(k)\}$ , for all  $i \in \mathcal{V}$ , be generated by Algorithm 2.1. In addition, let  $\mathbf{x}^* \in \mathcal{X}^*$  be a solution of problem (1.1). Let  $\alpha(k)$  be either  $\alpha(k) = 1/\sqrt{k+1}$  or  $\alpha(k) = 1/(k+1)$ . Then, we have

$$\|\bar{\mathbf{x}}(k+1) - \mathbf{x}^*\|^2 \le \|\bar{\mathbf{x}}(k) - \mathbf{x}^*\|^2 - \frac{2\alpha(k)}{n} \sum_{i=1}^n \mathbf{g}_i(\mathbf{x}_i(k))^T (\mathbf{x}_i(k) - \mathbf{x}^*) + \frac{2\left(4\sqrt{nd}\gamma + 3L2^b\right)}{\sqrt{n}(2^b - 1)} \alpha(k) \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^T\| + \frac{2(4\sqrt{nd}\gamma + 3L2^b)^2}{(2^b - 1)^2\sqrt{n}} \alpha^2(k) \cdot$$
(3.14)

3.2. Convergence Results of Convex Functions. We now present the first main result of this paper, which is the rate of convergence of Algorithm 2.1 to the optimal value of problem (1.1) when the local functions  $f_i$  are convex. Since the update of  $\bar{\mathbf{x}}(k)$  in Eq. (3.3) can be viewed as a variant of a centralized projected subgradient methods used to solve problem (1.1), we utilize standard techniques in the analysis of these methods to derive the rate of convergence of Algorithm 2.1. Specifically, at any time  $k \geq 0$  if each node  $i \in \mathcal{V}$  maintains a variable  $\mathbf{z}_i(k)$  to compute the time-weighted average of its estimate  $\mathbf{x}_i(k)$  and if the stepsize  $\alpha(k)$  decays as  $\alpha(k) = 1/\sqrt{k+1}$ , the objective function value f in Eq. (1.1) estimated at each  $\mathbf{z}_i(k)$  converges to the optimal value with a rate  $\mathcal{O}\left(\eta \ln(k+1)/\sqrt{k+1}\right)$ , where  $\eta$  is some constant depending on the algebraic connectivity  $1 - \sigma_2$  of the network, the number of quantized bits b, and the Lipschitz constants  $L_i$  of  $f_i$ . We also note that this condition on the stepsizes is also used to study the convergence rate of centralized subgradient methods [39]. The following theorem is used to show the convergence rate of Algorithm 2.1.

THEOREM 3.6. Suppose that Assumptions 1 and 2 hold. Let the sequence  $\{\mathbf{x}_i(k)\}$ , for all  $i \in \mathcal{V}$ , be generated by Algorithm 2.1. In addition, let  $\alpha(k) = 1/\sqrt{k+1}$ . Moreover, suppose that each node *i*, for all  $i \in \mathcal{V}$ , stores a variable  $\mathbf{z}_i \in \mathbb{R}^d$  initiated arbitrarily in  $\mathcal{X}$  and updated as

$$\mathbf{z}_{i}(k) = \frac{\sum_{t=0}^{k} \alpha(t) \mathbf{x}_{i}(t)}{\sum_{t=0}^{k} \alpha(t)}, \quad \forall i \in \mathcal{V}.$$
(3.15)

Then for all  $i \in \mathcal{V}$  and  $k \ge 0$  we have

$$f(\mathbf{z}_{i}(k)) - f^{*} \leq \frac{n \|\bar{\mathbf{x}}(0) - \mathbf{x}^{*}\|^{2}}{2\sqrt{k+1}} + \frac{\sqrt{n}(6\sqrt{nd\gamma} + 5L2^{b})^{2}}{(1 - \sigma_{2})(2^{b} - 1)^{2}} \frac{(\ln(k+1) + 1)}{\sqrt{k+1}}.$$
 (3.16)

*Proof.* For convenience, let  $\mathbf{r}(k) = \bar{\mathbf{x}}(k) - \mathbf{x}^*$ , where  $\mathbf{x}^* \in \mathcal{X}^*$  is a solution of problem (1.1). By Eq. (3.14) we have

$$\|\mathbf{r}(k+1)\|^{2} \leq \|\mathbf{r}(k)\|^{2} - \frac{2\alpha(k)}{n} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k))^{T}(\mathbf{x}_{i}(k) - \mathbf{x}^{*}) + \frac{2(4\sqrt{nd\gamma} + 3L2^{b})^{2}}{(2^{b} - 1)^{2}\sqrt{n}} \alpha^{2}(k) + \frac{2\left(4\sqrt{nd\gamma} + 3L2^{b}\right)}{\sqrt{n}(2^{b} - 1)} \alpha(k) \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\|,$$
10

which by the convexity of  $f_i$  yields

$$\|\mathbf{r}(k+1)\|^{2} \leq \|\mathbf{r}(k)\|^{2} + \frac{2(4\sqrt{nd\gamma} + 3L2^{b})^{2}}{(2^{b} - 1)^{2}\sqrt{n}}\alpha^{2}(k) + \frac{2\left(4\sqrt{nd\gamma} + 3L2^{b}\right)}{\sqrt{n}(2^{b} - 1)}\alpha(k)\|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| - \frac{2\alpha(k)}{n}\left(\sum_{i=1}^{2} f_{i}(\mathbf{x}_{i}(k)) - f_{i}(\mathbf{x}^{*})\right)\right), \qquad (3.17)$$

We now analyze the last term on the right-hand side of Eq. (3.17). Indeed, by Eq. (1.4) and using  $f = \sum_{i=1}^{n} f_i$  and  $f^* = f(x^*)$ , we have for a fixed  $\ell \in \mathcal{V}$ 

$$-\sum_{i=1}^{n} f_{i}(\mathbf{x}_{i}(k)) - f_{i}(\mathbf{x}^{*}) = -\sum_{i=1}^{n} \left( f_{i}(\mathbf{x}_{i}(k)) - f_{i}(\bar{\mathbf{x}}(k)) + f_{i}(\bar{\mathbf{x}}(k)) - f_{i}(\mathbf{x}^{*}) \right)$$

$$\leq \sum_{i=1}^{n} L_{i} |\mathbf{x}_{i}(k) - \bar{\mathbf{x}}(k)| - \left( f(\bar{\mathbf{x}}(k)) - f^{*} \right)$$

$$\leq L ||\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}|| - \left( f(\bar{\mathbf{x}}(k)) - f(\mathbf{x}_{\ell}(k)) + f(\mathbf{x}_{\ell}(k)) - f^{*} \right)$$

$$\leq 2L ||\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}|| - \left( f(\mathbf{x}_{\ell}(k)) - f^{*} \right), \qquad (3.18)$$

which when substituting into Eq. (3.17) yields

$$\|\mathbf{r}(k+1)\|^{2} \leq \|\mathbf{r}(k)\|^{2} + \frac{2(4\sqrt{nd\gamma} + 3L2^{b})^{2}}{(2^{b} - 1)^{2}\sqrt{n}}\alpha^{2}(k) - \frac{2}{n}\alpha(k)\Big(f(\mathbf{x}_{\ell}(k)) - f^{*}\Big) \\ + \left(\frac{2\left(4\sqrt{nd\gamma} + 3L2^{b}\right)}{\sqrt{n}(2^{b} - 1)}\Big) + \frac{4L}{n}\right)\alpha(k)\|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| \\ = \|\mathbf{r}(k)\|^{2} + \frac{2(4\sqrt{nd\gamma} + 3L2^{b})^{2}}{(2^{b} - 1)^{2}\sqrt{n}}\alpha^{2}(k) - \frac{2}{n}\alpha(k)\Big(f(\mathbf{x}_{\ell}(k)) - f^{*}\Big) \\ + \frac{2\left(4\sqrt{nd\gamma} + 5L2^{b}\right)}{\sqrt{n}(2^{b} - 1)}\alpha(k)\|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\|,$$
(3.19)

which when iteratively updating over  $k=0,\ldots,K$  for some  $K\geq 0$  we have

$$\|\mathbf{r}(K+1)\|^{2} \leq \|\mathbf{r}(0)\|^{2} + \frac{2(4\sqrt{nd\gamma} + 3L2^{b})^{2}}{(2^{b} - 1)^{2}\sqrt{n}} \sum_{k=0}^{K} \alpha^{2}(k) + \frac{2\left(4\sqrt{nd\gamma} + 5L2^{b}\right)}{\sqrt{n}(2^{b} - 1)} \sum_{k=0}^{K} \alpha(k) \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| - \frac{2}{n} \sum_{k=0}^{K} \alpha(k) \left(f(\mathbf{x}_{\ell}(k)) - f^{*}\right).$$
11

Since  $\alpha(k) = 1/\sqrt{k+1}$  we now use Eq. (3.13) into the preceding relation to have

$$\begin{aligned} \|\mathbf{r}(K+1)\|^2 &\leq \|\mathbf{r}(0)\|^2 + \frac{2(4\sqrt{nd}\gamma + 3L2^b)^2}{(2^b - 1)^2\sqrt{n}} (\ln(K+1) + 1) \\ &+ \frac{2\left(4\sqrt{nd}\gamma + 5L2^b\right)}{\sqrt{n}(2^b - 1)} \left(\frac{6\sqrt{nd}\gamma + 3L2^b}{(1 - \sigma_2)(2^b - 1)}\right) (\ln(K+1) + 1) \\ &- \frac{2}{n} \sum_{k=0}^K \alpha(k) \Big(f(\mathbf{x}_\ell(k)) - f^*\Big) \\ &\leq \|\mathbf{r}(0)\|^2 + \frac{4(6\sqrt{nd}\gamma + 5L2^b)^2}{(1 - \sigma_2)(2^b - 1)^2\sqrt{n}} (\ln(K+1) + 1) \\ &- \frac{2}{n} \sum_{k=0}^K \alpha(k) \Big(f(\mathbf{x}_\ell(k)) - f^*\Big). \end{aligned}$$

Rearranging the preceding relation and dropping the nonnegative  $\|\mathbf{r}(K+1)\|$  we obtain

$$\sum_{k=0}^{K} \alpha(k) \Big( f(\mathbf{x}_{\ell}(k)) - f^* \Big) \le \frac{n \|\mathbf{r}(0)\|^2}{2} + \frac{2\sqrt{n}(6\sqrt{nd\gamma} + 5L2^b)^2}{(1-\sigma_2)(2^b-1)^2} (\ln(K+1) + 1),$$

which by dividing both sides by  $\sum_{k=0}^{K} \alpha(k)$  and using the convexity of f gives Eq. (3.16), i.e.,

$$f(\mathbf{z}_{\ell}(K)) - f^* \le \frac{n \|\mathbf{r}(0)\|^2}{2\sqrt{K+1}} + \frac{2\sqrt{n}(6\sqrt{nd\gamma} + 5L2^b)^2}{(1-\sigma_2)(2^b-1)^2} \frac{(\ln(K+1)+1)}{\sqrt{K+1}},$$

where in the last inequality we use the integral test for  $K \ge 0$  to have

$$\sum_{k=0}^{K} \alpha(k) = \sum_{k=0}^{K} \frac{1}{\sqrt{k+1}} \ge \int_{t=0}^{K+1} \frac{1}{\sqrt{t+1}} dt = 2(\sqrt{K+2}-1) \ge \sqrt{K+1}. \qquad \Box$$

It is worth to mention that under the choice of  $\alpha(k) = 1 / (k+1)$ , for all  $k \ge 0$ , one can show that  $\mathbf{x}_i(k)$  asymptotically converges to  $\mathbf{x}^*$  for all  $i \in \mathcal{V}$ . This is a consequence of Lemmas 3.2 and 3.5, and some standard analysis. The following lemma states this result. The analysis is omitted and can be found in [16, Theorem 3].

LEMMA 3.7. Suppose that Assumptions 1 and 2 hold. Let the sequence  $\{\mathbf{x}_i(k)\}$ , for all  $i \in \mathcal{V}$ , be generated by Algorithm 2.1. Let  $\alpha(k) = 1 / (k+1)$ . Then we obtain

$$\lim_{k \to \infty} \mathbf{x}_i(k) = \mathbf{x}^*, \qquad \text{for all } i \in \mathcal{V}, \tag{3.20}$$

for some  $\mathbf{x}^*$  that is a solution of problem (1.1).

**3.3.** Convergence Results of Strongly Convex Case. In this section, our goal is to study the convergence rate of Algorithm 2.1 when the local functions  $f_i$  are strongly convex, that is, we make the following assumption on  $f_i$ 

ASSUMPTION 3. Each function  $f_i$  is strongly convex with some positive constant  $\mu_i$ , i.e., the condition (1.3) holds.

Under this assumption, we show that if each node  $i \in \mathcal{V}$  maintains a variable  $\mathbf{z}_i(k)$  to compute the time average of its estimate  $\mathbf{x}_i(k)$  and if the stepsize  $\alpha(k)$  decays as  $\alpha(k) = a / (k+1)$  for some properly chosen constant a, the variable  $\mathbf{z}_i(k)$  converges to the optimal solution  $x^*$  of problem (1.1) with a rate  $\mathcal{O}(\eta \ln(k+1) / (k+1))$ , where  $\eta$  is some constant depending on the algebraic connectivity  $1 - \sigma_2$  of the network, the number of quantized bits b, and the constants  $L_i$  and  $\mu_i$  of  $f_i$ . The following theorem is used to show the convergence rate of Algorithm 1 under Assumption 3.

THEOREM 3.8. Suppose that Assumptions 1 and 3 hold. Let the sequence  $\{\mathbf{x}_i(k)\}$ , for all  $i \in \mathcal{V}$ , be generated by Algorithm 2.1. We denote by  $\mu = \min_{i \in \mathcal{V}} \mu_i$ . In addition, let  $\{\alpha(k)\} = a / k + 1$  for some  $a \ge 1 / \mu$ . Moreover, suppose that each node *i*, for all  $i \in \mathcal{V}$ , stores a variable  $\mathbf{z}_i \in \mathbb{R}$  initiated arbitrarily in  $\mathcal{X}$  and updated as

$$\mathbf{z}_{i}(k) = \frac{\sum_{t=0}^{k} \mathbf{x}_{i}(t)}{k+1}, \quad \forall i \in \mathcal{V}.$$
(3.21)

Let  $\mathbf{x}^* \in \mathcal{X}^*$  be a solution of problem (1.1). Then for all  $i \in \mathcal{V}$  and  $k \ge 0$  we have

$$\|\mathbf{z}_{i}(k) - \mathbf{x}^{*}\|^{2} \leq \frac{4\sqrt{n\alpha(0)}(6\sqrt{nd\gamma} + 5L2^{b})^{2}}{(1 - \sigma_{2})(2^{b} - 1)^{2}} \frac{1 + \ln(k+1)}{k+1}.$$
 (3.22)

*Proof.* Let  $\mathbf{x}^*$  be a solution of problem (1.1). For convenience, let  $\mathbf{r}(k) = \bar{\mathbf{x}}(k) - \mathbf{x}^*$ . By Eq. (3.14) we have

$$\|\mathbf{r}(k+1)\|^{2} \leq \|\mathbf{r}(k)\|^{2} - \frac{2\alpha(k)}{n} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k))^{T}(\mathbf{x}_{i}(k) - \mathbf{x}^{*}) + \frac{2\left(4\sqrt{nd}\gamma + 3L2^{b}\right)}{\sqrt{n}(2^{b} - 1)} \alpha(k) \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| + \frac{2(4\sqrt{nd}\gamma + 3L2^{b})^{2}}{(2^{b} - 1)^{2}\sqrt{n}} \alpha^{2}(k) \leq \|\mathbf{r}(k)\| + \frac{2(4\sqrt{nd}\gamma + 3L2^{b})^{2}}{(2^{b} - 1)^{2}\sqrt{n}} \alpha^{2}(k) + \frac{2\left(4\sqrt{nd}\gamma + 3L2^{b}\right)}{\sqrt{n}(2^{b} - 1)} \alpha(k) \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| - \frac{2\alpha(k)}{n} \sum_{i=1}^{n} \left(f_{i}(\mathbf{x}_{i}(k)) - f_{i}(\mathbf{x}^{*}) + \frac{\mu_{i}}{2} \|\mathbf{x}_{i}(k) - \mathbf{x}^{*}\|^{2}\right), \quad (3.23)$$

where the last inequality is due to the strong convexity of  $f_i$ , i.e., Eq. (1.3). First, using the Jensen's inequality on quadratic function  $(\cdot)^2$  we have

$$-\frac{1}{n}\sum_{i=1}^{n}\mu_{i}\|\mathbf{x}_{i}(k)-\mathbf{x}^{*}\|^{2} \leq -\mu\frac{1}{n}\sum_{i=1}^{n}\|\mathbf{x}_{i}(k)-\mathbf{x}^{*}\|^{2} \leq -\mu\|\bar{\mathbf{x}}(k)-\mathbf{x}^{*}\|^{2} = -\mu\|\mathbf{r}(k)\|^{2}.$$

Fix some  $\ell \in \mathcal{V}$ . Then, substituting the preceding relation into Eq. (3.23) and using

Eq. (3.18) yield

$$\begin{aligned} \|\mathbf{r}(k+1)\|^{2} &\leq (1-\mu\alpha(k)) \|\mathbf{r}(k)\|^{2} + \frac{2(4\sqrt{nd\gamma}+3L2^{b})^{2}}{(2^{b}-1)^{2}\sqrt{n}}\alpha^{2}(k) \\ &+ \frac{2\left(4\sqrt{nd\gamma}+3L2^{b}\right)}{\sqrt{n}(2^{b}-1)}\alpha(k)\|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| \\ &- \frac{2\alpha(k)}{n}\sum_{i=1}^{n}\left(f_{i}(\mathbf{x}_{i}(k)) - f_{i}(\mathbf{x}^{*})\right) \\ &\stackrel{(3.18)}{\leq} (1-\mu\alpha(k)) \|\mathbf{r}(k)\|^{2} + \frac{2(4\sqrt{nd\gamma}+3L2^{b})^{2}}{(2^{b}-1)^{2}\sqrt{n}}\alpha^{2}(k) \\ &+ \frac{2\left(4\sqrt{nd\gamma}+5L2^{b}\right)}{\sqrt{n}(2^{b}-1)}\alpha(k)\|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| \\ &- \frac{2}{n}\alpha(k)\Big(f(\mathbf{x}_{\ell}(k)) - f^{*}\Big), \end{aligned}$$

Note that  $\alpha(k) = a / (k + 1)$  with  $a \ge 1 / \mu$ , implying  $\mu \alpha(k) \ge 1 / (k + 1)$ . Thus, the preceding equation gives

$$\|\mathbf{r}(k+1)\|^{2} \leq \frac{k}{k+1} \|\mathbf{r}(k)\|^{2} + \frac{2(4\sqrt{nd\gamma} + 3L2^{b})^{2}}{(2^{b} - 1)^{2}\sqrt{n}} \alpha^{2}(k) + \frac{2\left(4\sqrt{nd\gamma} + 5L2^{b}\right)}{\sqrt{n}(2^{b} - 1)} \alpha(k) \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| - \frac{2}{n} \alpha(k) \left(f(\mathbf{x}_{\ell}(k)) - f^{*}\right).$$

Multiplying both sides of the preceding equation by k+1, and using  $(k+1)\,/\,k\leq 2$  and  $\alpha(k)=\alpha(0)/(k+1)$  we have

$$(k+1) \|\mathbf{r}(k+1)\|^{2} \leq k \|\mathbf{r}(k)\|^{2} + \frac{4\alpha(0)(4\sqrt{nd}\gamma + 3L2^{b})^{2}}{(2^{b}-1)^{2}\sqrt{n}}\alpha(k) + \frac{2\left(4\sqrt{nd}\gamma + 5L2^{b}\right)}{\sqrt{n}(2^{b}-1)} \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| - \frac{2\alpha(0)}{n} \left(f(\mathbf{x}_{\ell}(k)) - f^{*}\right),$$
(3.24)

By Eq. (3.6) and using  $\|\Delta(t)\| \leq \alpha(t)$  we have

$$\sum_{k=0}^{K} \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| \le (3L+6) \sum_{k=0}^{K} \sum_{t=0}^{k-1} \sigma_{2}^{k-t} \alpha(t)$$
$$\le (3L+6) \sum_{k=0}^{K-1} \alpha(k) \sum_{t=k+1}^{K} \sigma_{2}^{t} \le \frac{3L+6}{1-\sigma_{2}} \sum_{k=0}^{K-1} \alpha(k).$$
(3.25)

Next, summing up both sides of Eq. (3.24) over  $k = 0, \ldots, K$  for some  $K \ge 0$ , using 14

the preceding relation, and rearranging we obtain

$$\frac{2\alpha(0)}{n} \sum_{k=0}^{K} \left( f(x_{\ell}(k)) - f^* \right) \le \frac{2\alpha(0)(6\sqrt{nd\gamma} + 5L2^b)^2}{\sqrt{n}(1 - \sigma_2)(2^b - 1)^2} \sum_{k=0}^{K} \alpha(k)$$
$$\le \frac{2\alpha(0)(6\sqrt{nd\gamma} + 5L2^b)^2}{\sqrt{n}(1 - \sigma_2)(2^b - 1)^2} (\ln(K+1) + 1),$$

which when dividing both sides by (K+1)/n and using the convexity of f yields

$$2\alpha(0) \left[ f(\mathbf{z}_{\ell}(K)) - f^* \right] \le \frac{2\sqrt{n\alpha(0)}(6\sqrt{nd\gamma} + 5L2^b)^2}{(1 - \sigma_2)(2^b - 1)^2} \frac{1 + \ln(K + 1)}{K + 1}$$

Since the functions  $f_i$  are strongly convex with constant  $\mu_i$ , f is strongly convex with constant  $\mu$ . Thus, using the preceding equation and  $\alpha(0) = a \ge 1/\mu$  gives Eq. (3.22), i.e.,

$$\|\mathbf{z}_{\ell}(K) - \mathbf{x}^{*}\|^{2} \leq \frac{2}{\mu} \left[ f(\mathbf{z}_{\ell}(K)) - f^{*} \right] \leq \frac{2\sqrt{n}\alpha(0)(6\sqrt{nd}\gamma + 5L2^{b})^{2}}{(1 - \sigma_{2})(2^{b} - 1)^{2}} \frac{1 + \ln(K + 1)}{K + 1} \cdot \Box$$

4. Simulations. In this section, we apply Algorithm 2.1 for solving linear regression problems, the most popular technique for data fitting [40, 41] in statistical machine learning, over a network of processors under random quantization. The goal of this problem is to find a linear relationship between a set of variables and some real value outcome. That is, given a training set  $S = \{(\mathbf{a}_i, b_i) \in \mathbb{R}^d \times \mathbb{R}\}$  for  $i = 1, \ldots, n$ , we want to learn a parameter  $\mathbf{x}$  that minimizes

$$\min_{\mathbf{x}\in\mathcal{X}}\sum_{i=1}^n f_i(\mathbf{x};\mathbf{a}_i,b_i),$$

where  $\mathcal{X} = [-1, 1]^d$  and d = 10, i.e.,  $\mathbf{x}, \mathbf{a}_i \in \mathbb{R}^{10}$ . Here,  $f_i$  are the loss functions defined over the dataset. For the purpose of our simulation, we will consider two loss functions, namely, quadratic loss and absolute loss functions. While the quadratic loss is strongly convex, the absolute loss is only convex.

First, when  $f_i$  are quadratic, we have the well-known least square problem

$$\min_{\mathbf{x}\in\mathcal{X}} \sum_{i=1}^{n} (\mathbf{a}_i^T \mathbf{x} - b_i)^2.$$

Second, regression problems with absolute loss functions (or L1 norm) is often referred to as robust regression, which is known to be robust to outliers [42], given as follows

$$\min_{\mathbf{x}\in\mathcal{X}} \sum_{i=1}^{n} |\mathbf{a}_{i}^{T}\mathbf{x} - b_{i}|.$$

We consider simulated training data sets, i.e.,  $(\mathbf{a}_i, b_i)$  are generated randomly with uniform distribution between [0, 1]. We consider the performance of the distributed subgradient methods on an undirected connected graph of 50 nodes, i.e.,  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and  $n = |\mathcal{V}| = 100$ . Our graph is generated as follows.

1. In each network, we first randomly generate the nodes' coordinates in the plane with uniform distribution.



Fig. 1: The convergence of function values using distributed subgradient methods without (--) [37], with random (--) [2], with time-varying (--) [16], and with adaptive (--) quantization for n = 100 and d = 10 are illustrated.

- 2. Then any two nodes are connected if their distance is less than a reference number r, e.g, r = 0.4 for our simulations.
- 3. Finally we check whether the network is connected. If not we return to step 1 and run the program again.

To implement our algorithm, the adjacency matrix  $\mathbf{A}$  is chosen as a lazy Metropolis matrix corresponding to  $\mathcal{G}$ , i.e.,

$$\mathbf{A} = [a_{ij}] = \begin{cases} \frac{1}{2(\max\{|\mathcal{N}_i|, |\mathcal{N}_j|\})}, & \text{if } (i,j) \in \mathcal{E} \\ 0, & \text{if } (i,j) \notin \mathcal{E} \text{ and } i \neq j \\ 1 - \sum_{j \in \mathcal{N}_i} a_{ij}, & \text{if } i = j \end{cases}$$

It is obvious to see that the lazy Metropolis matrix **A** satisfies Assumption 1.

4.1. Convergence of Function Values. In this simulation, we apply variants of distributed subgradient methods for solving the linear regression problems. In particular, we compare the performance of such methods for three different scenarios, namely, DSG with no quantization (i.e., Eq. (2.1)), DSG with time-varying quantization in [16], distributed stochastic approximation under random quantization [2], and the proposed Algorithm 2.1 with adaptive quantization. In addition, we use 8 bits as the size of the nodes' communication bandwidths. The plots in Fig. 1 show the convergence of these four methods for both quadratic and absolute loss functions.

Note that, DSG with time-varying quantization [16] achieves the same rate of convergence as the one with no quantization [37], but requires that the nodes eventually exchange an infinite number of bits. On the other hand, DSG with random quantization [2] only requires a finite number of bits, but achieves a slow rate of convergence. The adaptive quantization in this paper achieves both benefits of time-varying quantization [16] and random quantization [2], i.e., it achieves the same rate as the algorithm without quantization but only using a finite number of bits. In addition, as observed in Fig.1a for quadratic loss and in Fig. 1b for absolute loss, Algorithm 2.1 performs almost as well as the one without quantization [37], and significantly better than the



Fig. 2: The number of iterations as a function of b using distributed subgradient methods with adaptive quantization for n = 100 and d = 10 are illustrated.

algorithms in [2, 16]

**4.2. Impact of the Number of Bits**, *b*. We now consider the impact of the number of bits *b* on the performance of Algorithm 2.1. In Fig. 2 we plots the number of iterations, needed to obtain the relative error  $f(z_i(k)) - f^* / f^* \leq 0.05$ , as a function of *b*. We see that the more bits we use, the faster the algorithm converges. Moreover, even when only a very small number of bits, for example, b= 4 are used, the algorithm still works very well. Finally, these plots appear to describe the curve  $1/(2^b - 1)^2$  upto some constant in the upper bound of convergence rates given in Theorems 3.6 and 3.8. This implies that the simulation seems to agree with our results.

5. Concluding Remarks. In this paper, we consider distributed optimization over networks of nodes under finite bandwidths, and so information exchanged across the network must be quantized. For solving such problems, we consider distributed subgradient methods under quantization. Our main contribution is to propose a novel adaptive quantization, which quantizes the nodes' estimates based on the progress of the algorithm. Under this adaptive quantization, we show that the rates of convergence of DSG are unaffected by communication constraints. A natural question from this work is to ask whether the proposed adaptive quantization can be extended to study other distributed algorithms under finite bandwidths, such as, distributed primal-dual, ADMM, dual averaging, and mirror-descent. Indeed, it is not obvious whether we can meet the conditions presented in Section 2.2. We leave such an interesting question for our future research.

#### REFERENCES

- A. Nedić and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," IEEE Transactions on Automatic Control, vol. 54, no. 1, pp. 48–61, 2009.
- T. T. Doan, S. T. Maguluri, and J. Romberg, "Distributed stochastic approximation for solving network optimization problems under random quantization," Available at: https://arxiv. org/abs/1810.11568, 2018.
- [3] A. Reisizadeh, A. Mokhtari, H. Hassani, and R. Pedarsani, "Quantized Decentralized Consensus Optimization," Available at: https://arxiv.org/pdf/1806.11536.pdf, 2018.

- [4] A. Nedić, A. Olshevsky, and M. G. Rabbat, "Network topology and communicationcomputation tradeoffs in decentralized optimization," *Proceedings of the IEEE*, vol. 106, no. 5, pp. 953–976, 2018.
- [5] J. Tsitsiklis, D. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," *IEEE Transactions on Automatic Control*, vol. 31, no. 9, pp. 803–812, 1986.
- [6] 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.
- [7] G. Qu and N. Li, "Harnessing smoothness to accelerate distributed optimization," IEEE Transactions on Control of Network Systems, no. 99, 2017.
- [8] A. Nedíc, A. Olshevsky, and W. Shi, "Achieving geometric convergence for distributed optimization over time-varying graphs," *SIAM Journal on Optimization*, vol. 27, no. 4, pp. 2597–2633, 2017.
- P. D. Lorenzo and G. Scutari, "NEXT: in-network nonconvex optimization," *IEEE Trans.* Signal and Information Processing over Networks, vol. 2, no. 2, pp. 120–136, 2016.
- [10] C. Xi, V. S. Mai, R. Xin, E. H. Abed, and U. A. Khan, "Linear convergence in optimization over directed graphs with row-stochastic matrices," *IEEE Transactions on Automatic Control*, vol. 63, no. 10, pp. 3558–3565, 2018.
- [11] C. Xi, R. Xin, and U. A. Khan, "ADD-OPT: Accelerated distributed directed optimization," *IEEE Transactions on Automatic Control*, vol. 63, no. 5, pp. 1329–1339, 2018.
- [12] T. C. Aysal, M. J. Coates, and M. G. Rabbat, "Distributed average consensus with dithered quantization," *IEEE Transactions on Signal Processing*, vol. 56, no. 10, pp. 4905–4918, 2008.
- [13] J. Li, G. Chen, Z. Wu, and X. He, "Distributed subgradient method for multiagent optimization with quantized communication," *Mathematical Methods in the Applied Sciences*, vol. 40, no. 4, pp. 1201–1213, 2016.
- [14] A. Nedić, A. Olshevsky, A. Ozdaglar, and J. N. Tsitsiklis, "Distributed subgradient methods and quantization effects," in 2008 47th IEEE Conference on Decision and Control, Dec 2008, pp. 4177–4184.
- [15] Y. Pu, M. N. Zeilinger, and C. N. Jones, "Quantization design for distributed optimization," IEEE Transactions on Automatic Control, vol. 62, no. 5, pp. 2107–2120, May 2017.
- [16] T. T. Doan, S. T. Maguluri, and J. Romberg, "On the convergence of distributed subgradient methods under quantization," in 2018 56th Annual Allerton Conference on Communication, Control, and Computing (Allerton), 2018, pp. 567–574.
- [17] P. Yi and Y. Hong, "Quantized subgradient algorithm and data-rate analysis for distributed optimization," *IEEE Transactions on Control of Network Systems*, vol. 1, no. 4, pp. 380– 392, 2014.
- [18] R. W. Brockett and D. Liberzon, "Quantized feedback stabilization of linear systems," IEEE Transactions on Automatic Control, vol. 45, no. 7, pp. 1279–1289, 2000.
- [19] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, "Qsgd: Communication-efficient sgd via gradient quantization and encoding," in Advances in Neural Information Processing Systems 30, 2017, pp. 1709–1720.
- [20] D. Alistarh, T. Hoefler, M. Johansson, N. Konstantinov, S. Khirirat, and C. Renggli, "The convergence of sparsified gradient methods," in Advances in Neural Information Processing Systems 31, 2018, pp. 5977–5987.
- [21] S. U. Stich, J.-B. Cordonnier, and M. Jaggi, "Sparsified sgd with memory," in Advances in Neural Information Processing Systems 31, 2018, pp. 4452–4463.
- [22] T. Wu, K. Yuan, Q. Ling, W. Yin, and A. H. Sayed, "Decentralized consensus optimization with asynchrony and delays," *IEEE Trans. Signal and Information Processing over Networks*, vol. 4, no. 2, pp. 293–307, 2018.
- [23] Y. Tian, Y. Sun, B. Du, and G. Scutari, "ASY-SONATA: achieving geometric convergence for distributed asynchronous optimization," CoRR, vol. abs/1803.10359, 2018.
- [24] T. T. Doan, C. L. Beck, and R. Srikant, "Convergence rate of distributed subgradient methods under communication delays," in *Proceedings of American Control Conference (ACC)*, 2018.
- [25] —, "On the convergence rate of distributed gradient methods for finite-sum optimization under communication delays," *Proceedings ACM Meas. Anal. Comput. Syst.*, vol. 1, no. 2, pp. 37:1–37:27, Dec. 2017.
- [26] G. Lan, S. Lee, and Y. Zhou, "Communication-Efficient Algorithms for Decentralized and Stochastic Optimization," Available at: https://arxiv.org/abs/1708.03543, 2017.
- [27] K. Scaman, F. Bach, S. Bubeck, L. Massoulié, and Y. T. Lee, "Optimal algorithms for nonsmooth distributed optimization in networks," in Advances in Neural Information Pro-

cessing Systems 31, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., 2018, pp. 2740–2749.

- [28] E. Wei and A. Ozdaglar, "Distributed alternating direction method of multipliers," in 2012 IEEE 51st IEEE Conference on Decision and Control (CDC), Dec 2012, pp. 5445–5450.
- [29] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends* in Machine Learning, vol. 3, no. 1, pp. 1–22, 2011.
- [30] W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin, "On the linear convergence of the ADMM in decentralized consensus optimization," *IEEE Trans. Signal Processing*, vol. 62, no. 7, pp. 1750–1761, 2014.
- [31] T. 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.
- [32] T. Chang, M. Hong, W. Liao, and X. Wang, "Asynchronous distributed ADMM for large-scale optimizationpart I: Algorithm and convergence analysis," *IEEE Transactions on Signal Processing*, vol. 64, no. 12, pp. 3118–3130, 2016.
- [33] J. Duchi, A. Agarwal, and M. Wainwright, "Dual averaging for distributed optimization: Convergence analysis and network scaling," *IEEE Transactions on Automatic Control*, vol. 57, no. 3, pp. 592–606, 2012.
- [34] K. Tsianos, S. Lawlor, and M. Rabbat, "Push-sum distributed dual averaging for convex optimization," in *Proceedings of the 51st IEEE Conference on Decision and Control (CDC)*, Hawaii, USA, Dec 2012.
- [35] T. T. Doan, S. Bose, D. H. Nguyen, and C. L. Beck, "Convergence of the iterates in mirror descent methods," *IEEE Control Systems Letters*, vol. 3, no. 1, pp. 114–119, Jan 2019.
- [36] S. Shalev-Shwartz, "Online learning and online convex optimization," Foundations and Trends in Machine Learning, vol. 4, pp. 107–194, 2012.
- [37] A. Nedić, A. Ozdaglar, and P. A. Parrilo, "Constrained consensus and optimization in multiagent networks," *IEEE Transactions on Automatic Control*, vol. 55, no. 4, pp. 922–938, 2010.
- [38] R. Horn and C. Johnson, Matrix Analysis. Cambridge, U.K.: Cambridge Univ. Press, 1985.
- [39] Y. Nesterov, Introductory Lectures on Convex Optimization: A Basic Course. Norwell, MA: Kluwer Academic Publishers, 2004.
- [40] T. Hastie, T. Tibshirani, and J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, 2nd ed., ser. Springer Series in Statistics. New York: Springe-Verlag, 2009.
- [41] S. Shalev-Shwartz and S. Ben-David, Understanding Machine Learning: From Theory to Algorithms, 1st ed. Cambridge University Press, 2014.
- [42] O. J. Karst, "Linear curve fitting using least deviations," Journal of the American Statistical Association, vol. 53, no. 281, pp. 118–132, 1958.

### Appendix A. Proofs of Results in Section 3.1 .

#### A.1. Proof of Lemma 3.1.

*Proof.* For convenience, we use  $\mathbf{w}_i(k)$  to denote  $\sum_{j \in \mathcal{N}_i} a_{ij} \mathbf{x}_j(k)$  only in this proof. Since  $\mathbf{x}_i(k) \in \mathcal{X}$ , by convexity of  $\mathcal{X}$ , we have that  $\mathbf{w}_i(k) \in \mathcal{X}$ . Recall that  $\Delta_i(k) = \mathbf{x}_i(k) - \mathbf{q}_i(k)$ . By the definition of the projection and using Eqs. (1.4) and (3.1) we obtain Eq. (3.4), i.e.,

$$\|\boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))\| = \|\mathbf{v}_{i}(k) - [\mathbf{v}_{i}(k)]_{\mathcal{X}}\|^{2} \leq \|\mathbf{v}_{i}(k) - \mathbf{w}_{i}(k)\|^{2}$$
$$= \left\|\Delta_{i}(k) - \sum_{j \in \mathcal{N}_{i}} a_{ij}\Delta_{j}(k) - \alpha(k)\mathbf{g}_{i}(\mathbf{x}_{i}(k))\right\|^{2}$$
$$\leq \left\|\Delta_{i}(k) - \sum_{j \in \mathcal{N}_{i}} a_{ij}\Delta_{j}(k)\right\| + L_{i}\alpha(k)$$
$$\leq \sum_{j \in \mathcal{N}_{i}} a_{ij} \|\Delta_{i}(k) - \Delta_{j}(k)\| + L_{i}\alpha(k).$$

Using the preceding inbequality also yields Eq. (3.5), i.e.,

$$\sum_{i=1}^{n} \|\boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))\|^{2} \leq 2 \sum_{i=1}^{n} \left\| \Delta_{i}(k) - \sum_{j \in \mathcal{N}_{i}} a_{ij} \Delta_{j}(k) \right\|^{2} + 2 \sum_{i=1}^{n} L_{i}^{2} \alpha^{2}(k)$$
$$\leq 2 \sum_{i=1}^{n} \sum_{j \in \mathcal{N}_{i}} a_{ij} \|\Delta_{i}(k) - \Delta_{j}(k)\|^{2} + 2L^{2} \alpha^{2}(k)$$
$$\leq 8 \sum_{i=1}^{n} \|\Delta_{i}(k)\|^{2} + 2L^{2} \alpha^{2}(k),$$

where the second inequality follows from using Jensen's inequality.

### A.2. Proof of Lemma 3.2.

*Proof.* For convenience let  $\mathbf{W} = \mathbf{I} - 1/n\mathbf{1}\mathbf{1}^T$  and  $\mathbf{Y}(k)$  be defined as

$$\mathbf{Y}(k) = \mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^T = \mathbf{W}\mathbf{X}(k).$$

Using A1 = 1 and Eqs. (3.2) and (3.3) we consider

$$\begin{aligned} \mathbf{Y}(k+1) &= \mathbf{X}(k+1) - \mathbf{1}\bar{\mathbf{x}}(k+1)^T \\ &= \mathbf{A}\mathbf{X}(k) + (\mathbf{I} - \mathbf{A})(\mathbf{X}(k) - \mathbf{Q}(k)) - \alpha(k)\mathbf{G}(\mathbf{X}(k)) - \mathbf{\Xi}(\mathbf{V}(k)) \\ &- \mathbf{1}\bar{\mathbf{x}}(k)^T + \frac{\alpha(k)}{n}\mathbf{1}\sum_{i=1}^n \mathbf{g}_i(\mathbf{x}_i(k))^T + \mathbf{1}\bar{\boldsymbol{\xi}}(k)^T \\ &= \mathbf{A}\mathbf{W}\mathbf{X}(k) + (\mathbf{I} - \mathbf{A})\Delta(k) - \alpha(k)\mathbf{W}\mathbf{G}(\mathbf{X}(k)) - \mathbf{W}\mathbf{\Xi}(\mathbf{V}(k)), \end{aligned}$$

which by taking the Frobenius norm on both sides yields

$$\begin{aligned} \|\mathbf{Y}(k+1)\| &= \|\mathbf{AWX}(k) + (\mathbf{I} - \mathbf{A})\Delta(k) - \alpha(k)\mathbf{WG}(\mathbf{X}(k)) - \mathbf{W\Xi}(\mathbf{V}(k))\| \\ &\leq \|\mathbf{AWX}(k)\| + \|(\mathbf{I} - \mathbf{A})\Delta(k)\| + \|\alpha(k)\mathbf{WG}(\mathbf{X}(k)) - \mathbf{W\Xi}(\mathbf{V}(k))\| \\ &\leq \sigma_2 \|\mathbf{WX}(k)\| + 2 \|\Delta(k)\| + \alpha(k) \|\mathbf{G}(\mathbf{X}(k))\| + \|\mathbf{\Xi}(\mathbf{V}(k))\|, \quad (A.1) \end{aligned}$$

where the last inequality is because the largest singular values of  $\mathbf{W}$  and  $\mathbf{A}$  are smaller than 1 and using the Courant-Fisher theorem [38], i.e.,

$$\|\mathbf{AWX}(k)\| = \left\|\mathbf{A}\left(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^{T}\right)\mathbf{X}(k)\right\| \le \sigma_{2}\left\|\left(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^{T}\right)\mathbf{X}(k)\right\| = \sigma_{2}\left\|\mathbf{WX}(k)\right\|.$$

First, using  $L = \sum_{i=1}^{n} L_i$  Eq. (1.4) gives

$$\|\mathbf{G}(\mathbf{X}(k))\| \le \sqrt{\sum_{i=1}^{n} \|\mathbf{g}_i(\mathbf{x}_i(k))\|^2} \le \sqrt{\sum_{i=1}^{n} L_i^2} \le L.$$
 (A.2)

Second, Eq. (3.5) yields

$$\|\mathbf{\Xi}(\mathbf{V}(k))\| \le 2\sqrt{2} \|\Delta(k)\| + \sqrt{2}L\alpha(k).$$
(A.3)

Thus using Eqs. (A.2) and (A.3) into Eq. (A.1) yields Eq. (3.6), i.e.,

$$\begin{aligned} \|\mathbf{Y}(k+1)\| &\leq \sigma_2 \|\mathbf{Y}(k)\| + 6\|\Delta(k)\| + 3L\alpha(k) \\ &\leq \sigma_2^{k+1} \|\mathbf{Y}(0)\| + 6\sum_{t=0}^k \sigma_2^{k-t} \|\Delta(t)\| + 3L\sum_{t=0}^k \sigma_2^{k-t}\alpha(t) \\ &= 6\sum_{t=0}^k \sigma_2^{k-t} \|\Delta(t)\| + 3L\sum_{t=0}^k \sigma_2^{k-t}\alpha(t), \end{aligned}$$

where the last equality is due to  $\mathbf{x}_i(0) = \mathbf{x}_j(0)$  for all  $i, j \in \mathcal{V}$ , implying  $\mathbf{Y}(0) = \mathbf{0}$ .  $\Box$ 

## A.3. Proof of Lemma 3.3.

*Proof.* For convenience, just in this proof, we define  $\mathbf{w}_i(k)$  as

$$\mathbf{w}_{i}(k) = \left[\sum_{j \in \mathcal{N}_{i}} a_{ij} \mathbf{q}_{j}(k)\right]_{\mathcal{X}} = \sum_{j \in \mathcal{N}_{i}} a_{ij} \mathbf{q}_{j}(k) - \mathbf{p}_{i}(k).$$

where  $\mathbf{p}_i(k) \triangleq \sum_{j \in \mathcal{N}_i} a_{ij} \mathbf{q}_j(k) - \mathbf{w}_i(k)$ . Using the definition of the projection and since  $\sum_{j \in \mathcal{N}_i} a_{ij} \mathbf{x}_j(k) \in \mathcal{X}$  we have

$$\|\mathbf{p}_{i}(k)\| = \left\| \sum_{j \in \mathcal{N}_{i}} a_{ij} \mathbf{q}_{j}(k) - \mathbf{w}_{i}(k) \right\| \leq \left\| \sum_{j \in \mathcal{N}_{i}} a_{ij} \mathbf{x}_{j}(k) - \sum_{j \in \mathcal{N}_{i}} a_{ij} \mathbf{q}_{j}(k) \right\|$$
$$\leq \sum_{j \in \mathcal{N}_{i}} a_{ij} \|\Delta_{j}(k)\|, \qquad (A.4)$$

which implies

$$\|\mathbf{p}_{i}(k) - \bar{\mathbf{p}}(k)\| \leq \|\mathbf{P}(k) - \mathbf{1}\bar{\mathbf{p}}(k)^{T}\| \leq \|\mathbf{P}(k)\|$$
$$\leq \sqrt{\sum_{i=1}^{n} \left(\sum_{j \in \mathcal{N}_{i}} a_{ij} \|\Delta_{j}(k)\|\right)^{2}} \leq \|\Delta(k)\|.$$
(A.5)

Moreover, since  $\Delta_i(k) = \mathbf{x}_i(k) - \mathbf{q}_i(k)$  we have

$$\|\mathbf{q}_{i}(k) - \bar{\mathbf{q}}(k)\| \le \|\mathbf{Q}(k) - \mathbf{1}\bar{\mathbf{q}}(k)^{T}\| \le \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| + \|\Delta(k)\|.$$
(A.6)

Using the triangle inequality we now consider

$$\|\mathbf{x}_{i}(k+1) - \mathbf{q}_{i}(k)\| \leq \|\mathbf{x}_{i}(k+1) - \mathbf{w}_{i}(k)\| + \|\mathbf{w}_{i}(k) - \bar{\mathbf{w}}(k)\| + \|\bar{\mathbf{w}}(k) - \mathbf{q}_{i}(k)\|.$$
(A.7)

We now provide an upper bound for each term on the right-hand side of Eq. (A.7). First, by the nonexpansiveness of the projection, and Eqs. (1.4) and (2.4) we have

$$\|\mathbf{x}_{i}(k+1) - \mathbf{w}_{i}(k)\| \le \|\Delta_{i}(k) - \alpha(k)\mathbf{g}_{i}(\mathbf{x}_{i}(k))\| \le \|\Delta_{i}(k)\| + L_{i}\alpha(k).$$
(A.8)

Second, using Eqs. (A.5) and (A.6) we have

$$\|\mathbf{w}_{i}(k) - \bar{\mathbf{w}}(k)\| = \left\| \sum_{j \in \mathcal{N}_{i}} a_{ij} \mathbf{q}_{j}(k) - \bar{\mathbf{q}}(k) - \mathbf{p}_{i}(k) + \bar{\mathbf{p}}(k) \right\|$$
$$\leq \sum_{j \in \mathcal{N}_{i}} a_{ij} \|\mathbf{q}_{j}(k) - \bar{\mathbf{q}}(k)\| + \|\mathbf{p}_{i}(k) - \bar{\mathbf{p}}(k)\|$$
$$\leq \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| + 2\|\Delta(k)\|.$$
(A.9)

Third, using Eq. (A.6) we consider

$$\|\bar{\mathbf{w}}(k) - \mathbf{q}_{i}(k)\| = \|\bar{\mathbf{q}}(k) - \mathbf{q}_{i}(k) + \bar{\mathbf{p}}(k)\| \le \|\mathbf{q}_{i}(k) - \bar{\mathbf{q}}(k)\| + \|\bar{\mathbf{p}}(k)\| \le \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| + 2\|\Delta(k)\|.$$
(A.10)

Substituting Eqs. (A.8)–(A.10) into Eq. (A.7) yields

$$\|\mathbf{x}_{i}(k+1) - \mathbf{q}_{i}(k)\| \le 2\|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| + 5\|\Delta(k)\| + L\alpha(k),$$

which by Eq. (3.6) gives

$$\begin{aligned} \|\mathbf{x}_{i}(k+1) - \mathbf{q}_{i}(k)\| &\leq 12 \sum_{t=0}^{k-1} \sigma_{2}^{k-1-t} \|\Delta(t)\| + 6L \sum_{t=0}^{k-1} \sigma_{2}^{k-1-t} \alpha(t) \\ &+ 5 \|\Delta(k)\| + L\alpha(k) \\ &\leq 12 \sum_{t=0}^{k} \sigma_{2}^{k-t} \|\Delta(t)\| + 6L \sum_{t=0}^{k} \sigma_{2}^{k-t} \alpha(t). \end{aligned}$$
(A.11)

Recall that  $\gamma = 48(2+L)/(1-\sigma_2)$  and

$$\mathcal{R}_{i}(k+1) \triangleq \left[\mathbf{q}_{i}(k) - \frac{\gamma}{2}\alpha(k)\mathbf{1}, \, \mathbf{q}_{i}(k) + \frac{\gamma}{2}\alpha(k)\mathbf{1}\right].$$

We now show that  $\mathbf{x}_i(k+1) \in \mathcal{R}_i(k+1)$  by induction. First, when k = 0 we have  $\|\Delta_i(0)\| = 0$  since  $\mathbf{x}_i(0) = \mathbf{q}_i(0)$  for all  $i \in \mathcal{V}$ . By definition, we have  $\mathbf{x}_i(0) \in \mathcal{R}_i(0)$ . Suppose it is true for some k > 0, that is,  $\mathbf{x}_i(k) \in \mathcal{R}_i(k)$ . We now show that  $\mathbf{x}_i(k+1) \in \mathcal{R}_i(k+1)$ . Indeed, since  $\alpha(k)$  is nonincreasing, by the definition of  $\mathcal{R}_i(k)$  we have  $\Delta_i(k)$  is nonincreasing and

$$\Delta_{i}(t) \leq \frac{\gamma \alpha\left(t\right)}{2^{b} - 1} \mathbf{1}, \qquad \forall i \in \mathcal{V}, \ t \in [0, k].$$
22

Using Assumption 2, i.e.,  $\sqrt{nd}\gamma/(2^b-1) \leq 1$ , we have  $\|\Delta(t)\| \leq \alpha(t)$  for all  $t \in [0,k]$ . Thus, by using  $\alpha(k) \leq \alpha(t) \leq \alpha(0) = 1$  for  $t \in [0,k]$ , Eq. (A.11) gives

$$\begin{aligned} \|\mathbf{x}_i(k+1) - \mathbf{q}_i(k)\| &\leq 12 \sum_{t=0}^k \sigma_2^{k-t} \alpha(t) + 6L \sum_{t=0}^k \sigma_2^{k-t} \alpha(t) \\ &= 6(2+L) \left( \sum_{t=0}^{\lfloor k/2 \rfloor} \sigma_2^{k-t} \alpha(t) + \sum_{t=\lceil k/2 \rceil}^k \sigma_2^{k-t} \alpha(t) \right) \\ &\leq 6(2+L) \left( \frac{\alpha(0)\sigma_2^{\lceil k/2 \rceil}}{1-\sigma_2} + \frac{\alpha(\lceil k/2 \rceil)}{1-\sigma_2} \right) \\ &\leq \frac{24(2+L)}{1-\sigma_2} \alpha(k) = \frac{\gamma}{2} \alpha(k), \end{aligned}$$

where in the last inequality is due to  $\sigma_2^k \leq \alpha(k), \alpha(0) = 1, \alpha(\lceil k/2 \rceil) \leq 2\alpha(k)$  since we only consider  $\alpha(k) = 1/(k+1)$  or  $\alpha(k) = 1/\sqrt{k+1}$ . This concludes our proof.  $\Box$ 

# A.4. Proof of Lemma 3.4.

Proof. First, Eq. (3.9) yields

$$\|\Delta(k)\| \le \frac{\sqrt{nd\gamma}}{2^b - 1}\alpha(k),$$

which using Eq. (3.6) gives

$$\begin{split} \|\mathbf{X}(k+1) - \mathbf{1}\bar{\mathbf{x}}(k+1)^T\| &\leq \frac{6\sqrt{nd\gamma}}{2^b - 1} \sum_{t=0}^k \sigma_2^{k-t} \alpha(t) + 3L \sum_{t=0}^k \sigma_2^{k-t} \alpha(t) \\ &\leq \left(\frac{6\sqrt{nd\gamma} + 3L2^b}{2^b - 1}\right) \left(\sum_{t=0}^{\lfloor k/2 \rfloor} \sigma_2^{k-t} \alpha(t) + \sum_{t=\lceil k/2 \rceil}^k \sigma_2^{k-t} \alpha(t)\right) \\ &\leq \left(\frac{6\sqrt{nd\gamma} + 3L2^b}{2^b - 1}\right) \left(\sum_{t=0}^{\lfloor k/2 \rfloor} \sigma_2^{k-t} + \alpha(\lceil k/2 \rceil) \sum_{t=\lfloor k/2 \rfloor}^k \sigma_2^{k-t}\right) \\ &\leq \left(\frac{6\sqrt{nd\gamma} + 3L2^b}{2^b - 1}\right) \left(\frac{1}{1 - \sigma_2} \sigma_2^{\lceil k/2 \rceil} + \frac{1}{1 - \sigma_2} \alpha(\lceil k/2 \rceil)\right), \end{split}$$

which since  $\lim_{k\to\infty} \alpha(k) = 0$  gives Eq. (3.10).

Suppose now that the condition (3.11) is held. Then, for some  $K \ge 0$  we have Eq. (3.12), i.e.,

$$\begin{split} \sum_{k=0}^{K} \alpha(k) \| \mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T} \| &\leq \left( \frac{6\sqrt{nd}\gamma + 3L2^{b}}{2^{b} - 1} \right) \sum_{k=0}^{K} \alpha(k) \sum_{t=0}^{k-1} \sigma_{2}^{k-1-t} \alpha(t) \\ &\leq \left( \frac{6\sqrt{nd}\gamma + 3L2^{b}}{2^{b} - 1} \right) \sum_{k=0}^{K} \sum_{t=0}^{k-1} \sigma_{2}^{k-1-t} \alpha^{2}(t) = \left( \frac{6\sqrt{nd}\gamma + 3L2^{b}}{2^{b} - 1} \right) \sum_{t=0}^{K} \alpha^{2}(t) \sum_{k=t+1}^{K} \sigma_{2}^{k} \\ &\leq \left( \frac{6\sqrt{nd}\gamma + 3L2^{b}}{(1 - \sigma_{2})(2^{b} - 1)} \right) \sum_{t=0}^{K} \alpha^{2}(t) \sum_{k=0}^{(3.11)} \infty. \end{split}$$

Suppose now that  $\alpha(k) = 1/\sqrt{k+1}$ . Then by the inequality above we have Eq. (3.13), i.e.,

$$\begin{split} \sum_{k=0}^{K} \alpha(k) \| \mathbf{X}(k) - \mathbf{1} \bar{\mathbf{x}}(k)^{T} \| &\leq \left( \frac{6\sqrt{nd}\gamma + 3L2^{b}}{(1 - \sigma_{2})(2^{b} - 1)} \right) \sum_{t=0}^{K} \frac{1}{t+1} \\ &\leq \left( \frac{6\sqrt{nd}\gamma + 3L2^{b}}{(1 - \sigma_{2})(2^{b} - 1)} \right) (1 + \ln(K + 1)), \end{split}$$

where we use the integral test in the last inequality to have

$$\sum_{t=0}^{K-1} \frac{1}{t+1} \le 1 + \int_0^K \frac{1}{t+1} dt \le 1 + \ln(K+1).$$

### A.5. Proof of Lemma 3.5.

*Proof.* Let  $\mathbf{x}^*$  be a solution of problem (1.1). For convenience, let  $\mathbf{r}(k) = \bar{\mathbf{x}}(k) - \mathbf{x}^*$ . First, recall from Eq. (3.5) that

$$\|\mathbf{\Xi}(\mathbf{V}(k))\|^2 \le 8\|\Delta(k)\|^2 + 2L^2\alpha^2(k).$$
(A.12)

Second, by the definition of the projection we have

$$\boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))^{T}(\mathbf{x}^{*}-\mathbf{v}_{i}(k)) \leq -\|\boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))\|^{2}.$$
(A.13)

We now use Eq. (3.3) to have

$$\|\mathbf{r}(k+1)\|^{2} = \left\| \bar{\mathbf{x}}(k) - \mathbf{x}^{*} - \frac{\alpha(k)}{n} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k)) - \bar{\mathbf{\xi}}(k) \right\|^{2}$$

$$= \|\mathbf{r}(k)\|^{2} - 2(\bar{\mathbf{x}}(k) - \mathbf{x}^{*})^{T} \bar{\mathbf{\xi}}(k) - \frac{2\alpha(k)}{n} (\bar{\mathbf{x}}(k) - \mathbf{x}^{*})^{T} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k))$$

$$+ \left\| \frac{\alpha(k)}{n} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k)) + \bar{\mathbf{\xi}}(k) \right\|^{2}$$

$$\leq \|\mathbf{r}(k)\|^{2} - 2(\bar{\mathbf{x}}(k) - \mathbf{x}^{*})^{T} \bar{\mathbf{\xi}}(k) - \frac{2\alpha(k)}{n} (\bar{\mathbf{x}}(k) - \mathbf{x}^{*})^{T} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k))$$

$$+ 2 \left\| \frac{\alpha(k)}{n} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k)) \right\|^{2} + 2 \left\| \bar{\mathbf{\xi}}(k) \right\|^{2}$$

$$\overset{(A.12)}{\leq} \|\mathbf{r}(k)\|^{2} - 2(\bar{\mathbf{x}}(k) - \mathbf{x}^{*})^{T} \bar{\mathbf{\xi}}(k) - \frac{2\alpha(k)}{n} (\bar{\mathbf{x}}(k) - \mathbf{x}^{*})^{T} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k))$$

$$+ \frac{16\Delta^{2}(k) + 6L^{2}\alpha^{2}(k)}{n}. \qquad (A.14)$$

We now analyze the second term on the right-hand side of Eqs. (A.14) by using Eqs.

(A.12) and (A.13)

$$-2(\bar{\mathbf{x}}(k) - \mathbf{x}^{*})^{T} \bar{\boldsymbol{\xi}}(k) = -\frac{2}{n} \sum_{i=1}^{n} \boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))^{T}(\bar{\mathbf{x}}(k) - \mathbf{x}^{*})$$

$$= -\frac{2}{n} \sum_{i=1}^{n} \boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))^{T}(\bar{\mathbf{x}}(k) - \mathbf{v}_{i}(k) + \mathbf{v}_{i}(k) - \mathbf{x}^{*})$$

$$\leq \frac{2}{n} \sum_{i=1}^{n} \|\boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))\| \|\bar{\mathbf{x}}(k) - \mathbf{v}_{i}(k)\| - \frac{2}{n} \sum_{i=1}^{n} \boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))^{T}(\mathbf{v}_{i}(k) - \mathbf{x}^{*})$$

$$\stackrel{(A.12)}{\leq} \frac{2(2\sqrt{2}\|\Delta(k)\| + \sqrt{2}L\alpha(k))}{n} \sum_{i=1}^{n} \|\bar{\mathbf{x}}(k) - \mathbf{v}_{i}(k)\| - \frac{2}{n} \sum_{i=1}^{n} \|\boldsymbol{\xi}_{i}(\mathbf{v}_{i}(k))\|^{2}$$

$$\leq \frac{2(2\sqrt{2}\|\Delta(k)\| + \sqrt{2}L\alpha(k))}{\sqrt{n}} \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\|$$

$$+ \frac{2(2\sqrt{2}\|\Delta(k)\| + \sqrt{2}L\alpha(k))(2\sqrt{n}\|\Delta(k)\| + L\alpha(k))}{n}, \quad (A.15)$$

where the last inequality is due to

$$\begin{split} &\sum_{i=1}^{n} \left\| \bar{\mathbf{x}}(k) - \mathbf{v}_{i}(k) \right\| \\ &\leq \sum_{i=1}^{n} \left\| \bar{\mathbf{x}}(k) - \sum_{j=1}^{n} a_{ij} \mathbf{x}_{j}(k) + \Delta_{i}(k) - \sum_{j \in \mathcal{N}_{i}} a_{ij} \Delta_{j}(k) - \alpha(k) \mathbf{g}_{i}(\mathbf{x}_{i}(k)) \right\| \\ &\leq \sqrt{n} \left\| \mathbf{X}(k) - \mathbf{1} \bar{\mathbf{x}}(k)^{T} \right\| + 2\sqrt{n} \left\| \Delta(k) \right\| + L\alpha(k), \end{split}$$

which uses the Jensen's inequality. Next, we analyze the third term on the right-hand side of Eq. (A.14)  $\,$ 

$$-\frac{2\alpha(k)}{n}(\bar{\mathbf{x}}(k) - \mathbf{x}^*)^T \sum_{i=1}^n \mathbf{g}_i(\mathbf{x}_i(k))$$

$$= -\frac{2\alpha(k)}{n} \sum_{i=1}^n \mathbf{g}_i(x_i(k))^T(\bar{\mathbf{x}}(k) - \mathbf{x}_i(k)) - \frac{2\alpha(k)}{n} \sum_{i=1}^n \mathbf{g}_i(\mathbf{x}_i(k))^T(\mathbf{x}_i(k) - \mathbf{x}^*)$$

$$\leq \frac{2L\alpha(k)}{n} \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^T\| - \frac{2\alpha(k)}{n} \sum_{i=1}^n \mathbf{g}_i(\mathbf{x}_i(k))^T(\mathbf{x}_i(k) - \mathbf{x}^*). \quad (A.16)$$

Substituting Eqs. (A.15) and (A.16) into (A.14) we obtain

$$\begin{aligned} \|\mathbf{r}(k+1)\|^{2} &\leq \|\mathbf{r}(k)\|^{2} - \frac{2\alpha(k)}{n} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k))^{T}(\mathbf{x}_{i}(k) - \mathbf{x}^{*}) \\ &+ \frac{2\left(2\sqrt{2}\|\Delta(k)\| + \sqrt{2}L\alpha(k)\right)}{\sqrt{n}} \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| \\ &+ \frac{2\left(2\sqrt{2}\|\Delta(k)\| + \sqrt{2}L\alpha(k)\right)\left(2\sqrt{n}\|\Delta(k)\| + L\alpha(k)\right)}{n} \\ &+ \frac{2L\alpha(k)}{n} \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| + \frac{16\Delta^{2}(k) + 6L^{2}\alpha^{2}(k)}{n}, \end{aligned}$$

which gives us Eq. (3.14), i.e.,

$$\begin{split} \|\mathbf{r}(k+1)\|^{2} &\leq \|\mathbf{r}(k)\|^{2} - \frac{2\alpha(k)}{n} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k))^{T}(\mathbf{x}_{i}(k) - \mathbf{x}^{*}) \\ &+ \frac{2\left(4\|\Delta(k)\| + 3L\alpha(k)\right)}{\sqrt{n}} \|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| \\ &+ \frac{32\|\Delta(k)\|^{2} + 10L^{2}\alpha^{2}(k) + 12L\alpha(k)\|\Delta(k)\|}{\sqrt{n}} \\ &\leq \|\mathbf{r}(k)\|^{2} - \frac{2\alpha(k)}{n} \sum_{i=1}^{n} \mathbf{g}_{i}(\mathbf{x}_{i}(k))^{T}(\mathbf{x}_{i}(k) - \mathbf{x}^{*}) \\ &+ \frac{2\left(4\sqrt{nd}\gamma + 3L2^{b}\right)}{\sqrt{n}(2^{b} - 1)} \alpha(k)\|\mathbf{X}(k) - \mathbf{1}\bar{\mathbf{x}}(k)^{T}\| \\ &+ \frac{2(4\sqrt{nd}\gamma + 3L2^{b})^{2}}{(2^{b} - 1)^{2}\sqrt{n}}\alpha^{2}(k), \end{split}$$

where the last inequality we use Eq. (3.9) to have

$$\|\Delta(k)\| \le \frac{\sqrt{nd\gamma}}{2^b - 1} \alpha(k).$$