

# Convergence rate analysis of a subgradient averaging algorithm for distributed optimisation with different constraint sets

Licio Romao, Kostas Margellos, Giuseppe Notarstefano, and Antonis Papachristodoulou

**Abstract**—We consider a multi-agent setting with agents exchanging information over a network to solve a convex constrained optimisation problem in a distributed manner. We analyse a new algorithm based on local subgradient exchange under undirected time-varying communication. First, we prove asymptotic convergence of the iterates to a minimum of the given optimisation problem for time-varying step-sizes of the form  $c(k) = \frac{\eta}{k+1}$ , for some  $\eta > 0$ . We then restrict attention to step-size choices  $c(k) = \frac{\eta}{\sqrt{k+1}}$ ,  $\eta > 0$ , and establish a convergence rate of  $\mathcal{O}\left(\frac{\ln(k)}{\sqrt{k}}\right)$  in objective value. Our algorithm extends currently available distributed subgradient/proximal methods by: (i) accounting for different constraint sets at each node, and (ii) enhancing the convergence speed thanks to a subgradient averaging step performed by the agents. A numerical example demonstrates the efficacy of the proposed algorithm.

## I. INTRODUCTION

We focus on distributed algorithms to solve convex optimisation problems. They are motivated by applications, such as sensor networks, robust estimation and source localisation [1], that require distribution of computational power and possibly data to alleviate the burden caused by data size. The main challenge is to devise fast and efficient algorithms that converge to an optimal solution of the centralised problem without requiring global information.

In the past decade, motivated by [2], [3], distributed optimisation has drawn the attention of the community because of its relevance to important real-world problems. Indeed, [2], [3] proposed a distributed projected subgradient algorithm that converges under time-varying network for problems with a common constraint set known by all agents. More recently, [4] extended these results to the time-varying directed case by relying on a push-sum consensus protocol [5]. They showed convergence rates of  $\mathcal{O}\left(\frac{\ln(k)}{\sqrt{k}}\right)$  for the function value at the running average of the local iterates. Another contribution in this direction was made by [6], which proposed a proximal algorithm in which the local iterates maintained by the agents converge to a point in the optimal set under time-varying network and for problems with different constraint sets.

L. Romao is supported by the Coordination for the Improvement of Higher Education Personnel (CAPES) - Brazil. The work of K. Margellos and A. Papachristodoulou is supported in part by the Engineering and Physical Sciences Research Council (EPSRC) UK under grants EP/P03277X/1 and EP/M002454/1, respectively. Giuseppe Notarstefano is supported by the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme (grant 638992-OPT4SMART).

L. Romao, K. Margellos and A. Papachristodoulou are with the Department of Engineering Science, University of Oxford, Parks Road, Oxford OX1 3PJ, UK. {licio.romao, kostas.margellos, antonis}@eng.ox.ac.uk.

G. Notarstefano is with the Department of Electrical, Electronic and Information Engineering, University of Bologna, Bologna, Italy. giuseppe.notarstefano@unibo.it.

A new research direction involves the use of gradient tracking, mainly because sharp convergence results can be obtained for directed and undirected communication networks [7], [8], [9]. To achieve this improved performance, each agent maintains an additional local variable that tracks asymptotically the (sub-)gradient of the global function. Agents then use this additional information, which provides a more accurate direction towards minimising the overall objective function, to update their estimate of the solution.

The contribution of this paper is twofold: 1) unlike the aforementioned literature, we propose a new algorithm based on subgradient averaging that can simultaneously cope with non-differentiable local objective functions, and different constraint sets, while accounting for a time-varying communication network. By showing convergence in iterates for a step-size of the form  $c(k) = \frac{\eta}{k+1}$ ,  $\eta > 0$ , we set a new framework accounting for the presence of different constraint sets and subgradient exchanges. As a consequence of this result, we expect faster practical convergence when compared to standard projected subgradient algorithms because we use an additional information that better approximate the subgradient of the global function; 2) We build upon the results of [6] and establish a convergence rate of  $\mathcal{O}\left(\frac{\ln(k)}{\sqrt{k}}\right)$ , when the step-size is  $c(k) = \frac{\eta}{\sqrt{k+1}}$ ,  $\eta > 0$ . Even though similar bounds have appeared in the literature, the present analysis offers the first convergence rate for the particular subgradient averaging scheme with the same rate as for standard distributed subgradient methods. We highlight that in our results we allow for different constraint sets, thus extending the scope of existing algorithms in the literature. Note that lifting constraints in the objective via characteristic functions, although possible, is not amenable to algorithms like [3], [10], [11], as this would render the subgradient of the resulting objective unbounded.

The paper is organised as follows. In Section II we present the problem statement, the main assumptions, as well as a description of the proposed algorithm. Section III contains the main results of this paper related to convergence in iterates and a convergence rate analysis as far as the optimal value is concerned. Section IV provides a numerical example to demonstrate the main algorithmic features of our scheme. Finally, some concluding remarks are provided in Section V. All omitted and simplified proofs can be found in [12].

## II. PROBLEM STATEMENT

### A. Problem set-up and Assumptions

Consider the following optimisation problem

$$\begin{aligned} & \underset{x}{\text{minimise}} && f(x) = \sum_{i=1}^m f_i(x) \\ & \text{subject to} && x \in \cap_{i=1}^m X_i \end{aligned} \quad (1)$$

where  $x \in \mathbb{R}^n$  is the global decision vector, and  $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$  and  $X_i \subset \mathbb{R}^n$ , for all  $i = 1, \dots, m$ , constitute the local objective function and constraint set for agent  $i$ , respectively. We suppose that each agent  $i$  possesses as private information the triple  $(x_i, f_i, X_i)$ , where the first component  $x_i$  is a local copy of the global variable  $x$ .

The goal is for all agents to agree on the local variables, that is,  $x_i = x^*$ , for all  $i = 1, \dots, m$ , where  $x^*$  belongs to the optimal set of (1), i.e., the subset of  $\mathbb{R}^n$  with the property that  $f(x^*) \leq f(x)$  for all  $x \in \cap_{i=1}^m X_i$ . In this paper, we suppose the following assumptions on  $f_i$  and  $X_i$ .

*Assumption 1:* (Convexity, compactness and non-emptiness of the interior)

- i) For all  $i = 1, \dots, m$ , the function  $f_i$  is proper and convex (see [13, Chapter 1] for a definition).
- ii) The set  $X_i \subset \mathbb{R}^n$  is compact and convex for all  $i = 1, \dots, m$ , and  $\cap_{i=1}^m X_i$  has a non-empty interior. We also assume that  $X_i \subset \cap_{i=1}^m \text{int}(\text{dom} f_i)$  for each  $i = 1, \dots, m$ , where  $\text{int}(A)$  indicates the interior of the set  $A$ .
- iii) The distance between the set  $\cup_{i=1}^m X_i$  and the complement of the interior of the domain of  $f$  (which is closed and convex) is strictly greater than zero.

As a consequence<sup>1</sup> of Assumption 1, we have that  $\cup_{x \in \text{conv}(\cup_{i=1}^m X_i)} \partial f(x)$  is a bounded set, that is,  $\|g\| \leq L$ , where  $g \in \partial f(x)$  for any  $x \in \cup_{i=1}^m X_i$ . This result is formally stated in the next Lemma.

*Lemma 1:* Under Assumption 1, we have that

- i) The set  $\text{conv}(\cup_{i=1}^m X_i)$  is compact, where  $\text{conv}(A)$  is the convex hull of the set  $A$ ;
- ii) The set  $\cup_{x \in \text{conv}(\cup_{i=1}^m X_i)} \partial f(x)$  is non-empty and bounded;
- iii) The function  $f$  is Lipschitz continuous over  $\cap_{i=1}^m X_i$ , i.e., there exists a positive scalar  $L$  such that

$$|f(x) - f(y)| \leq L\|x - y\|, \quad \forall x, y \in \cap_{i=1}^m X_i.$$

*Proof:* The proof is omitted for brevity (see [14, Theorem 24.7], or [12] for an alternative argument.) ■

## B. Proposed algorithm

The pseudocode of the proposed scheme is shown in Algorithm 1. We initialise each agent's local variable with an arbitrary  $x_i(0) \in X_i$ ,  $i = 1, \dots, m$ ; such points are not required to belong to  $\cap_{i=1}^m X_i$ .

At iteration  $k$ , agent  $i$  receives  $x_j$  from the neighbouring agents, and averages them through  $A(k)$ , which captures the communication network, to obtain  $z_i(k)$ . Here we represent the element of the  $j$ -th row and  $i$ -th column of matrix  $A(k)$  by  $[A(k)]_j^i$ . Agent  $i$  then calculates a subgradient,  $g_i$ , of its own objective function evaluated at  $z_i(k)$  and sends this information back to its neighbours. In the sequel, agent  $i$  averages the received  $g_j(z_j(k))$  in order to compose a proxy for a subgradient of  $f(x)$  (Step 3), called  $\tilde{z}_i(k)$ . Finally, at Step 4, agents use variables  $\tilde{z}_i(k)$  and  $z_i(k)$  to update their local estimates by projecting  $z_i(k) - c(k)\tilde{z}_i(k)$  onto the local set. Indeed, note that Step 4 can be rewritten as

$$x_i(k+1) = \mathcal{P}_{X_i}[z_i(k) - c(k)\tilde{z}_i(k)]$$

where  $\mathcal{P}_{X_i}$  denotes the projection operator onto the set  $X_i$ .

<sup>1</sup>A thorough discussion of Assumption 1 is given in the extended version of this paper [12]. Note that the bound on  $\|g\|$  coincides with the Lipschitz constant in Lemma 1, item iii).

---

## Algorithm 1 Proposed Scheme

---

**Require:**  $x_i(0)$ ,  $i = 1, \dots, m$

1: **while** Until convergence **do**

2:  $z_i(k) = \sum_{j=1}^m [A(k)]_j^i x_j(k), \quad \forall i = 1, \dots, m$

3:  $\tilde{z}_i(k) = \sum_{j=1}^m [A(k)]_j^i g_j(z_j(k)), \quad \forall i = 1, \dots, m$

4:  $x_i(k+1) = \text{argmin}_{\xi \in X_i} \tilde{z}_i(k)^T \xi + \frac{1}{2c(k)} \|z_i(k) - \xi\|_2^2, \quad \forall i = 1, \dots, m$

5:  $k \leftarrow k + 1$

6: **end while**

---

We now characterise  $A(k)$  that encodes the network in Algorithm 1. To this end, let  $\mathcal{G}(k) = (\mathcal{N}, \mathcal{E}(k))$  be a undirected graph, where  $\mathcal{N} = \{1, \dots, m\}$  is the number of agents and  $\mathcal{E}(k) \subset \mathcal{N} \times \mathcal{N}$  is the set of edges at iteration  $k$ , that is, if node  $(j, i) \in \mathcal{E}(k)$  then node  $j$  sends information to node  $i$  at iteration  $k$ . We associate the time-varying matrix  $A(k)$  to the edge set  $\mathcal{E}(k)$ , with  $[A(k)]_j^i \neq 0$  if  $(j, i) \in \mathcal{E}(k)$  at time  $k$ . As the graph is undirected, matrix  $A(k)$  can be chosen to be symmetric. We also define the graph  $\mathcal{G}_\infty = (\mathcal{N}, \mathcal{E}_\infty)$ , in which  $(j, i) \in \mathcal{E}_\infty$  if agent  $j$  communicates with agent  $i$  infinitely many times. Then, we impose the following assumption.

*Assumption 2:* (Network Properties)

- i) The graph  $(\mathcal{N}, \mathcal{E}_\infty)$  is strongly connected. Moreover, there exists a uniform upper bound on the communication time for all  $(j, i) \in \mathcal{E}_\infty$ .
- ii) There exists an  $\eta \in (0, 1)$  such that  $[A(k)]_j^i \geq \eta$  and that if  $[A(k)]_j^i > 0$  then we have  $[A(k)]_j^i \geq \eta$ , for all  $k \in \mathbb{N}$  and for all  $i, j = 1, \dots, m$ .
- iii) Matrix  $A(k)$  is doubly stochastic.

These are standard hypotheses in the distributed optimisation literature (see [2], [3], [10], [6] for more details).

The presented analysis is divided into two parts. First, we prove asymptotic convergence of the local variables,  $x_i$ , to some optimal solution of the centralised problem counterpart under square-summable but not summable step-sizes, e.g.,  $c(k) = \frac{\eta}{k+1}$ ,  $\eta > 0$ . We then show convergence rates of  $\mathcal{O}\left(\frac{\ln(k)}{\sqrt{k}}\right)$  in terms of the function value for the time-varying step-sizes of the form  $c(k) = \frac{\eta}{\sqrt{k+1}}$ ,  $\eta > 0$ .

## III. ALGORITHM ANALYSIS

### A. Convergence in iterates

In this subsection, we impose the following assumption on the step-size.

*Assumption 3:* (Non-increasing, square-summable step-size) Let  $(c(k))_{k \in \mathbb{N}}$  be the sequence of step-sizes adopted in Step 4 of Algorithm 1. We impose that

- i)  $c(k) \geq 0$ , and  $c(k) \geq c(r)$ , for all  $k, r \in \mathbb{N}$  with  $r \geq k$ ,
- ii)  $\sum_{k=1}^{\infty} c(k) = \infty$  and  $\sum_{k=1}^{\infty} c(k)^2 < \infty$ .

A sequence that satisfies Assumption 3 is  $c(k) = \frac{\eta}{k+1}$ , for some  $\eta > 0$ . Assumption 3 is necessary to prove one of the main results of this paper, namely, the asymptotic convergence for the sequences  $(x_i(k))_{k \in \mathbb{N}}$ , for all  $i = 1, \dots, m$ , to a point in the optimal set of (1).

*Theorem 1:* Let  $(x_i(k))_{k \in \mathbb{N}}$  be the sequences generated by Algorithm 1, for all  $i = 1, \dots, m$ . Under Assumptions 1-3, we have that for some minimiser  $x^*$  in the optimal set of (1),

$$\lim_{k \rightarrow \infty} \|x_i(k) - x^*\| = 0, \quad \forall i = 1, \dots, m.$$

*Proof:* See Appendix for the main steps of the proof and [12] for a complete proof. ■

Theorem 1 extends the result in [6] by allowing an agent to communicate subgradient information to neighbouring agents, a feature that can speed up practical convergence.

### B. Convergence in value and convergence rate

We impose now the following assumption on the step-size.

*Assumption 4:* The sequence  $(c(k))_{k \in \mathbb{N}}$  used in Step 4 of Algorithm (1) is  $c(k) = \frac{\eta}{\sqrt{k+1}}$ , for some  $\eta > 0$ .

Our convergence rate results build on the running average of the iterates generated by Algorithm 1, that is, the sequence

$$\hat{x}_i(k+1) = \frac{c(k+1)x_i(k+1) + S(k)\hat{x}_i(k)}{S(k+1)}, \quad (2)$$

where  $S(k) = \sum_{r=1}^k c(r)$ , and  $(x_i(k))_{k \in \mathbb{N}}$  for all  $i = 1, \dots, m$  are the sequences generated by Algorithm 1, with the initial conditions  $\hat{x}_i(0) = x_i(0)$  and  $S(0) = 1$ . By rewriting expression (2) as

$$\hat{x}_i(k) = \frac{1}{S(k)} \sum_{r=1}^k c(r)x_i(r),$$

we observe that the running average can be interpreted as a convex combination of previous iterates. The next theorem establishes a convergence rate for the function value along the running average defined in (2).

*Theorem 2:* Consider the running average defined in (2). Under Assumptions 1, 2, and 4 the following inequality holds for all  $k \in \mathbb{N} \setminus \{0\}$

$$\sum_{i=1}^m f_i(\hat{x}_i(k)) - f(x^*) \leq B_1 \frac{1}{\sqrt{k}} + B_2 \frac{\ln(k)}{\sqrt{k}}. \quad (3)$$

where  $B_1, B_2 > 0$  and defined in the Appendix.

*Proof:* See Appendix for a sketch of the proof. A self-contained argument of this result is presented in [12]. ■

Theorem 2 asserts convergence of the function value along the running average  $\hat{x}_i(k)$ , i.e., all limit point of  $(\hat{x}_i(k))_{k \in \mathbb{N}}$  are optimal. We point out that the result of Theorem 2 further extends the work presented in [6] not only by allowing agents to communicate their (sub-)gradients, but by also unveiling how to adapt the proof line in that paper to come up with convergence results that recover traditional rates for distributed subgradient methods.

## IV. NUMERICAL EXAMPLES

We now demonstrate the results through a numerical example. We consider problem (1) in which the functions  $f_i(x)$  are given by

$$f_i(x) = \max \left\{ |x^{(1)}|, \max_{2 \leq \ell \leq n} |x^{(\ell)} - (i+1)x^{(\ell-1)}| \right\},$$

where  $x^{(\ell)}$ ,  $\ell = 1, \dots, n$ , represents the  $\ell$ -th component of the vector<sup>2</sup>  $x$ . This example was analysed in [15], and was originally adapted from [16] to the distributed case. We consider the case where  $n = 20$  and  $m = 12$ . Note that the optimal solution and optimal value for this problem are  $x^* = 0$  and  $f(x^*) = 0$ , respectively.

We consider a time-invariant undirected network (notice that the relevant matrices do not depend on the iteration index  $k$ ) whose topology is given by a line graph. Given the topology, we generate a doubly stochastic  $A$  such that,

$$[A]_j^i = \frac{1}{1 + \max\{\mathcal{N}_i, \mathcal{N}_j\}}, \quad i \neq j,$$

where  $\mathcal{N}_i$  and  $\mathcal{N}_j$  are the number of neighbours of agent  $i$  and  $j$ , respectively. The diagonal elements of  $A$  are defined as  $[A]_i^i = 1 - \sum_{j \neq i} [A]_j^i$ . Note that matrix  $A$  is symmetric and doubly stochastic, as the network is undirected. Time-varying networks could be accommodated by the proposed algorithm.

Agents' constraint sets are different, and for each one we assumed that the constraint set is encoded by componentwise upper and lower limits on  $x$ . These limits were randomly generated from a uniform distribution. In our first numerical investigation, we set the step-sizes to be  $c(k) = \frac{1}{k+1}$ , aligned with Assumption 3. To investigate the statement of Theorem 1, we monitor the evolution

$$\text{Res}_x(k) = \sum_{i=1}^{12} \|x_i(k) - x^*\|$$

for  $k = 1, \dots, 10,000$  iterations. This result is shown in Figure 1, where the solid blue line corresponds to the iterates generated by Algorithm 1, initialised at the optimal solution  $x^* = 0$ . Observe that the iterates do not stay at the optimal solution as the function is not differentiable at the origin, implying that there exists a nonzero subgradient such that the iterate sequence escapes from the optimal solution. However, we also observe that after some initial perturbation, the iterates are steered back towards the optimal solution, thus supporting the results presented in Theorem 1.

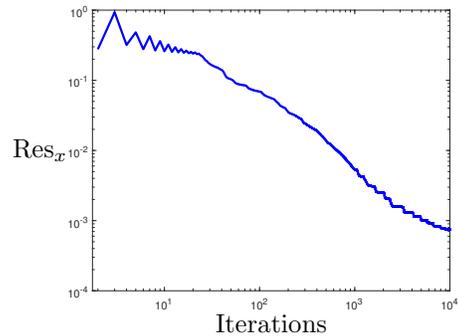


Fig. 1. Evolution of  $\text{Res}_x(k)$  for Algorithm 1 (solid, blue line). Both axes are in logarithmic scale.

In the second part of our numerical investigation, we choose the time-varying step-size according to Assumption 4,

<sup>2</sup>Variable  $x^{(\ell)}$  should not be related to  $x_i$  which corresponds to a local copy of  $x$  maintained by agent  $i$ , rather than to a particular component.

i.e.,  $c(k) = \frac{100}{\sqrt{k+1}}$ . To investigate the statement of Theorem 2, we monitor the evolution of

$$\text{Res}_f(k) = \sum_{i=1}^{12} f_i(\hat{x}_i(k)) - f(x^*),$$

where  $\hat{x}_i(k)$  is defined in (2), for  $k = 1, \dots, 100,000$  iterations. We initialized Algorithm 1 with  $x_i^{(\ell)}(0) = 0.1$  for  $\ell = 1, \dots, 19$ , and  $x_i^{(20)}(0) = 1$ , for all  $i = 1, \dots, 12$ . The results are illustrated in Figure 2. The theoretical bound  $\mathcal{O}\left(\frac{\ln(k)}{\sqrt{k}}\right)$  is depicted as the dashed-dot black line. The solid blue curve is the sequence  $\sum_{i=1}^m f_i(\hat{x}_i(k))$  for the iterates generated by Algorithm 1. Observe that the results comply with the theoretical bound of Theorem 2.

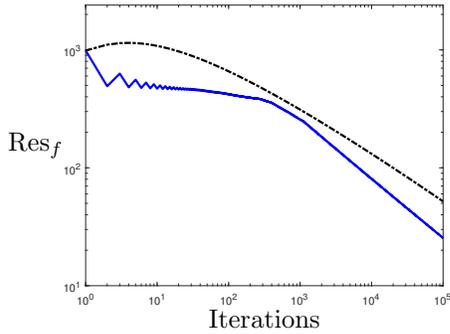


Fig. 2. Evolution of  $\text{Res}_f(k)$  for Algorithm 1 (solid, blue line). The solid black line represents an estimate of the theoretical bound (up to a constant factor) predicted by means of Theorem 2. Both axes are in logarithmic scale.

## V. CONCLUSION

We proposed a new algorithm based on subgradient averaging and proved two results for this scheme: (1) we have shown asymptotic convergence to the optimal set of the centralised problem for a time-varying step-size of the form  $c(k) = \frac{\eta}{k+1}, \eta > 0$ ; (2) for time-varying step-sizes of the form  $c(k) = \frac{\eta}{\sqrt{k+1}}, \eta > 0$ , we established a convergence rate of  $\mathcal{O}\left(\frac{\ln(k)}{\sqrt{k}}\right)$  as far as the function value of the running average of the local iterates is concerned. With these results, we recovered standard convergence rates of distributed subgradient methods; however, we extended them to the more general case in which agents are allowed to have their own constraint set. A numerical example has been presented to demonstrate the obtained results.

## APPENDIX

### A. Auxiliary results and proofs of Section III-A

Let

$$v(k) = \frac{1}{m} \sum_{i=1}^m x_i(k), \quad (4)$$

be the average of the estimates at time  $k$ . Since  $v(k)$  might not necessarily belong to the feasible set  $\cap_{i=1}^m X_i$ , we define

$$\bar{v}(k) = \frac{\rho}{\epsilon(k) + \rho} v(k) + \frac{\epsilon(k)}{\epsilon(k) + \rho} \bar{x}, \quad (5)$$

where  $\bar{x}$  is a point in the interior of the feasible set (which is non-empty by Assumption 1), and  $\epsilon(k) =$

$\sum_{i=1}^m \text{dist}(v(k), X_i)$ . As shown in [3], the point  $\bar{v}(k)$  is in  $\cap_{i=1}^m X_i$  for all  $k \in \mathbb{N}$ . Define  $e_i(k+1) = x_i(k+1) - z_i(k)$ , and note that Step 2 of Algorithm 1 can be written as

$$x_i(k+1) = \sum_{j=1}^m [A(k)]_j^i x_j(k) + e_i(k+1), \quad (6)$$

which can be interpreted as perturbed consensus protocol.

*Lemma 2:* The following relations hold.

- i)* Let  $(x_i(k))_{k \in \mathbb{N}}$  for all  $i = 1, \dots, m$  be the sequences generated by Algorithm 1, and  $(v(k))_{k \in \mathbb{N}}$  and  $(\bar{v}(k))_{k \in \mathbb{N}}$  be defined as in (4) and (5), respectively. Then, under Assumption 1,

$$\sum_{i=1}^m \|x_i(k+1) - \bar{v}(k)\| \leq \mu \sum_{i=1}^m \|x_i(k) - v(k)\|,$$

where  $\mu = \frac{2}{\rho} mD + 1$ , and  $D$  is the diameter of the set  $\cup_{i=1}^m X_i$  (which is well-defined by Lemma 1, item *i*).

- ii)* Let  $(x_i(k))_{k \in \mathbb{N}}$  for all  $i = 1, \dots, m$  and  $(v(k))_{k \in \mathbb{N}}$  be as in item *i*), and consider the definition of the error  $e_i(k)$ . Then, under Assumption 2, we have that

$$\begin{aligned} \|x_i(k+1) - v(k+1)\| &\leq \lambda q^k \sum_{j=1}^m \|x_j(0)\| + \|e_i(k+1)\| \\ &+ \sum_{r=0}^{k-1} \lambda q^{k-r-1} \sum_{j=1}^m \|e_j(r+1)\| + \frac{1}{m} \sum_{j=1}^m \|e_j(k+1)\|, \end{aligned}$$

where  $\lambda = 2(1 + \eta^{-(m-1)T}) / (1 - \eta^{(m-1)T}) \in \mathbb{R}_+$  and  $q = (1 - \eta^{(m-1)T})^{\frac{1}{(m-1)T}} \in (0, 1)$ , holds for all  $k \in \mathbb{N}$  and for all  $i = 1, \dots, m$ , with  $T$  being the uniform bound of Assumption 2, item *i*).

- iii)* Let  $(c(k))_{k \in \mathbb{N}}$  be a non-increasing and non-negative sequence, and  $\bar{L}$  be a positive scalar. If Assumption 2 holds, then for all  $k \in \mathbb{N}$  we have that

$$\begin{aligned} 2\bar{L} \sum_{k=0}^N c(k) \sum_{i=1}^m \|x_i(k+1) - \bar{v}(k+1)\| \\ < \beta_1 \sum_{k=0}^N \sum_{i=1}^m \|e_i(k+1)\|^2 + \beta_2 \sum_{k=0}^N c(k)^2 + \beta_3, \end{aligned}$$

where  $\beta_1 \in (0, 1)$ , and  $\beta_2$  and  $\beta_3$  are positive constants.

*Proof:* The proof of *i)* is presented in [6, Lemma 1]. For *ii)*, see [6, Lemma 2]. Finally, the proof of *iii)* follows the line of [6, Lemma 3]. ■

Notice that in Lemma 2, item *iii)*, we can choose any  $\beta_1 \in (0, 1)$ , at the price of modifying the values of  $\beta_2$  and  $\beta_3$ . For the presented analysis, the specific positive values for  $\beta_2$  and  $\beta_3$  are irrelevant.

Item *ii)* of the following lemma is a novel derivation, allowing the auxiliary sequences  $\alpha_1(k)$  and  $\alpha_2(k)$  to be iteration dependent. This is instrumental for proving Theorem 2.

*Lemma 3:* Let  $(x_i(k))_{k \in \mathbb{N}}, (z_i(k))_{k \in \mathbb{N}}$  and  $(\tilde{z}_i(k))_{k \in \mathbb{N}}, i = 1, \dots, m$ , be the sequences generated by Algorithm 1, and  $x^*$  by any point in the set of optimal solutions of problem (1). Suppose Assumptions 1 and 2 hold. Then,

- i)* For all  $k \in \mathbb{N}$  we have that

$$2c(k) \sum_{i=1}^m \tilde{z}_i(k)^T (x_i(k+1) - x^*) + \sum_{i=1}^m \|e_i(k+1)\|^2$$

$$+ \sum_{i=1}^m \|x_i(k+1) - x^*\|^2 \leq \sum_{i=1}^m \|x_i(k) - x^*\|^2. \quad (7)$$

ii) For any  $\beta_1 \in (0, 1)$ , there exist sequences  $(\alpha_1(k))_{k \in \mathbb{N}}$  and  $(\alpha_2(k))_{k \in \mathbb{N}}$  such that  $1 - \beta_1 - \alpha_1(k) - \alpha_2(k) \geq 0$  for all  $k \in \mathbb{N}$  and that

$$\begin{aligned} & 2 \sum_{k=0}^N c(k) \sum_{i=1}^m (f_i(\bar{v}(k+1)) - f_i(x^*)) \\ & + \sum_{k=0}^N (1 - \alpha_1(k) - \alpha_2(k) - \beta_1) \sum_{i=1}^m \|e_i(k+1)\|^2 \\ & + \sum_{k=0}^N \sum_{i=1}^m \|x_i(k+1) - x^*\|^2 \leq \sum_{k=0}^N \sum_{i=1}^m \|x_i(k) - x^*\|^2 \\ & + \sum_{k=0}^N \left( mL^2 \frac{\alpha_1(k) + \alpha_2(k)}{\alpha_1(k)\alpha_2(k)} + \beta_2 \right) c(k)^2 + \beta_3. \end{aligned}$$

*Proof:* The proof of *i*) is omitted for brevity, as it follows from the arguments in [6, Lemma 5]. Paper [12] also contains the derivation of this result.

To prove *ii*) we use *i*). Indeed, consider the first term of the left-hand side in inequality (7), and rewrite it as

$$\begin{aligned} & \underbrace{2c(k) \sum_{i=1}^m \tilde{z}_i(k)^T (x_i(k+1) - \bar{v}(k+1))}_{V_1} \\ & + \underbrace{2c(k) \sum_{i=1}^m \tilde{z}_i(k)^T (\bar{v}(k+1) - x^*)}_{V_2}, \end{aligned} \quad (8)$$

by adding and subtracting  $\bar{v}(k+1)$ . Now, let us consider the terms of the right hand-side of (8) separately. First,

$$V_1 \geq -2c(k)L \sum_{i=1}^m \|x_i(k+1) - \bar{v}(k+1)\|, \quad (9)$$

by Cauchy-Schwarz, triangle inequality, and  $L = \max_{\xi \in \cup_{j=1}^m X_j} \|g_j(\xi)\|$ , which is well-defined by Lemma 1. Second, we use the definition of  $\tilde{z}_i(k)$  – Step 3 in Algorithm 1 – into  $V_2$  to obtain (via double stochasticity of  $A$ )

$$V_2 = 2c(k) \sum_{i=1}^m g_i(z_i(k))^T (\bar{v}(k+1) - x^*). \quad (10)$$

Moreover, if we add and subtract  $x_i(k+1)$  and  $z_i(k)$  for all  $i = 1, \dots, m$  into the right-hand side of (10) we obtain

$$\begin{aligned} & \underbrace{2c(k) \sum_{i=1}^m g_i(z_i(k))^T (\bar{v}(k+1) - x_i(k+1))}_{V_{2,1}} \\ & + \underbrace{2c(k) \sum_{i=1}^m g_i(z_i(k))^T (x_i(k+1) - z_i(k))}_{V_{2,2}} \\ & + \underbrace{2c(k) \sum_{i=1}^m g_i(z_i(k))^T (z_i(k) - x^*)}_{V_{2,3}}. \end{aligned} \quad (11)$$

Let us focus on the right-hand side of (11). For the first term,

$$V_{2,1} \geq -2c(k)L \sum_{i=1}^m \|\bar{v}(k+1) - x_i(k+1)\|, \quad (12)$$

by Cauchy-Schwarz. As for the middle term, we have that

$$V_{2,2} \geq -\alpha_1(k) \sum_{i=1}^m \|e_i(k+1)\|^2 - m \frac{L^2}{\alpha_1(k)} c(k)^2 \quad (13)$$

by using Cauchy-Schwarz in the term  $V_{2,2}$  and then applying the definition  $e_i(k)$  given in (6) followed by the relation  $2xy \leq x^2 + y^2$  with  $x = \frac{L}{\sqrt{\alpha_1(k)}} c(k)$  and  $y = \sqrt{\alpha_1(k)} \|e_i(k+1)\|$  for some  $\alpha_1(k) \in (0, 1)$  for all  $k \in \mathbb{N}$ . Similarly, the right-most term of (11) yields

$$\begin{aligned} & V_{2,3} \geq -\alpha_2(k) \sum_{i=1}^m \|e_i(k+1)\|^2 \\ & - m \frac{L^2}{\alpha_2(k)} c(k)^2 - 2c(k)L \sum_{i=1}^m \|x_i(k+1) - \bar{v}(k+1)\| \\ & + 2c(k) \sum_{i=1}^m (f_i(\bar{v}(k+1)) - f_i(x^*)) \end{aligned} \quad (14)$$

for some sequence  $\alpha_2(k) \in (0, 1)$  for all  $k \in \mathbb{N}$ . The details that led to inequality (14) resemble the ones in (13), and are omitted for brevity. Substituting inequalities (9), (12), (13) and (14) into (7), and then using the result of Lemma 2, item *iii*), to the term involving  $\sum_{i=1}^m \|x_i(k+1) - \bar{v}(k+1)\|$  with  $\bar{L} = 3L$ , we obtain inequality of item *ii*). This concludes the proof of the Lemma. See [12] for the omitted steps. ■

Two immediate consequences of Lemma 3 are presented in the following proposition.

*Proposition 1:* Consider the result of Lemma 3, item *ii*), and suppose Assumptions 1–3 hold. Then

- i*) We have that  $\sum_{k=0}^{\infty} \sum_{i=1}^m \|e_i(k)\|^2 < \infty$ ;
- ii*) The error sequence  $(e_i(k))_{k \in \mathbb{N}}$  converges to zero for all  $i = 1, \dots, m$ ;
- iii*) For all  $i = 1, \dots, m$ ,

$$\lim_{k \rightarrow \infty} \|x_i(k) - v(k)\| = 0.$$

*Proof:* See [6, Propositions 2, 3], or [12]. ■

To prove Theorem 1, we leverage on a deterministic version of the supermartingale theorem.

*Lemma 4 ([17]):* Given non-negative scalar sequences  $(\ell(k))_{k \in \mathbb{N}}$ ,  $(u(k))_{k \in \mathbb{N}}$  and  $(\zeta(k))_{k \in \mathbb{N}}$  that obey the recursion

$$\ell(k+1) \leq \ell(k) - u(k) + \zeta(k).$$

If  $\sum_{k=1}^{\infty} \zeta(k) < \infty$ , then the sequence  $(\ell(k))_{k \in \mathbb{N}}$  converges and the sequence  $(u(k))_{k \in \mathbb{N}}$  is summable.

*Proof of Theorem 1*

In the view of item *ii*) in Lemma 3, fix a  $\beta_1 \in (0, 1)$  and choose  $\alpha_1(k) = \alpha_2(k) = \alpha$  such that  $(1 - 2\alpha - \beta_1) > 0$ . We now substitute inequalities (9), (12)–(14) into the inequality of Lemma 3, item *i*), to obtain the inequality of Lemma 4 with<sup>3</sup>

$$\ell(k) = \sum_{i=1}^m \|x_i(k) - x^*\|_2^2, \quad u(k) = 2c(k)(f(\bar{v}(k+1)) - f(x^*))$$

<sup>3</sup>The assumptions of Lemma 4 are satisfied because the sequence of step-sizes is square-summable under Assumption 3 and because of the relation given by Lemma 2, item *iii*), and the result of Proposition 1, item *i*).

$$\zeta(k) = \frac{2mL^2}{\alpha} c(k)^2 + 6c(k)L \sum_{i=1}^m \|x_i(k+1) - \bar{v}(k+1)\|, \quad (15)$$

As a consequence of Lemma 4, we have that the sequence  $\sum_{i=1}^m \|x_i(k) - x^*\|$  converges and that  $\sum_{k=0}^{\infty} 2c(k)(f(\bar{v}(k+1)) - f(x^*)) < \infty$ . Moreover, using the arguments presented in [12], we can further conclude that, in fact, sequence  $(\|x_i(k) - x^*\|)_{k \in \mathbb{N}}$  converges to zero, thus concluding the proof. ■

### B. Proofs of Section III-B

We only present here a sketch of the proof of Theorem 2. The reader is referred to [12] for a detailed argument. We start the proof by writing

$$\begin{aligned} \sum_{i=1}^m f_i(\hat{x}_i(k+1)) - f(x^*) &\leq f(\hat{v}(k+1)) - f(x^*) \\ &+ L \sum_{i=1}^m \|\hat{x}_i(k+1) - \hat{v}(k+1)\|, \quad (16) \end{aligned}$$

which follows from the relation in Lemma 1, item *iii*), and by defining  $\hat{v}(k)$  similarly as  $\hat{x}_i(k)$  in Theorem 2 from the sequence  $(\bar{v}(k))_{k \in \mathbb{N}}$ .

We split the argument into two parts: we first consider that (17) and (18) are satisfied for positive constants  $d_1, d_2, d_3$  and  $d_4$ , and prove the statement of Theorem 2. Then we return to (17) and (18), and prove the existence of such constants. To this end, consider

$$f(\hat{v}(k+1)) - f(x^*) \leq d_1 \frac{1}{S(k+1)} + d_2 \frac{\sum_{r=0}^k c(r)^2}{S(k+1)} \quad (17)$$

$$L \sum_{i=1}^m \|\hat{x}_i(k+1) - \hat{v}(k+1)\| \leq \frac{d_3}{S(k+1)} + d_4 \frac{\sum_{r=0}^k c(r)^2}{S(k+1)}. \quad (18)$$

Notice the following bounds on  $S(k+1)$  and  $\sum_{r=0}^k c(r)^2$ :

$$S(k+1) = \sum_{r=1}^{k+1} \frac{1}{\sqrt{r+1}} \geq \int_2^{k+3} \frac{1}{\sqrt{x}} dx \geq 0.5858\sqrt{k+1}, \quad (19)$$

$$\sum_{r=0}^k c(r)^2 = \sum_{r=0}^k \frac{1}{r+1} \leq \int_1^{k+1} \frac{1}{x} dx + 1 \leq \ln(k+1) + 1. \quad (20)$$

The result of the Theorem 2 would then follow by substituting (19) and (20) into (17) and (18) with constants  $B_1 = \sum_{i=1}^4 \frac{d_i}{\nu}$  and  $B_2 = \frac{d_2}{\nu} + \frac{d_4}{\nu}$ . This concludes the proof of Theorem 2, provided inequalities (17) and (18) hold.

Observe that since  $(\bar{v}(k))_{k \in \mathbb{N}} \subset \cap_{i=1}^m \text{int}(X_i)$  and that  $\hat{v}(k)$  is a convex combination of  $\bar{v}(k)$ , the relation given by (18) gives us the rate with which the generate solution becomes feasible.

The main steps of the second part of the proof are omitted for brevity, see [12] for the complete argument. However, we mention here that to prove (17) we invoke Lemma 3, item *ii*), with  $\alpha_1(k) = \alpha_2(k) = \alpha(k)$ , where  $\alpha(k) = a \left(1 - \frac{1}{\sqrt{k+1}}\right)$ , with  $a = (1 - \beta_1)/2$ . After some algebraic manipulations, we show that inequality (17) holds with constants

$$d_1 = 4mD^2 + \beta_3, \quad d_2 = \beta_2 + \frac{8mL^2}{a},$$

where  $D$  is defined as in Lemma 2, item *i*), and  $\beta_2$  and  $\beta_3$  are the constants obtained from Lemma 2, item *iii*). Similarly, we use some facts about convolution of sequences to prove inequality (18) with constants

$$\begin{aligned} d_3 &= 2mDL\mu c(1) \left(1 + \frac{m\lambda}{2(1-q)}\right) \\ &+ L\mu \left(1 + \frac{m\lambda}{2(1-q)}\right) \frac{4mD^2 + \beta_3}{1 - \beta_1 - 2\alpha}, \end{aligned}$$

$$d_4 = L\mu \left(1 + \frac{m\lambda}{2(1-q)}\right) \left(1 + \frac{1}{1 - \beta_1 - 2\alpha} \left(mL^2 \frac{2}{\alpha} + \beta_2\right)\right),$$

where  $\alpha \in (0, 1)$ . ■

### REFERENCES

- [1] M. G. Rabbat and R. Nowak, "Distributed Optimization in Sensor Networks," in *Inf. Process. Sens. Networks*, 2004, pp. 20–27.
- [2] A. Nedić and A. Ozdaglar, "Distributed Subgradient Methods for Multi-Agent Optimization," *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [3] A. Nedić, A. Ozdaglar, and P. A. Parrilo, "Constrained Consensus and Optimization in Multi-Agent Networks," *IEEE Trans. Automat. Contr.*, vol. 55, no. 4, pp. 922–938, 2010.
- [4] A. Nedić and A. Olshevsky, "Distributed optimization over time-varying directed graphs," *IEEE Trans. Automat. Contr.*, vol. 60, no. 3, pp. 601–615, 2015.
- [5] D. Kempe, A. Dobra, and J. Gehrke, "Gossip-Based Computation of Aggregate Information David," in *IEEE Symp. Found. Comput. Sci.* IEEE Computer Society, 2003, p. 482.
- [6] K. Margellos, A. Falsone, S. Garatti, and M. Prandini, "Distributed Constrained Optimization and Consensus in Uncertain Networks via Proximal Minimization," *IEEE Trans. Automat. Contr.*, vol. 63, no. 5, pp. 1372–1387, 2018.
- [7] W. Shi, Q. Ling, G. Wu, and W. Yin, "EXTRA: An Exact First-Order Algorithm for Decentralized Consensus Optimization," *SIAM J. Optim.*, vol. 25, no. 2, pp. 944–966, 2015.
- [8] A. Nedić, A. Olshevsky, and W. Shi, "Achieving Geometric Convergence for Distributed Optimization over Time-Varying Graphs," *SIAM J. Optim.*, vol. 55, no. 2, pp. 664–690, 2017.
- [9] G. Qu and N. Li, "Harnessing smoothness to accelerate distributed optimization," *IEEE Trans. Control Netw. Syst.*, vol. 5, no. 3, pp. 1245–1260, 2018.
- [10] J. C. Duchi, A. Agarwal, and M. J. Wainwright, "Dual averaging for distributed optimization: Convergence analysis and network scaling," *IEEE Trans. Automat. Contr.*, vol. 57, no. 3, pp. 592–606, 2012.
- [11] K. I. Tsianos, S. Lawlor, and M. G. Rabbat, "Push-Sum Distributed Dual Averaging for convex optimization," *2012 IEEE 51st IEEE Conf. Decis. Control*, pp. 5453–5458, 2012.
- [12] L. Romao, K. Margellos, G. Notarstefano, and A. Papachristodoulou, "Subgradient averaging for multi-agent optimisation with different constraint sets," *Tech. Report, ArXiv*, 2019. [Online]. Available: <http://arxiv.org/abs/1909.04351>
- [13] D. P. Bertsekas, *Convex Optimization Theory*. Athena Scientific, 2009.
- [14] R. T. Rockafellar, *Convex Analysis*. Princeton University Press, 1972.
- [15] S. Liang, L. Wang, and G. Yin, "Distributed quasi-monotone subgradient algorithm for nonsmooth convex optimization over directed graphs," *Automatica*, vol. 101, pp. 175–181, 2019.
- [16] Y. Nesterov and V. Shikhman, "Quasi-monotone Subgradient Methods for Nonsmooth Convex Minimization," *J. Optim. Theory Appl.*, vol. 165, no. 3, pp. 917–940, 2015.
- [17] D. P. Bertsekas, A. Nedić, and A. Ozdaglar, *Convex Analysis and Optimization*. Athena Scientific, 2003.