

# Systematic Analysis of Distributed Optimization Algorithms over Jointly-Connected Networks

Bryan Van Scoy

Laurent Lessard

To Appear, IEEE Conference on Decision and Control, 2020

## Abstract

We consider the distributed optimization problem, where a group of agents work together to optimize a common objective by communicating with neighboring agents and performing local computations. For a given algorithm, we use tools from robust control to systematically analyze the performance in the case where the communication network is *time-varying*. In particular, we assume only that the network is jointly connected over a finite time horizon (commonly referred to as  $B$ -connectivity), which does *not* require connectivity at each time instant. When applied to the distributed algorithm DIGing, our bounds are orders of magnitude tighter than those available in the literature.

## 1 Introduction

Many recent and emerging applications in multi-agent systems require groups of agents to cooperatively solve problems. Examples of agents include computing nodes, robots, or mobile sensors connected in a network. In this paper, we consider the *distributed optimization problem*, where each agent  $i \in \{1, \dots, n\}$  has a local function  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ , and agents cooperate to minimize the sum of the functions over all agents,

$$\underset{y \in \mathbb{R}^d}{\text{minimize}} \sum_{i=1}^n f_i(y). \quad (1)$$

Each agent is capable of evaluating its local gradient  $\nabla f_i$ , communicating information with neighboring agents, and performing local computations. This problem is relevant in applications such as large-scale machine learning [2], distributed spectrum sensing [1], and sensor networks [7].

We are interested in cases where the communication network among agents is *time varying*, which occurs in mobile agents with range-limited communication and systems with noisy and unreliable communication.

---

B. Van Scoy and L. Lessard were both with the University of Wisconsin–Madison, Madison, WI 53706, USA at the time of initial submission. B. Van Scoy is now with the Department of Electrical and Computer Engineering at Miami University, Oxford, OH 45056, USA, and L. Lessard is now with the Department of Mechanical and Industrial Engineering at Northeastern University, Boston, MA 02115, USA.  
bvanscoy@miamioh.edu, l.lessard@northeastern.edu

This material is based upon work supported by the National Science Foundation under Grants No. 1750162 and 1936648.

In the past several years, numerous algorithms have been proposed for distributed optimization (see [8] and the reference therein). While some algorithms have been studied in the time-varying scenario [5, 9], the analysis is typically performed on a case-by-case basis, resulting in lengthy convergence proofs and conservative bounds. On the other hand, the recent work [8] provides a systematic framework for deriving convergence bounds for a large class of distributed algorithms; this approach yields a straightforward comparison between various algorithms, but the analysis requires the network to be connected at *every* iteration, which is often unrealistic.

In this work, we provide a systematic analysis in the time-varying scenario for distributed algorithms whose state update is synchronous, homogeneous across agents, time invariant, and linear in the state, local gradient, and a weighted average among neighbors. Unlike [8], we only assume that the union of every  $B$  consecutive networks is connected; this assumption is common in the consensus literature and is often referred to as  $B$ -connectivity [3, 4]. Our main contributions are the following.

- **Generality.** Our analysis applies to a large class of distributed algorithms, allowing for straightforward comparisons without deriving lengthy convergence proofs for each algorithm individually.
- **Tightness.** For the gradient tracking algorithm DIGing [5, 6], we improve on existing bounds. While the available bounds scale with the number of agents  $n$  (and become vacuous as  $n \rightarrow \infty$ ), our bounds are orders of magnitude tighter and independent of  $n$ .

In Section 2, we state our assumptions on the local functions and communication network, as well as describe the class of algorithms considered. We then describe our analysis and present our main result in Section 3. We conclude with a case study for DIGing in Section 4, where we compare our results with those in the literature.

**Notation.** We use  $\mathbf{1}$  and  $\mathbf{0}$  to denote the  $n \times 1$  vectors of all ones and zeros, and  $\mathbb{R}^{m \times m}$  ( $\mathbb{S}^m$ ) to denote the set of  $m \times m$  real (symmetric) matrices. A matrix  $A \in \mathbb{S}^m$  is denoted positive (semi)definite by  $A \succ 0$  ( $A \succeq 0$ ). We denote the vertical concatenation of a list of matrices or vectors by  $\text{vcat}(A_1, \dots, A_n)^\top = [A_1^\top \ \dots \ A_n^\top]$ . Subscripts  $i$  and  $j$  refer to agents, and index  $k$  denotes the discrete time index. For a signal  $x_i(k)$  on agent  $i$  at time  $k$ , we denote the aggregation over all agents as

$$x(k) = \text{vcat}(x_1(k), \dots, x_n(k)).$$

## 2 Problem setup

We now discuss the objective functions, communication networks, and algorithms that we consider in this paper.

### 2.1 Objective function

We assume that the objective function has the form (1), where each agent  $i$  can evaluate its local gradient  $\nabla f_i$ . Furthermore, we assume that the gradient of each local function satisfies the following *sector bound*.

**Assumption 1.** There exist parameters  $0 < m \leq L$  such that each local function  $f_i$  is continuously differentiable and satisfies the inequality

$$\begin{aligned} & (\nabla f_i(y) - \nabla f_i(y_{\text{opt}}) - m(y - y_{\text{opt}}))^\top \\ & \quad \times (\nabla f_i(y) - \nabla f_i(y_{\text{opt}}) - L(y - y_{\text{opt}})) \leq 0 \end{aligned}$$

for all  $y \in \mathbb{R}^d$ , where  $y_{\text{opt}} \in \mathbb{R}^d$  is the optimizer of (1).

**Remark.** One way to satisfy Assumption 1 is if each  $\nabla f_i$  is  $L$ -Lipschitz continuous and each  $f_i$  is  $m$ -strongly convex, though in general, Assumption 1 is much weaker.

The condition ratio  $\kappa := L/m$  captures how much the curvature of the objective function varies. If  $f_i$  is twice differentiable, then  $\kappa$  is an upper bound on the condition number of the Hessian  $\nabla^2 f_i$ . In general, as  $\kappa \rightarrow \infty$ , the functions become poorly conditioned and are therefore more difficult to optimize using first-order methods.

### 2.2 Communication network

We represent the communication network among agents as a time-varying directed graph. Each agent corresponds to a node in the graph, and a directed edge from node  $i$  to node  $j$  indicates that agent  $i$  sends information to agent  $j$ . Each agent processes the communicated data by computing a weighted sum of the information from its neighbors (the set of agents from which it receives information). We characterize this diffusion process by a *gossip matrix*  $W(k) \in \mathbb{R}^{n \times n}$ , where the discrete time index  $k$  denotes the iteration of the algorithm. We make the following assumptions on the gossip matrices.

**Assumption 2.** The set of gossip matrices  $\{W(k)\}_{k=0}^\infty$  satisfies the following properties at each iteration  $k$ .

1. **Graph sparsity:**  $W_{ij}(k) = 0$  if agent  $i$  does not receive information from agent  $j$  at time  $k$ .
2. **Weight-balanced:**  $W(k)\mathbf{1} = W(k)^\top \mathbf{1} = \mathbf{1}$ .
3. **Spectrum property:**  $\|\frac{1}{n}\mathbf{1}\mathbf{1}^\top - W(k)\|_2 \leq 1$ .
4. **Joint-spectrum property:** There exists a positive integer  $B$  and a scalar  $\sigma \in [0, 1)$ , called the *spectral gap*, such that

$$\left\| \frac{1}{n}\mathbf{1}\mathbf{1}^\top - \prod_{\ell=k}^{k+B-1} W(\ell) \right\|_2^{1/B} \leq \sigma.$$

**Remark.** Our assumption on the gossip matrices does not require the graph to be connected at each iteration if  $B > 1$  and is common in the consensus and distributed optimization literature; see [4, 5].

### 2.3 Algorithm

We now describe a broad class of algorithms that may be used to have the group of agents solve the distributed optimization problem (1), where each agent may perform local computations and communicate with neighboring agents. At each iteration  $k$ , each agent  $i$  has a local state variable  $x_i(k) \in \mathbb{R}^{s \times d}$  that it updates as follows:

$$\begin{bmatrix} x_i(k+1) \\ y_i(k) \\ z_i(k) \end{bmatrix} = \begin{bmatrix} A & B_u & B_v \\ C_y & D_{yu} & D_{yv} \\ C_z & D_{zu} & D_{zv} \end{bmatrix} \begin{bmatrix} x_i(k) \\ u_i(k) \\ v_i(k) \end{bmatrix}, \quad (2a)$$

$$u_i(k) = \nabla f_i(y_i(k)), \quad (2b)$$

$$v_i(k) = \sum_{j=1}^n W_{ij}(k) z_j(k). \quad (2c)$$

The local gradient  $\nabla f_i$  is evaluated<sup>1</sup> at  $y_i(k) \in \mathbb{R}^{1 \times d}$  in (2b), and the quantity  $z_i(k) \in \mathbb{R}^{c \times d}$  is transmitted to neighboring agents in (2c). We also allow for linear state-input invariants to be enforced with

$$\sum_{j=1}^n (F_x x_j(k) + F_u u_j(k)) = 0. \quad (2d)$$

Such invariants typically arise from requiring a particular initialization for the algorithm.

The matrices  $A \in \mathbb{R}^{s \times s}$ ,  $D_{yu} \in \mathbb{R}^{1 \times 1}$ , and  $D_{zv} \in \mathbb{R}^{c \times c}$  are square with the other matrices having compatible dimensions. Here,  $s$  is the number of local states on each agent and  $c$  is the number of variables that each agent communicates with its neighbors at each iteration.

We want each agent's trajectory of algorithm (2) to converge to the optimizer of the distributed optimization problem (1), that is,  $y_i(k) \rightarrow y_{\text{opt}}$  for all  $i \in \{1, \dots, n\}$ . To obtain this, we need (i) the algorithm to have a fixed point corresponding to the optimal solution, and (ii) the trajectory to converge to this fixed point. Existence of such a fixed point places requirements on the structure of the algorithm, which we characterize in the following proposition; we prove the result in the Appendix.

**Proposition 1.** An algorithm of the form (2) has a fixed point such that  $y_1^* = \dots = y_n^*$  and  $\sum_{i=1}^n \nabla f_i(y_i^*) = 0$  for any local functions and any set of weight-balanced gossip matrices if and only if the algorithm matrices satisfy

$$\text{null} \left( \begin{bmatrix} A - I & B_v \\ C_z & D_{zv} - I \\ F_x & 0 \end{bmatrix} \right) \cap \text{row}([C_y \ D_{yv}]) \neq \{0\} \quad (3a)$$

$$\text{and} \quad \begin{bmatrix} B_u \\ D_{yu} \\ D_{zu} \end{bmatrix} \in \text{col} \left( \begin{bmatrix} A - I \\ C_y \\ C_z \end{bmatrix} \right). \quad (3b)$$

<sup>1</sup>We interpret the gradient  $\nabla f_i$  as a mapping from  $\mathbb{R}^{1 \times d}$  to  $\mathbb{R}^{1 \times d}$ .



### 3.2 Main result

Given an algorithm of the form (2), parameters  $(m, L)$  from Assumption 1 and  $(\sigma, B)$  from Assumption 2, and a prospective convergence rate  $\rho \in (0, 1)$ , we now construct the consensus and disagreement LMIs used to find the matrices  $P$  and  $Q$  in (5) and then state our main result.

**Map from basis to iterates.** To construct the LMIs, we first define a set of matrices that map the basis

$$\eta_i(k) = \text{vcat}(x_i(k), u_i(k), \dots, u_i(k+B-1), v_i(k), \dots, v_i(k+B-1), w_i(k+2), \dots, w_i(k+B)) \quad (6)$$

to the corresponding iterates of the algorithm. The basis has size  $b \times d$ , where  $b = s - c + B(2c + 1)$  (recall that  $s$  is the number of states on each agent and  $c$  is the number of variables communicated per iteration). In particular, we define the sets of matrices

$$\begin{aligned} \mathbf{u}(\ell) &\in \mathbb{R}^{1 \times b} & \ell \in \{0, \dots, B-1\} \\ \mathbf{v}(\ell) &\in \mathbb{R}^{c \times b} & \ell \in \{0, \dots, B-1\} \\ \mathbf{w}(\ell) &\in \mathbb{R}^{c \times b} & \ell \in \{0, \dots, B\} \\ \mathbf{x}(\ell) &\in \mathbb{R}^{s \times b} & \ell \in \{0, \dots, B\} \\ \mathbf{y}(\ell) &\in \mathbb{R}^{1 \times b} & \ell \in \{0, \dots, B-1\} \\ \mathbf{z}(\ell) &\in \mathbb{R}^{c \times b} & \ell \in \{0, \dots, B-1\} \end{aligned}$$

such that the concatenated matrix

$$\text{vcat}(\mathbf{x}(0), \mathbf{u}(0), \dots, \mathbf{u}(B-1), \mathbf{v}(0), \dots, \mathbf{v}(B-1), \mathbf{w}(2), \dots, \mathbf{w}(B))$$

is the  $b \times b$  identity matrix,  $\mathbf{w}(0) = \mathbf{z}(0)$  and  $\mathbf{w}(1) = \mathbf{v}(0)$ , and the matrices satisfy the algorithm update

$$\begin{bmatrix} \mathbf{x}(\ell+1) \\ \mathbf{y}(\ell) \\ \mathbf{z}(\ell) \end{bmatrix} = G \begin{bmatrix} \mathbf{x}(\ell) \\ \mathbf{u}(\ell) \\ \mathbf{v}(\ell) \end{bmatrix}, \quad \ell \in \{0, \dots, B-1\},$$

where  $G$  is defined in (4). These matrices are constructed such that multiplying each matrix on the right by the basis vector (6) yields the corresponding iterate.

**Lyapunov function.** Using these matrices, we define the matrices mapping the basis vector to the current and next state of the Lyapunov function as

$$\begin{aligned} \boldsymbol{\xi} &= \text{vcat}(\mathbf{x}(0), \mathbf{u}(0), \dots, \mathbf{u}(B-2), \mathbf{v}(0), \dots, \mathbf{v}(B-2)) \\ \boldsymbol{\xi}_+ &= \text{vcat}(\mathbf{x}(1), \mathbf{u}(1), \dots, \mathbf{u}(B-1), \mathbf{v}(1), \dots, \mathbf{v}(B-1)) \end{aligned}$$

with dimensions  $a \times b$ , where  $a = s + (c + 1)(B - 1)$ .

**Consensus LMI.** Let the matrix  $\Psi$  be a basis for

$$\begin{aligned} &\text{null} \left( \begin{bmatrix} I_B \otimes F_x & I_B \otimes F_u \end{bmatrix} \begin{bmatrix} \text{vcat}(\mathbf{x}(0), \dots, \mathbf{x}(B-1)) \\ \text{vcat}(\mathbf{u}(0), \dots, \mathbf{u}(B-1)) \end{bmatrix} \right) \\ &\cap \text{null} \left( \text{vcat}(\mathbf{v}(0) - \mathbf{z}(0), \dots, \mathbf{v}(B) - \mathbf{z}(B)) \right) \\ &\cap \text{null} \left( \text{vcat}(\mathbf{w}(0) - \mathbf{w}(1), \dots, \mathbf{w}(B) - \mathbf{w}(B+1)) \right). \end{aligned}$$

The consensus LMI is then

$$0 \succeq X(P, \lambda) \quad (7a)$$

$$0 \prec P \quad (7b)$$

$$0 \leq \lambda(\ell) \quad \ell \in \{0, \dots, B-1\} \quad (7c)$$

with variables  $P \in \mathbb{S}^a$  and  $\lambda(\ell) \in \mathbb{R}$ , where the symmetric matrix  $X$  is given by

$$\begin{aligned} X(P, \lambda) &= \Psi^T \left( \boldsymbol{\xi}_+^T P \boldsymbol{\xi}_+ - \rho^2 (\boldsymbol{\xi}^T P \boldsymbol{\xi}) \right. \\ &\quad \left. + \sum_{\ell=0}^{B-1} \lambda(\ell) \begin{bmatrix} \mathbf{y}(\ell) \\ \mathbf{u}(\ell) \end{bmatrix}^T M_0 \begin{bmatrix} \mathbf{y}(\ell) \\ \mathbf{u}(\ell) \end{bmatrix} \right) \Psi \end{aligned}$$

$$\text{with } M_0 = \begin{bmatrix} -2mL & L+m \\ L+m & -2 \end{bmatrix}.$$

**Disagreement LMI.** The disagreement LMI is

$$0 \succeq Y(Q, R, S, \lambda) \quad (8a)$$

$$0 \prec Q \quad (8b)$$

$$0 \preceq R \quad (8c)$$

$$0 \preceq S(\ell) \quad \ell \in \{0, \dots, B-1\} \quad (8d)$$

$$0 \leq \lambda(\ell) \quad \ell \in \{0, \dots, B-1\} \quad (8e)$$

with variables  $Q \in \mathbb{S}^a$ ,  $R \in \mathbb{S}^c$ ,  $S(\ell) \in \mathbb{S}^{2c}$ , and  $\lambda(\ell) \in \mathbb{R}$ , where the  $b \times b$  symmetric matrix  $Y$  is given by

$$\begin{aligned} Y(Q, R, S, \lambda) &= \boldsymbol{\xi}_+^T Q \boldsymbol{\xi}_+ - \rho^2 (\boldsymbol{\xi}^T Q \boldsymbol{\xi}) \\ &\quad + \sum_{\ell=0}^{B-1} \lambda(\ell) \begin{bmatrix} \mathbf{y}(\ell) \\ \mathbf{u}(\ell) \end{bmatrix}^T M_0 \begin{bmatrix} \mathbf{y}(\ell) \\ \mathbf{u}(\ell) \end{bmatrix} \\ &\quad + \begin{bmatrix} \mathbf{w}(0) \\ \mathbf{w}(B) \end{bmatrix}^T (M_1 \otimes R) \begin{bmatrix} \mathbf{w}(0) \\ \mathbf{w}(B) \end{bmatrix} \\ &\quad + \sum_{\ell=0}^{B-1} \begin{bmatrix} \mathbf{z}(\ell) \\ \mathbf{w}(\ell) \\ \mathbf{v}(\ell) \\ \mathbf{w}(\ell+1) \end{bmatrix}^T (M_2 \otimes S(\ell)) \begin{bmatrix} \mathbf{z}(\ell) \\ \mathbf{w}(\ell) \\ \mathbf{v}(\ell) \\ \mathbf{w}(\ell+1) \end{bmatrix} \end{aligned}$$

$$\text{with } M_1 = \begin{bmatrix} \sigma^{2B} & 0 \\ 0 & -1 \end{bmatrix} \text{ and } M_2 = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}.$$

**Remark.** The consensus and disagreement LMIs are coupled through the variable  $\lambda(\ell)$ .

We now use the consensus and disagreement LMIs to state our main result, which characterizes the worst-case convergence rate of algorithm (2); we prove the result in the Appendix.

**Theorem 1 (Main result).** Consider the optimization problem (1) solved using a distributed algorithm of the form (2) that satisfies the fixed-point conditions (3), and suppose that Assumptions 1–2 hold. If the consensus and disagreement LMIs in (7) and (8) are feasible for some scalar  $\rho > 0$ , then for any initial condition, there exists a constant  $\gamma > 0$  such that

$$\|x_i(k) - x_i^*\| \leq \gamma \rho^k \quad (9)$$

for all agents  $i \in \{1, \dots, n\}$  and all iterations  $k \geq 0$ , where  $x_i^*$  is a fixed point corresponding to the optimal solution of (1).

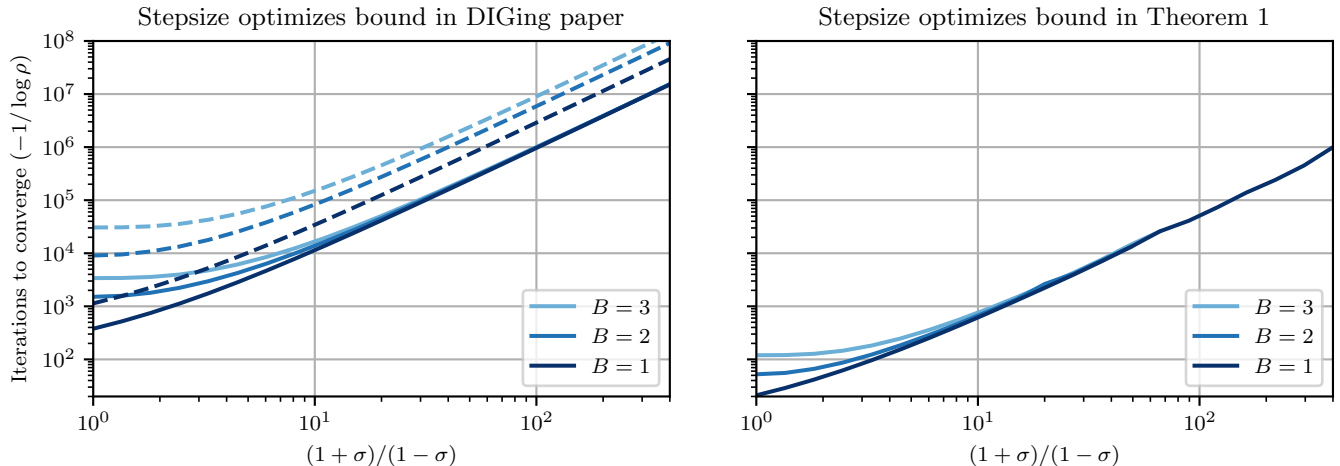


Figure 2: Upper bound on the number of iterations to converge for DIGing as a function of the spectral gap  $\sigma$  with condition ratio  $\kappa = 10$  and connectivity parameter  $B \in \{1, 2, 3\}$  using Theorem 1 (solid lines) and the bound from the original DIGing paper [5, Theorem 3.14] (dashed lines) with the stepsize  $\alpha$  chosen to optimize the bound in the DIGing paper (left) and Theorem 1 (right). Note that the bound in the DIGing paper depends on the number of agents; we use  $n = 2$  for the plot on the left, and the bound is vacuous for the stepsizes used in the plot on the right.

#### 4 Case study: DIGing

To illustrate our results, we applied our analysis to the gradient tracking algorithm DIGing [5,6], which has been analyzed under the same assumptions<sup>2</sup>.

The DIGing algorithm is given by the recursion

$$\begin{aligned} x(k+1) &= W(k)x(k) - \alpha y(k) \\ y(k+1) &= W(k)y(k) + \nabla f(x(k+1)) - \nabla f(x(k)) \end{aligned}$$

with initial condition  $y_i(0) = \nabla f_i(x_i(0))$  and stepsize  $\alpha$ .

If we define the state as  $\text{vcat}(x_i(k), y_i(k), \nabla f_i(x_i(k)))$ , then DIGing is equivalent to our algorithm form (2) with

$$\begin{bmatrix} A & B_u & B_v \\ C_y & D_{yu} & D_{yv} \\ C_z & D_{zu} & D_{zv} \\ F_x & F_u & \end{bmatrix} = \begin{bmatrix} 0 & -\alpha & 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & -\alpha & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & & \end{bmatrix}.$$

From the dimensions of the matrices, we see that each agent has  $s = 3$  state variables and communicates  $c = 2$  variables to neighbors at each iteration.

We compare our convergence bound from Theorem 1 with that from the original DIGing paper [5, Thm. 3.14] in Figure 2. In the plot on the left, we use the stepsize

$$\alpha = \frac{1.5(\sqrt{J^2 - (1 - \delta^2)J} - \delta J)^2}{mJ(J+1)^2}$$

with  $J = 3\kappa B^2(1 + 4\sqrt{n\kappa})$ , which optimizes the worst-case linear rate  $\rho = (1 - \alpha m/1.5)^{1/2B}$  from the DIGing

<sup>2</sup>While the authors of [5] do not explicitly assume the spectrum property from Assumption 2, they make use of this property in [5, Equation (14)] to prove their convergence bound.

paper. While our bound depends on the spectral gap  $\sigma$  and condition ratio  $\kappa$ , this bound also depends on the number of agents  $n$  and is vacuous in the limit as  $n \rightarrow \infty$ . In the plot on the right, we use the (much larger) stepsize which optimizes our bound from Theorem 1, for which the bound from the DIGing paper is vacuous.

To summarize, our analysis is tighter than previous bounds for DIGing (for both small and large stepsizes), is independent of the number of agents, and is applicable to any algorithm of the form (2).

#### References

- [1] J. A. Bazerque and G. B. Giannakis. Distributed spectrum sensing for cognitive radio networks by exploiting sparsity. *IEEE Transactions on Signal Processing*, 58(3):1847–1862, 2009.
- [2] P. A. Forero, A. Cano, and G. B. Giannakis. Consensus-based distributed support vector machines. *Journal of Machine Learning Research*, 11:1663–1707, 2010.
- [3] A. Jadbabaie, J. Lin, and A. S. Morse. Coordination of groups of mobile autonomous agents using nearest neighbor rules. *IEEE Transactions on Automatic Control*, 48(6):988–1001, 2003.
- [4] A. Nedić, A. Olshevsky, A. Ozdaglar, and J. N. Tsitsiklis. On distributed averaging algorithms and quantization effects. *IEEE Transactions on Automatic Control*, 54(11):2506–2517, 2009.
- [5] A. Nedić, A. Olshevsky, and W. Shi. Achieving geometric convergence for distributed optimization over time-varying graphs. *SIAM Journal on Optimization*, 27(4):2597–2633, 2017.
- [6] G. Qu and N. Li. Harnessing smoothness to accelerate distributed optimization. *IEEE Transactions on Control of Network Systems*, 5(3):1245–1260, 2018.

- [7] M. Rabbat and R. Nowak. Distributed optimization in sensor networks. In *Proceedings of the 3rd International Symposium on Information Processing in Sensor Networks*, pages 20–27. ACM, 2004.
- [8] A. Sundararajan, B. Van Scoy, and L. Lessard. Analysis and design of first-order distributed optimization algorithms over time-varying graphs. *IEEE Transactions on Control of Network Systems*, 2020. (to appear).
- [9] J. Xu, S. Zhu, Y. C. Soh, and L. Xie. Convergence of asynchronous distributed gradient methods over stochastic networks. *IEEE Transactions on Automatic Control*, 63(2):434–448, 2018.

## Appendix

**Proof of Proposition 1.** Suppose that the algorithm matrices satisfy the conditions in (3), and let  $y_i^* = y_{\text{opt}} \in \mathbb{R}^{1 \times d}$  for all  $i$  with  $\sum_{i=1}^n \nabla f_i(y_{\text{opt}}) = 0$ . Then there exist matrices  $\bar{x} \in \mathbb{R}^{s \times d}$ ,  $\hat{x} \in \mathbb{R}^s$ , and  $\bar{v} \in \mathbb{R}^{1 \times d}$  such that

$$\begin{bmatrix} 0 \\ y_{\text{opt}} \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} A - I & B_v \\ C_y & D_{yv} \\ C_z & D_{zv} - I \\ F_x & 0 \end{bmatrix} \begin{bmatrix} \bar{x} \\ \bar{v} \end{bmatrix}, \quad \begin{bmatrix} B_u \\ D_{yu} \\ D_{zu} \end{bmatrix} = \begin{bmatrix} A - I \\ C_y \\ C_z \end{bmatrix} \hat{x}.$$

For all agents  $i \in \{1, \dots, n\}$ , use these to define

$$\begin{aligned} x_i^* &= \bar{x} - \hat{x} \nabla f_i(y_{\text{opt}}), & y_i^* &= y_{\text{opt}}, & z_i^* &= \bar{v}, \\ u_i^* &= \nabla f_i(y_{\text{opt}}), & v_i^* &= \bar{v}, \end{aligned}$$

which is a fixed point of (2) corresponding to  $y_{\text{opt}}$ .

Now suppose  $(x_i^*, y_i^*, z_i^*, u_i^*, v_i^*)$  is a fixed point of (2) such that  $y_i^* = y_{\text{opt}}$  and  $\sum_{i=1}^n \nabla f_i(y_{\text{opt}}) = 0$ . Define the average state as  $\bar{x} = (1/n) \sum_{i=1}^n x_i^*$ , and similarly for the other points. Then  $\bar{u} = 0$ , so the concatenated matrix  $\text{vcat}(\bar{x}, \bar{v})$  must be nonzero and in the space (3a). Now let  $q$  be any nonzero vector such that  $q^\top \mathbf{1} = 0$ . For the fixed point to not depend on the sequence of gossip matrices,  $v_i = \bar{v}$  must be in consensus. Then multiplying (2a) by  $q_i$  and summing over  $i \in \{1, \dots, n\}$ ,

$$0 = \begin{bmatrix} A - I \\ C_y \\ C_z \end{bmatrix} \sum_{i=1}^n (q_i x_i^*) + \begin{bmatrix} B_u \\ D_{yu} \\ D_{zu} \end{bmatrix} \sum_{i=1}^n (q_i u_i^*).$$

This must hold for all objective functions, which implies the condition in (3b).  $\square$

**Proof of Theorem 1.** Let  $(x_i, y_i, z_i, u_i, v_i)$  denote a trajectory of algorithm (2). From Proposition 1, there exists a fixed point  $(x_i^*, y_i^*, z_i^*, u_i^*, v_i^*)$  with  $y_i^* = y_{\text{opt}}$  for all  $i$ , where  $y_{\text{opt}} \in \mathbb{R}^{1 \times d}$  is the optimizer of (1). We denote the error coordinates as  $\tilde{x}_i(k) := x_i(k) - x_i^*$ , and similarly for the other signals.

From the invariant (2c) and the gossip matrix being weight-balanced, there exists  $\tilde{s}(k)$  such that

$$\Psi \tilde{s}(k) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{\eta}_i(k),$$

where  $\eta_i(k)$  is the basis in (6) and  $\tilde{\eta}_i(k)$  the corresponding error signal. Multiplying  $X$  in the consensus LMI on the

right and left by  $\tilde{s}(k)$  and its transpose, respectively, and defining  $\Pi = \frac{1}{n} \mathbf{1}\mathbf{1}^\top$ , we obtain the consensus inequality

$$\begin{aligned} 0 &\geq \langle \tilde{\xi}(k+1), (\Pi \otimes P) \tilde{\xi}(k+1) \rangle \\ &\quad - \rho^2 \langle \tilde{\xi}(k), (\Pi \otimes P) \tilde{\xi}(k) \rangle \\ &\quad + \sum_{\ell=0}^{B-1} \lambda(\ell) \left\langle \begin{bmatrix} \tilde{y}(k) \\ \tilde{u}(k) \end{bmatrix}, (M_0 \otimes \Pi) \begin{bmatrix} \tilde{y}(k) \\ \tilde{u}(k) \end{bmatrix} \right\rangle, \end{aligned} \quad (10a)$$

where  $\langle A, B \rangle = \text{tr}(A^\top B)$  is the Frobenius inner product. Now choose vectors  $q_2, \dots, q_n \in \mathbb{R}^n$  so that the matrix  $[\mathbf{1}/\sqrt{n} \ q_2 \ \dots \ q_n]$  is orthonormal, and multiply the matrix  $Y$  in the disagreement LMI on the right and left by the weighted sum  $\sum_{i=1}^n (q_m)_i \tilde{\xi}_i(k)$  and its transpose, respectively, and sum over  $m \in \{2, \dots, n\}$ . This results in the disagreement inequality

$$\begin{aligned} 0 &\geq \langle \tilde{\xi}(k+1), ((I - \Pi) \otimes Q) \tilde{\xi}(k+1) \rangle \\ &\quad - \rho^2 \langle \tilde{\xi}(k), ((I - \Pi) \otimes Q) \tilde{\xi}(k) \rangle \\ &\quad + \sum_{\ell=0}^{B-1} \lambda(\ell) \left\langle \begin{bmatrix} \tilde{y}(k) \\ \tilde{u}(k) \end{bmatrix}, (M_0 \otimes (I - \Pi)) \begin{bmatrix} \tilde{y}(k) \\ \tilde{u}(k) \end{bmatrix} \right\rangle \\ &\quad + \left\langle \begin{bmatrix} \tilde{w}(k) \\ \tilde{w}(k+B) \end{bmatrix}, (M_1 \otimes (I - \Pi) \otimes R) \begin{bmatrix} \tilde{w}(k) \\ \tilde{w}(k+B) \end{bmatrix} \right\rangle \\ &\quad + \sum_{\ell=0}^{B-1} \left\langle \begin{bmatrix} \tilde{z}(\ell) \\ \tilde{w}(\ell) \\ \tilde{v}(\ell) \\ \tilde{w}(\ell+1) \end{bmatrix}, (M_2 \otimes (I - \Pi) \otimes S(\ell)) \begin{bmatrix} \tilde{z}(\ell) \\ \tilde{w}(\ell) \\ \tilde{v}(\ell) \\ \tilde{w}(\ell+1) \end{bmatrix} \right\rangle \end{aligned} \quad (10b)$$

where we used that  $\{q_m\}_{m=1}^n$  form an orthonormal basis for  $\mathbb{R}^n$ . Summing the inequalities in (10), we obtain

$$\begin{aligned} 0 &\geq V(\tilde{\xi}(k+1)) - \rho^2 V(\tilde{\xi}(k)) \\ &\quad + \sum_{\ell=0}^{B-1} \lambda(\ell) \left\langle \begin{bmatrix} \tilde{y}^k \\ \tilde{u}^k \end{bmatrix}, (M_0 \otimes I) \begin{bmatrix} \tilde{y}^k \\ \tilde{u}^k \end{bmatrix} \right\rangle \\ &\quad + \left\langle \begin{bmatrix} \tilde{w}(k) \\ \tilde{w}(k+B) \end{bmatrix}, (M_1 \otimes (I - \Pi) \otimes R) \begin{bmatrix} \tilde{w}(k) \\ \tilde{w}(k+B) \end{bmatrix} \right\rangle \\ &\quad + \sum_{\ell=0}^{B-1} \left\langle \begin{bmatrix} \tilde{z}(\ell) \\ \tilde{w}(\ell) \\ \tilde{v}(\ell) \\ \tilde{w}(\ell+1) \end{bmatrix}, (M_2 \otimes (I - \Pi) \otimes S(\ell)) \begin{bmatrix} \tilde{z}(\ell) \\ \tilde{w}(\ell) \\ \tilde{v}(\ell) \\ \tilde{w}(\ell+1) \end{bmatrix} \right\rangle \end{aligned}$$

where the Lyapunov function  $V$  is defined in (5). Each quadratic form in the first summation is nonnegative from Assumption 1, the following term is nonnegative from the joint spectrum property in Assumption 2, and each term in the last summation is nonnegative from the spectrum property in Assumption 2; see [8, Prop. 15–16]. Then using the slight abuse of notation  $V(k) = V(\tilde{\xi}(k))$ , we have the decrease condition  $V(k+1) \leq \rho^2 V(k)$ , which implies  $V(k) \leq \rho^{2k} V(0)$ . Now define the matrix  $T := \Pi \otimes P + (I - \Pi) \otimes Q$ , and note that  $T \succ 0$  since  $P$  and  $Q$  are positive definite. Letting  $\text{cond}(T)$  denote the condition number of  $T$ , we have the bound

$$\|x_i(k) - x_i^*\|^2 \leq \text{cond}(T) V(k) \leq \rho^{2k} \text{cond}(T) V(0),$$

so (9) holds with  $\gamma = \sqrt{\text{cond}(T) V(0)}$ .  $\square$