

# Newton-Raphson Consensus for Distributed Convex Optimization

Damiano Varagnolo, *Member, IEEE*, Filippo Zanella, *Member, IEEE*, Angelo Cenedese, *Member, IEEE*,  
Gianluigi Pillonetto, *Member, IEEE*, and Luca Schenato, *Senior Member, IEEE*

**Abstract**—We address the problem of distributed unconstrained convex optimization under separability assumptions, i.e., the framework where each agent of a network is endowed with a local private multidimensional convex cost, is subject to communication constraints, and wants to collaborate to compute the minimizer of the sum of the local costs. We propose a design methodology that combines average consensus algorithms and separation of time-scales ideas. This strategy is proved, under suitable hypotheses, to be globally convergent to the true minimizer. Intuitively, the procedure lets the agents distributedly compute and sequentially update an approximated Newton-Raphson direction by means of suitable average consensus ratios. We show with numerical simulations that the speed of convergence of this strategy is comparable with alternative optimization strategies such as the Alternating Direction Method of Multipliers. Finally, we propose some alternative strategies which trade-off communication and computational requirements with convergence speed.

**Index Terms**—Consensus, distributed optimization, multi-agent systems, Newton-Raphson methods, smooth functions, unconstrained convex optimization.

## I. INTRODUCTION

OPTIMIZATION is a pervasive concept underlying many aspects of modern life [1]–[5], and it also includes the management of distributed systems, i.e., artifacts composed by a multitude of interacting entities often referred to as “agents”. Examples are transportation systems, where the agents are both the vehicles and the traffic management devices (traffic lights), and smart electrical grids, where the agents are the energy producers-consumers and the power transformers-transporters. Here we consider the problem of distributed optimization, i.e., the class of algorithms suitable for networked systems

and characterized by the absence of a centralized coordination unit [6]–[8]. Distributed optimization tools have received an increasing attention over the last years, concurrently with the research on networked control systems. Motivations comprise the fact that the former methods let the networks self-organize and adapt to surrounding and changing environments, and that they are necessary to manage extremely complex systems in an autonomous way with only limited human intervention. In particular we focus on unconstrained convex optimization, although there is a rich literature also on distributed constrained optimization such as Linear Programming [9].

**Literature Review:** The literature on distributed unconstrained convex optimization is extremely vast and a first taxonomy can be based whether the strategy uses or not the Lagrangian framework, see, e.g., [5, Ch. 5].

Among the distributed methods exploiting Lagrangian formalism, the most widely known algorithm is Alternating Direction Method of Multipliers (ADMM) [10], whose roots can be traced back to [11]. Its efficacy in several practical scenarios is undoubted, see, e.g., [12] and references therein. A notable size of the dedicated literature focuses on the analysis of its convergence performance and on the tuning of its parameters for optimal convergence speed, see, e.g., [13] for Least Squares (LS) estimation scenarios, [14] for linearly constrained convex programs, and [15] for more general ADMM algorithms. Even if proved to be an effective algorithm, ADMM suffers from requiring synchronous communication protocols, although some recent attempts for asynchronous and distributed implementations have appeared [16]–[18].

On the other hand, among the distributed methods not exploiting Lagrangian formalisms, the most popular ones are the Distributed Subgradient Methods (DSMs) [19]. Here the optimization of non-smooth cost functions is performed by means of subgradient based descent/ascent directions. These methods arise in both primal and dual formulations, since sometimes it is better to perform dual optimization. Subgradient methods have been exploited for several practical purposes, e.g., to optimally allocate resources in Wireless Sensor Networks (WSNs) [20], to maximize the convergence speeds of gossip algorithms [21], to manage optimality criteria defined in terms of ergodic limits [22]. Several works focus on the analysis of the convergence properties of the DSM basic algorithm [23]–[25] (see [26] for a unified view of many convergence results). We can also find analyses for several extensions of the original idea, e.g., directions that are computed combining information from other agents [27], [28] and stochastic errors in the evaluation of the subgradients [29]. Explicit characterizations can also show trade-offs between desired accuracy and number of iterations [30].

Manuscript received April 10, 2014; revised November 4, 2014, May 31, 2015, and June 9, 2015; accepted June 18, 2015. This work was supported by the Framework Programme for Research and Innovation Horizon 2020 under the grant agreement 636834 “DISIRE”, the Swedish research council Norrbottens Forskningsråd, by the University of Padova under the Progetto di Ateneo CPDA147754/14-New statistical learning approach for multi-agents adaptive estimation and coverage control, by the Italian Ministry of Education under Grant SCN 00398, and by Smart & safe Energy-aware Assisted Living. Recommended by Associate Editor A. Papachristodoulou.

D. Varagnolo is with the Department of Computer Science, Electrical and Space Engineering, Luleå University of Tec, 971 87 Luleå, Sweden (e-mail: damiano.varagnolo@ltu.se).

F. Zanella, A. Cenedese, G. Pillonetto, and L. Schenato are with the Department of Information Engineering, Università di Padova, 35131 Padova, Italy (e-mail: fzanella@dei.unipd.it; angelo.cenedese@dei.unipd.it; giapi@dei.unipd.it; schenato@dei.unipd.it).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TAC.2015.2449811

These methods have the advantage of being easily distributed, to have limited computational requirements and to be inherently asynchronous as shown in [31]–[33]. However they suffer from low convergence rate since they require the update steps to decrease to zero as  $1/t$  (being  $t$  the time) therefore as a consequence the rate of convergence is sub-exponential. In fact, one of the current trends is to design strategies that improve the convergence rate of DSMs. For example, a way is to accelerate the convergence of subgradient methods by means of multi-step approaches, exploiting the history of the past iterations to compute the future ones [34]. Another is to use Newton-like methods, when additional smoothness assumptions can be used. These techniques are based on estimating the Newton direction starting from the Laplacian of the communication graph. More specifically, distributed Newton techniques have been proposed in dual ascent scenarios [35]–[37]. Since the Laplacian cannot be computed exactly, the convergence rates of these schemes rely on the analysis of inexact Newton methods [38]. These Newton methods are shown to have super-linear convergence under specific assumptions, but can be applied only to specific optimization problems such as network flow problems.

Recently, several alternative approaches to ADMM and DSM have appeared. For example, in [39], [40] the authors construct contraction mappings by means of cyclic projections of the estimate of the optimum onto the constraints. A similar idea based on contraction maps is used in F-Lipschitz methods [41] but it requires additional assumptions on the cost functions. Other methods are the control-based approach [42] which exploits distributed consensus, the distributed randomized Kaczmarz method [43] for quadratic cost functions, and distributed dual sub-gradient methods [44].

*Statement of Contributions:* Here we propose a distributed Newton-Raphson optimization procedure, named Newton-Raphson Consensus (NRC), for the exact minimization of smooth multidimensional convex separable problems, where the global function is a sum of private local costs. With respect to the classification proposed before, the strategy exploits neither Lagrangian formalisms nor Laplacian estimation steps. More specifically, it is based on average consensus techniques [45] and on the principle of separation of time scales [46, Ch. 11]. The main idea is that agents compute and keep updated, by means of average consensus protocols, an approximated Newton-Raphson direction that is built from suitable Taylor expansions of the local costs. Simultaneously, agents move their local guesses towards the Newton-Raphson direction. It is proved that, if the costs satisfy some smoothness assumptions and the rate of change of the local update steps is sufficiently slow to allow the consensus algorithm to converge, then the NRC algorithm exponentially converges to the global minimizer.

The main contribution of this work is to propose an algorithm that extends Newton-Raphson ideas in a distributed setting, thus being able to exploit second order information to speed up converge rate. By using singular perturbation theory we formally show that under suitable assumptions the convergence of the algorithm is exponential (linear in logspace). Differently, DSM algorithms have sublinear convergence rate even if the cost functions are smooth [39], [47], although they are easy to implement and can be employed also for non-smooth cost func-

tions and for constrained optimization. We also show by means of numerical simulations on real-world database benchmarks that the proposed algorithm exhibits faster convergence rates (in number of communications) than standard implementations of distributed ADMM algorithms [12], probably due to the second-order information embedded into the Newton-Raphson consensus. Although we have no theoretical guarantee of the superiority of the proposed algorithmic in terms of convergence rate, these simulations suggest that it is at least a potentially competitive algorithm. Moreover, one of the promising features of the NRC is that it is essentially based on average consensus algorithms, for which there exist robust implementations that encompass asynchronous communications, time-varying network topologies [48], directed graphs [49], and packet-losses effects.

*Structure of the Paper:* The paper is organized as follows: Section II collects the notation used through the whole paper, while Section III formulates the considered problem and provides some ancillary results that are then used to study the convergence properties of the main algorithm. Section IV proposes the main optimization algorithm, provides convergence results and describes some strategies to trade-off communication and computational complexities with convergence speed. Section V compares, via numerical simulations, the performance of the proposed algorithm with several distributed optimization strategies available in the literature. Finally, Section VI collects some final observations and suggests future research directions. We collect all the proofs in the Appendix.

## II. NOTATION

We model the communication network as a graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$  whose vertices  $\mathcal{N} := \{1, 2, \dots, N\}$  represent the agents and whose edges  $(i, j) \in \mathcal{E}$  represent the available communication links. We assume that the graph is undirected and connected, and that the matrix  $P \in \mathbb{R}^{N \times N}$  is stochastic, i.e., its elements are non-negative, it is s.t.  $P\mathbb{1} = \mathbb{1}$  (where  $\mathbb{1} := [1 \ 1 \ \dots \ 1]^T \in \mathbb{R}^N$ ), symmetric, i.e.,  $P = P^T$  and consistent with the graph  $\mathcal{G}$ , in the sense that each entry  $p_{ij}$  of  $P$  is  $p_{ij} > 0$  only if  $(i, j) \in \mathcal{E}$ . We recall that if  $P$  is stochastic, symmetric, and includes all edges (i.e.,  $p_{ij} > 0$  if and only if  $(i, j) \in \mathcal{E}$ ) then  $\lim_{k \rightarrow \infty} P^k = (1/N)\mathbb{1}\mathbb{1}^T$ . Such  $P$ 's are also often referred to as *average consensus matrices*. We will indicate with  $\rho(P) := \max_{i, \lambda_i \neq 1} |\lambda_i(P)|$  the spectral radius of  $P$ , with  $\sigma(P) := 1 - \rho(P)$  its spectral gap.

We use fraction bars to indicate also Hadamard divisions, e.g., if  $\mathbf{a} = [a_1, \dots, a_N]^T$  and  $\mathbf{b} = [b_1, \dots, b_N]^T$  then  $\mathbf{a}/\mathbf{b} := [(a_1/b_1) \ \dots \ (a_N/b_N)]^T$ . Fraction bars like the previous ones may also indicate pre-multiplication with inverse matrices, i.e., if  $b_i$  is a matrix then  $a_i/b_i$  indicates  $b_i^{-1}a_i$ . We indicate with  $n$  the dimensionality of the domains of the cost functions,  $k$  a discrete time index,  $t$  a continuous time index. For notational simplicity we denote differentiation with  $\nabla$  operators, so that  $\nabla f = \partial f / \partial x$  and  $\nabla^2 f = \partial^2 f / \partial x^2$ . With a little abuse of notation, we will define  $\chi = (x, Z)$ , where  $x \in \mathbb{R}^n$  and  $Z \in \mathbb{R}^{\ell \times q}$  as the vector obtained by stacking in a column both the vector  $x$  and the vectorized matrix  $Z$ . We indicate with  $\|\cdot\|$  Frobenius norms. With an other abuse of notation we also define the norm of the pair  $\chi = (x, Z)$  where  $x$  is a vector and  $Z$  a matrix with  $\|\chi\|^2 = \|x\|^2 + \|Z\|^2$ .

When using plain italic fonts with a subscript (usually  $i$ , e.g.,  $x_i \in \mathbb{R}^n$ ) we refer to the local decision variable of the specific agent  $i$ . When using bold italic fonts, e.g.,  $\mathbf{x}$ , we instead refer to the collection of the decision variables of all the various agents, e.g.,  $\mathbf{x} := [x_1^T, \dots, x_N^T]^T \in \mathbb{R}^{nN}$ . To indicate special variables we will instead consider the following notation:

$$\begin{aligned}\bar{\mathbf{x}} &:= \frac{1}{N} \sum_{i=1}^N x_i \quad \mathbb{R}^n \\ \mathbf{x}^{\parallel} &:= \mathbb{1}_N \otimes \bar{\mathbf{x}} \quad \mathbb{R}^{nN} \\ \mathbf{x}^{\perp} &:= \mathbf{x} - \mathbf{x}^{\parallel} \quad \mathbb{R}^{nN}\end{aligned}$$

As in [46, p. 116], we say that a function  $V$  is a *Lyapunov function* for a specific dynamics if  $V$  is continuously differentiable and satisfies  $V(0)=0$ ,  $V(x) > 0$  for  $x \neq 0$ , and  $\dot{V}(x) \leq 0$ .

### III. PROBLEM FORMULATION AND PRELIMINARY RESULTS

#### A. Structure of the Section

Our main contribution is to characterize the convergence properties of the distributed Newton-Raphson (NR) scheme proposed in Section IV. In doing so we both exploit standard singular perturbation analysis tools [46, Ch. 11] [50] and a set of ancillary results, collected for readability in this section.

The logical flow of these ancillary results is the following: Section III-C claims that, under suitable assumptions, forward-Euler discretizations of stable continuous dynamics lead to stable discrete dynamics. This basic result enables reasoning on continuous-time systems. Then, Section III-D and E respectively claim that single- and multi-agent continuous-time NR dynamics satisfy these discretization assumptions. Section III-F and G then generalize these dynamics by introducing perturbation terms that mimic the behavior of the proposed main optimization algorithm, and characterize their stability properties. Summarizing, the ancillary results characterize the stability properties of systems that are progressive approximations of the dynamics under investigation.

#### B. Problem Formulation

We assume that the  $N$  agents of the network are endowed with cost functions  $f_i : \mathbb{R}^n \mapsto \mathbb{R}$  so that

$$\bar{f} : \mathbb{R}^n \mapsto \mathbb{R}, \quad \bar{f}(x) := \frac{1}{N} \sum_{i=1}^N f_i(x) \quad (1)$$

is a well-defined global cost. We assume that the aim of the agents is to cooperate and distributedly compute the minimizer of  $\bar{f}$ , namely

$$x^* := \arg \min_{x \in \mathbb{R}^n} \bar{f}(x). \quad (2)$$

We now enforce the following simplifying assumptions, valid throughout the rest of the paper.

**Assumption 1 (Convexity):** The local costs  $f_i$  in (1) are of class  $\mathcal{C}^3$ . Moreover the global cost  $\bar{f}$  has bounded positive definite Hessian, i.e.,  $0 < cI \leq \nabla^2 \bar{f}(x) \leq mI$  for some  $c, m \in \mathbb{R}_+$  and  $\forall x \in \mathbb{R}^n$ . Moreover, w.l.o.g., we assume  $\bar{f}(x^*) = 0$ ,  $c \leq 1$  and  $m \geq 1$ .

The scalar  $c$  is assumed to be known by all the agents *a-priori*. Assumption 1 ensures that  $x^*$  in (2) exists and is unique. The strictly positive definite Hessian is moreover a mild sufficient

condition to guarantee that the minimum  $x^*$  defined in (2) will be globally exponentially stable under the continuous and discrete Newton-Raphson dynamics described in the following Theorem 3. We also notice that, for the subsequent Theorems 2 and 3, in principle just the average function  $\bar{f}$  needs to have specific properties, and thus no conditions for the single  $f_i$ 's are required (that for example might be even non convex). For the convergence of the distributed NR scheme we will nonetheless enforce the more restrictive Assumptions 5 and 9, not presented now for readability issues. In the rest of this section, in order to simplify notation, we will consider, without loss of generality, the following translated cost functions:

$$f'_i(x) = f_i(x + x^*), \quad \bar{f}'(x) = \frac{1}{N} \sum_{i=1}^N f'_i(x) \quad (3)$$

so that the origin becomes the minimizer of the averaged cost function  $\bar{f}'(x)$ , i.e.,  $\bar{f}'(0) = 0$ .

#### C. Stability of Discretized Dynamics

This subsection aims to show that, under suitable assumptions, forward-Euler discretization of suitable exponentially stable continuous-time dynamics maintains the same global exponential stability properties.

**Theorem 2:** Let the continuous-time system

$$\dot{x} = \phi(x) \quad (4)$$

admit  $x = 0 \in \mathbb{R}^n$  as an equilibrium, and let  $V(x) : \mathbb{R}^n \mapsto \mathbb{R}$  be a Lyapunov function for (4) for which there exist positive scalars  $a_1, a_2, a_3, a_4$  s.t.,  $\forall x \in \mathbb{R}^n$

$$a_1 I \leq \nabla^2 V(x) \leq a_2 I \quad (5a)$$

$$\frac{\partial V(x)}{\partial x} \phi(x) \leq -a_3 \|x\|^2 \quad (5b)$$

$$\|\phi(x)\| \leq a_4 \|x\|. \quad (5c)$$

Then

- a) for system (4) the origin is globally exponentially stable;
- b) for the following forward-Euler discretization of system (4):

$$x(k+1) = x(k) + \varepsilon \phi(x(k)) \quad (6)$$

there exists a positive scalar  $\bar{\varepsilon}$  such that for every  $\varepsilon \in (0, \bar{\varepsilon})$  the origin is globally exponentially stable.

#### D. Stability of Single-Agent NR Dynamics

This subsection shows that the results of Section III-C apply to continuous NR dynamics, i.e., that forward-Euler discretizations maintain global exponential stability properties.<sup>1</sup>

**Theorem 3:** Let

$$\phi_{\text{NR}}(x) := -\bar{h}'(x)^{-1} \nabla \bar{f}'(x) \quad (7)$$

be defined by a generic function  $\bar{h}'(x) \in \mathbb{R}^{n \times n}$  that satisfies the positive definiteness conditions  $cI \leq \bar{h}'(x) = \bar{h}'(x)^T \leq mI$  for all  $x \in \mathbb{R}^n$  where  $c$  and  $m$  are defined in Assumption 1. Let (7) define both the dynamics

$$\dot{x} = \phi_{\text{NR}}(x) \quad (8)$$

$$x(k+1) = x(k) + \varepsilon \phi_{\text{NR}}(x(k)). \quad (9)$$

<sup>1</sup>We notice that other asymptotic properties of continuous time NR methods are available in the literature, e.g., [51], [52].



286 Then, under Assumption 1:

287 a)

$$V_{\text{NR}}(x) := \bar{f}'(x) \quad (10)$$

288 is a Lyapunov function for (8);

289 b) there exist positive scalars  $b_1, b_2, b_3, b_4$  s.t.,  $\forall x \in \mathbb{R}^n$

$$\begin{aligned} b_1 I &\leq \nabla^2 V_{\text{NR}}(x) \leq b_2 I & (11a) \\ \frac{\partial V_{\text{NR}}}{\partial x} \phi_{\text{NR}}(x) &\leq -b_3 \|x\|^2 & (11b) \\ \|\phi_{\text{NR}}(x)\| &\leq b_4 \|x\|, & (11c) \end{aligned}$$

290 i.e., Theorem 2 applies to dynamics (8) and (9).

291 For suitable choices of  $\bar{h}'(x)$  the dynamics (8) corresponds to  
292 continuous versions of well known descent dynamics. Indeed,  
293 the correspondences are

$$\bar{h}'(x) = \begin{cases} \nabla^2 \bar{f}'(x) & \rightarrow \text{Newton-Raphson descent} & (12a) \\ \text{diag}[\nabla^2 \bar{f}'(x)] & \rightarrow \text{Jacobi descent} & (12b) \\ I & \rightarrow \text{Gradient descent} & (12c) \end{cases}$$

294 where  $\text{diag}[A]$  is a diagonal matrix containing the main diago-  
295 nal of  $A$ . Note that for every choice of  $\bar{h}'(x)$  as in (12a)–(12c),  
296 Assumption 1 ensures the hypotheses<sup>2</sup> of Theorem 3, therefore  
297 by combining Theorem 3 with Theorem 2 we are guaranteed  
298 that both continuous and discrete generalized NR dynamics  
299 induced by (7) are globally exponentially stable.

300 *Lemma 4:* Under Assumption 1, the origin is a globally  
301 exponentially stable point for dynamics (8). Moreover there  
302 exists  $\bar{\varepsilon} > 0$  such that the origin is a globally exponentially  
303 stable point also for dynamics (9) for all  $\varepsilon < \bar{\varepsilon}$ .

304 The previous lemma and theorems do not require  $\bar{h}'(x)$   
305 to be differentiable. However, differentiability may be used  
306 to linearize the system dynamics and obtain explicit rates of  
307 convergence. In fact, the linearized dynamics around the origin  
308 is given by

$$F(0) := \frac{\partial \phi_{\text{NR}}(0)}{\partial x} = -\bar{h}'(0)^{-1} \nabla^2 \bar{f}'(0) - \frac{\partial \bar{h}'(0)^{-1}}{\partial x} \nabla \bar{f}'(0).$$

309 In particular, for the NR descent it holds that  $\bar{h}'(x) = \nabla^2 \bar{f}'(x)$ .  
310 Thus in this case  $F(0) = -I$ , since  $\nabla \bar{f}'(0) = 0$ , and this says  
311 that the linearized continuous time NR dynamics is  $\dot{x} = -x$ ,  
312 independent of the cost  $\bar{f}'(x)$  and whose rate of convergence is  
313 unitary and uniform along any direction.

#### 314 E. Stability of Multi-Agent NR Dynamics

315 We now generalize (8) by considering  $N$  coupled dynamical  
316 systems that, when starting at the very same initial condition,  
317 behave like  $N$  decoupled systems (8). This novel dynamics  
318 is the core of the slow-dynamics embedded in the main algo-  
319 rithm presented in Section IV. In this section we also include  
320 additional assumptions to show that the generalization of (8)  
321 presented here preserves global exponential stability and some  
322 other additional properties.

<sup>2</sup>For the Jacobi descent, clearly  $\min_{\|x\|=1} x^T \text{diag}[\nabla^2 \bar{f}'(x)] x = \min_{x \in \{e_1, \dots, e_n\}} x^T \text{diag}[\nabla^2 \bar{f}'(x)] x = \min_{x \in \{e_1, \dots, e_n\}} x^T \nabla^2 \bar{f}'(x) x \geq \min_{\|x\|=1} x^T \nabla^2 \bar{f}'(x) x = c$ , where  $e_i$  is the  $n$ -dimensional vector with all zeros except for a one in the  $i$ -th entry.

To this aim we introduce some additional notation: let  $h'_i(x) : \mathbb{R}^n \mapsto \mathbb{R}^{n \times n}$ ,  $i = 1, \dots, N$  be defined according to one of the possible three cases

$$\begin{aligned} h'_i(x) &= \begin{cases} \nabla^2 f'_i(x) & (13a) \\ \text{diag}[\nabla^2 f'_i(x)] & (13b) \\ I & (13c) \end{cases} \end{aligned}$$

so that  $h'_i(x) = h'_i(x)^T$  for all  $x$ . Moreover let

$$\begin{aligned} h'(\mathbf{x}) &:= [h'_1(x_1), \dots, h'_N(x_N)]^T \quad \mathbb{R}^{nN} \mapsto \mathbb{R}^{nN \times nN} \\ \bar{h}'(\mathbf{x}) &:= \frac{1}{N} \sum_{i=1}^N h'_i(x_i) \quad \mathbb{R}^{nN} \mapsto \mathbb{R}^{n \times n} \\ \bar{\bar{h}}'(\bar{x}) &:= \frac{1}{N} \sum_{i=1}^N h'_i(\bar{x}) \quad \mathbb{R}^n \mapsto \mathbb{R}^{n \times n} \end{aligned}$$

be additional composite functions defined starting from the  $h'_i$ 's (recall that  $\mathbf{x} := [x_1^T, \dots, x_N^T]^T \in \mathbb{R}^{nN}$  and that  $\bar{x} := (1/N) \sum_{i=1}^N x_i \in \mathbb{R}^n$ ). Let moreover

$$g'_i(x) := h'_i(x)x - \nabla f'_i(x) \quad \mathbb{R}^n \mapsto \mathbb{R}^n \quad (14)$$

and  $g'(\mathbf{x})$ ,  $\bar{g}'(\mathbf{x})$ ,  $\bar{\bar{g}}'(\bar{x})$  be defined accordingly as for  $h'_i$ .

The definitions of  $h'_i$  and  $g'_i$  are instrumental to generalize the NR dynamics (8) to the distributed case. Indeed, let

$$\psi(\mathbf{x}) := \bar{h}'(\mathbf{x})^{-1} \bar{g}'(\mathbf{x}) \quad \mathbb{R}^{nN} \mapsto \mathbb{R}^n \quad (15)$$

(with the existence of  $\bar{h}'(\mathbf{x})^{-1}$  guaranteed by the following Assumption 5). It is easy to verify that the previous functions satisfy the following properties:

$$\begin{aligned} \bar{h}'(\mathbf{x}^{\parallel}) &= \bar{h}'(\bar{x}) & (16a) \\ \bar{g}'(\mathbf{x}^{\parallel}) &= \bar{g}'(\bar{x}) = \bar{h}'(\bar{x})\bar{x} - \nabla \bar{f}'(\bar{x}) & (16b) \\ \psi(\mathbf{x}^{\parallel}) &= \bar{x} - \bar{h}'(\bar{x})^{-1} \nabla \bar{f}'(\bar{x}). & (16c) \end{aligned}$$

Consider then

$$\dot{\mathbf{x}} = \phi_{\text{PNR}}(\mathbf{x}) := -\mathbf{x} + \mathbb{1}_N \otimes \psi(\mathbf{x}) \quad (17)$$

that can be also equivalently written as

$$\dot{x}_i = -x_i + \psi(\mathbf{x}), \quad i = 1, \dots, N$$

i.e., as the combination of  $N$  independent dynamical systems that are driven by the same forcing term  $\psi(\mathbf{x})$ .

As mentioned above, this dynamics embeds the centralized generalized NR dynamics since, under identical initial conditions  $x_i(0) = \bar{x}(0) \in \mathbb{R}^n$  for all  $i$ , the trajectories coincide, i.e.,  $x_i(t) = \bar{x}(t)$ ,  $\forall i, \forall t \geq 0$ . Moreover, due to (16c)

$$\begin{aligned} \dot{\bar{x}} &= -\bar{x} + \psi(\mathbb{1}_N \otimes \bar{x}) \\ &= -\bar{x} + \bar{x} - \bar{h}'(\bar{x})^{-1} \nabla \bar{f}'(\bar{x}) = \phi_{\text{NR}}(\bar{x}) \end{aligned} \quad (18)$$

i.e., we obtain dynamics (7), that is, thanks to Theorem 3 and the assumption that  $\bar{h}'(\mathbf{x})$  is invertible, globally exponentially stable.

The question is then whether dynamics (17) is exponentially stable also in the general case where the  $x_i(0)$ 's may not be identical. To characterize this case we assume some additional global properties.

351 *Assumption 5 (Global Properties):* The local costs  $f'_1, \dots,$   
 352  $f'_N$  in (1) are s.t. there exist positive scalars  $m_g, a_g, a_h, a_\psi$  s.t.,  
 353  $\forall x, x' \in \mathbb{R}^n$  and  $\forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^{nN}$

$$\begin{aligned} & \begin{cases} cI \leq \bar{h}'(\mathbf{x}) \leq mI \\ \|\bar{g}'(\mathbf{x})\| \leq m_g \\ \|g'_i(x) - g'_i(x')\| \leq a_g \|x - x'\| \\ \|h'_i(x) - h'_i(x')\| \leq a_h \|x - x'\| \\ \|\psi(\mathbf{x}) - \psi(\mathbf{x}')\| \leq a_\psi \|\mathbf{x} - \mathbf{x}'\| \end{cases} \end{aligned} \quad \begin{aligned} (19a) \\ (19b) \\ (19c) \\ (19d) \\ (19e) \end{aligned}$$

354 with  $c$  and  $m$  from Assumption 1.

355 Note that Assumption 5 implies

$$\begin{aligned} & \begin{cases} \|\bar{g}'(\mathbf{x}) - \bar{g}'(\mathbf{x}')\| \leq a_g \|\mathbf{x} - \mathbf{x}'\| \\ \|\bar{h}'(\mathbf{x}) - \bar{h}'(\mathbf{x}')\| \leq a_h \|\mathbf{x} - \mathbf{x}'\| \\ \|g'(\mathbf{x}) - g'(\mathbf{x}')\| \leq a_g \|\mathbf{x} - \mathbf{x}'\| \\ \|h'(\mathbf{x}) - h'(\mathbf{x}')\| \leq a_h \|\mathbf{x} - \mathbf{x}'\|. \end{cases} \end{aligned} \quad \begin{aligned} (20a) \\ (20b) \\ (20c) \\ (20d) \end{aligned}$$

356 Using the previous assumptions we can now prove global  
 357 stability of dynamics (17).

358 *Theorem 6:* Under Assumptions 1 and 5, and for a suitable  
 359 positive scalar  $\eta$ ,

360 a)

$$V_{\text{PNR}}(\mathbf{x}) := V_{\text{NR}}(\bar{\mathbf{x}}) + \frac{1}{2}\eta\|\mathbf{x}^\perp\|^2 = \bar{f}'(\bar{\mathbf{x}}) + \frac{1}{2}\eta\|\mathbf{x}^\perp\|^2 \quad (21)$$

361 is a Lyapunov function for (17);

362 b) there exist positive scalars  $b_5, b_6, b_7, b_8$  s.t.,  $\forall \mathbf{x} \in \mathbb{R}^{nN}$

$$\begin{aligned} & \begin{cases} b_5 I \leq \nabla^2 V_{\text{PNR}}(\mathbf{x}) \leq b_6 I \\ \frac{\partial V_{\text{PNR}}}{\partial \mathbf{x}} \phi_{\text{PNR}}(\mathbf{x}) \leq -b_7 \|\mathbf{x}\|^2 \\ \|\phi_{\text{PNR}}(\mathbf{x})\| \leq b_8 \|\mathbf{x}\|. \end{cases} \end{aligned} \quad \begin{aligned} (22a) \\ (22b) \\ (22c) \end{aligned}$$

363 As in Lemma 4, combining Theorem 6 with Theorem 2 it is  
 364 possible to claim that (17) and its discrete-time counterpart are  
 365 globally exponentially stable.

### 366 F. Multi-Agent NR Dynamics Under Vanishing Perturbations

367 We now aim to generalize the dynamics  $\phi_{\text{PNR}}(\mathbf{x})$  by consid-  
 368 ering some perturbation term, that will be described by the vari-  
 369 able  $\chi$ . Let then  $\chi^y := (\chi_1^y, \dots, \chi_N^y)$  where  $\chi_i^y \in \mathbb{R}^n$ ,  $\chi^z :=$   
 370  $(\chi_1^z, \dots, \chi_N^z)$  where  $\chi_i^z \in \mathbb{R}^{n \times n}$ , and  $\chi := (\chi^y, \chi^z)$ .  
 371 We also define the operator  $[\cdot]_c : \mathbb{R}^{nN \times n} \mapsto \mathbb{R}^{nN \times n}$ , which  
 372 indicates the component-wise matrix-operation

$$[\mathbf{z}]_c = \begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix}_c := \begin{bmatrix} z'_1 \\ \vdots \\ z'_N \end{bmatrix} \quad z'_i = \begin{cases} z_i & \text{if } z_i \geq \frac{c}{2}I \\ \frac{c}{2}I & \text{otherwise.} \end{cases} \quad (23)$$

373 Consider then the perturbed version of the multi-agent NR  
 374 dynamics (17)

$$\dot{\mathbf{x}} = \phi_x(\mathbf{x}, \chi) := -\mathbf{x} - \mathbb{1}_N \otimes x^* + \frac{\chi^y + \mathbb{1}_N \otimes (\bar{g}'(\mathbf{x}) + \bar{h}'(\mathbf{x})x^*)}{[\chi^z + \mathbb{1}_N \otimes \bar{h}'(\mathbf{x})]_c} \quad (24)$$

375 where the division is a Hadamard division, as recalled in  
 376 Section II. Direct inspection of dynamics (24) then shows that

$$\phi_x(\mathbf{x}, \mathbf{0}) = \phi_{\text{PNR}}(\mathbf{x}). \quad (25)$$

377 The next lemma provides perturbations interconnection bounds  
 378 that will be used in Theorem 12.

*Lemma 7:* Under Assumptions 1 and 5 there exist positive  
 scalars  $a_x, a_\Delta$  s.t., for all  $\mathbf{x}$  and  $\chi$

$$\begin{aligned} & \begin{cases} \|\phi_x(\mathbf{x}, \chi)\| \leq a_x (\|\mathbf{x}\| + \|\chi\|) \\ \|\phi_x(\mathbf{x}, \chi) - \phi_{\text{PNR}}(\mathbf{x})\| \leq a_\Delta \|\chi\|. \end{cases} \end{aligned} \quad \begin{aligned} (26a) \\ (26b) \end{aligned}$$

### G. Multi-Agent NR Dynamics Under Non-Vanishing Perturbations

Let us now consider some additional properties of the flow  
 (24) for some specific non-vanishing perturbation. Consider  
 then the perturbations  $\xi^y \in \mathbb{R}^n$  and  $\xi^z \in \mathbb{R}^{n \times n}$ , and their multi-  
 agents versions  $\xi^y = \mathbb{1}_N \otimes \xi^y$ ,  $\xi^z = \mathbb{1}_N \otimes \xi^z$ . Consider also  
 the shorthand  $\xi = (\xi^y, \xi^z)$ . The equilibrium points of the dy-  
 namics induced by  $\phi_x(\mathbf{x}, \xi)$  are characterized by the following  
 theorem.

*Theorem 8:* Let  $\xi^y \in \mathbb{R}^n$ ,  $\xi^z \in \mathbb{R}^{n \times n}$ ,  $\xi = (\xi^y, \xi^z)$ ,  $\xi^y = \mathbb{1}_N \otimes$   
 $\xi^y$ ,  $\xi^z = \mathbb{1}_N \otimes \xi^z$ ,  $\xi = (\xi^y, \xi^z)$ , and consider the equation

$$\phi_x(\mathbf{x}, \xi) = 0$$

defining the equilibrium points of the dynamics  $\dot{\mathbf{x}} = \phi_x(\mathbf{x}, \xi)$ .  
 Then, under Assumptions 1 and 5 there exist a positive scalar  
 $r > 0$  and a unique continuously differentiable function  $\mathbf{x}^{\text{eq}} : \mathcal{B}_r$   
 $\mathcal{B}_r \rightarrow \mathbb{R}^{nN}$  where  $\mathcal{B}_r := \{\xi \mid \|\xi\| \leq r\}$  such that

$$\phi_x(\mathbf{x}^{\text{eq}}(\xi), \xi) = 0, \quad \mathbf{x}^{\text{eq}}(0) = 0. \quad (27)$$

Moreover,  $\mathbf{x}^{\text{eq}}(\xi) = \mathbb{1}_N \otimes x^{\text{eq}}(\xi)$ , with

$$x^{\text{eq}}(\xi) = \left( \bar{h}'(x^{\text{eq}}(\xi)) + \xi^z \right)^{-1} \left( \bar{g}'(x^{\text{eq}}(\xi)) + \xi^y - \xi^z x^* \right). \quad (28)$$

Theorem 8 allows to define

$$\phi'_x(\mathbf{x}, \xi) := \phi_x(\mathbf{x} + \mathbb{1}_N \otimes x^{\text{eq}}(\xi), \xi) \quad (29)$$

and the corresponding dynamics

$$\dot{\mathbf{x}} = \phi'_x(\mathbf{x}, \xi) \quad (30)$$

which corresponds to the translated version of the original  
 perturbed system  $\phi_x(\mathbf{x}, \xi)$ , which has now the property that the  
 origin is an equilibrium point, i.e.,  $\phi'_x(\mathbf{0}, \xi) = 0, \forall \|\xi\| \leq r$ .

To prove the global exponential stability of (30) we need the  
 flow  $\phi'_x$  to satisfy a global Lipschitz condition.

*Assumption 9 (Global Lipschitz Perturbation):* There exist  
 positive scalars  $a_\xi$  and  $r$  such that, for all  $\mathbf{x} \in \mathbb{R}^{nN}$  and  $\xi$   
 satisfying  $\|\xi\| \leq r$

$$\|\phi'_x(\mathbf{x}, \xi) - \phi'_x(\mathbf{x}, 0)\| \leq a_\xi \|\xi\| \|\mathbf{x}\|.$$

With these assumptions we can prove that the origin is a  
 globally exponentially stable equilibrium for dynamics (30).

*Theorem 10:* Under Assumptions 1, 5, and 9,

a)  $V_{\text{PNR}}(\mathbf{x})$  defined in (21) is a Lyapunov function for (30);

b) there exist positive scalars  $r, b'_7, b'_8$  s.t., for all  $\mathbf{x} \in \mathbb{R}^{nN}$   
 and  $\xi$  satisfying  $\|\xi\| \leq r$

$$\begin{aligned} & \begin{cases} \frac{\partial V_{\text{PNR}}}{\partial \mathbf{x}} \phi'_x(\mathbf{x}, \xi) \leq -b'_7 \|\mathbf{x}\|^2 \\ \|\phi'_x(\mathbf{x}, \xi)\| \leq b'_8 \|\mathbf{x}\|. \end{cases} \end{aligned} \quad \begin{aligned} (31a) \\ (31b) \end{aligned}$$

Again, as in Lemma 4, combining Theorem 10 with  
 Theorem 2 it is possible to claim that (30) and its discrete-time  
 counterpart are globally exponentially stable.

#### 416 H. Quadratic Functions

417 Before presenting the main algorithm, we show that  
418 quadratic costs satisfy all the previous assumptions. In fact, let  
419 us consider then

$$f_i(x) = \frac{1}{2}(x - d_i)^T A_i (x - d_i) + e_i, \quad A_i = A_i^T.$$

420 Based on this definition we have the following result.

421 *Theorem 11:* Quadratic costs that satisfy

$$A := \sum_i A_i > 0$$

422 satisfy Assumptions 1, 5, and 9 for  $h'_i(x) = \nabla^2 f'_i(x)$ .

#### 423 IV. NEWTON-RAPHSON CONSENSUS

424 In this section we provide an algorithm to distributively  
425 compute the minimizer of the function  $x^*$  defined in (2).  
426 The algorithm will be shown to converge to  $x^*$  even if  $x^* \neq 0$ .  
427 The proof of convergence will be based on the results derived  
428 in the previous sections via a suitable translation of the argu-  
429 ment of the cost functions, which basically reduces the problem  
430 to the special case  $x^* = 0$ .

431 Consider then Algorithm 1, where  $g(x(-1)) = \mathbf{0}$  and  
432  $h(x(-1)) = \mathbf{0}$  in the initialization step should be intended as  
433 initialization of suitable registers and not as operations involv-  
434 ing the quantity  $x(-1)$ .

---

#### 435 Algorithm 1 Fast Newton-Raphson Consensus (NRC)

---

436 *(storage allocation and constraints on the parameters)*

437 1:  $x_i(k), y_i(k) \in \mathbb{R}^n$  and  $z_i(k) \in \mathbb{R}^{n \times n}$  for all  $k$  and  $i =$   
438  $1, \dots, N; \varepsilon \in (0, 1], c > 0$

439 *(initialization)*

440 2:  $x_i(0) = 0; y_i(0) = g_i(x_i(-1)) = 0; z_i(0) = h_i(x_i(-1)) = 0$   
441 *(main algorithm)*

442 3: **for**  $k = 1, 2, \dots$  **do**

443 4:   **for**  $i = 1, \dots, N$  **do**

444 5:      $x_i(k) = (1 - \varepsilon)x_i(k-1) + \varepsilon[z_i(k-1)]_c^{-1}y_i(k-1)$

445 6:      $y_i(k) = \sum_{j=1}^N p_{ij}(y_j(k-1) + g_j(x_j(k-1)) -$   
446  $g_j(x_j(k-2)))$

447 7:      $z_i(k) = \sum_{j=1}^N p_{ij}(z_j(k-1) + h_j(x_j(k-1)) -$   
448  $h_j(x_j(k-2)))$

449 8:   **end for**

450 9: **end for**

---

451 Intuitively, the algorithm functions as follows: if the dynam-  
452 ics of the  $x_i(k)$ s is sufficiently slow w.r.t. the dynamics of the  
453  $y_i(k)$ s and  $z_i(k)$ s, then the two latter quantities tend to reach  
454 consensus. Then, the more these quantities reach consensus,  
455 the more the products  $[z_i(k)]_c^{-1}y_i(k)$  exhibit these two specific  
456 characteristics: *i*) being the same among the various agent;  
457 *ii*) representing Newton descent directions. Thus, the more the  
458  $y_i(k)$ s and  $z_i(k)$ s in Algorithm 1 are sufficiently close, the  
459 more the various  $x_i(k)$ s are driven by the same forcing term,  
460 that makes them converge to the same value, equal to the  
461 optimum  $x^*$ .

We now characterize the convergence properties of  
Algorithm 1. Let us define

$$\xi^y := \frac{1}{N} \sum_{i=1}^N (y_i(0) - g_i(x_i(-1)))$$

$$\xi^z := \frac{1}{N} \sum_{i=1}^N (z_i(0) - h_i(x_i(-1)))$$

then we have the following theorem.

*Theorem 12:* Consider the dynamics defined by Algorithm 1  
with possibly nonzero initial conditions. If  $\xi^y = 0$  and  $\xi^z = 0$ ,  
then under Assumptions 1 and 5 there exists a positive scalar  
 $\bar{\varepsilon} > 0$  such that Theorem 2 holds, i.e., the algorithm can be  
considered a forward-Euler discretization of a globally expo-  
nentially stable continuous dynamics. Thus the local estimates  
 $x_i(k)$  produced by the algorithm exponentially converge to the  
global minimizer, i.e.,

$$\lim_{k \rightarrow \infty} x_i(k) = x^* \quad \forall i = 1, \dots, N$$

for all  $\varepsilon \in (0, \bar{\varepsilon})$  and  $x_i(0) \in \mathbb{R}^n$ .

Consider now that, due to finite-precision issues, the quan-  
tities  $\xi^y$  and  $\xi^z$  may be non-null. Non-null initial  $\xi^y$  and  $\xi^z$   
will make the proposed algorithm converge to a point that,  
in general does not coincide with the global optimum  $x^*$ .  
Nonetheless in this case the computed solution, as a function  
of the initial conditions, is a smooth function and thus small  
errors in the initial conditions do not produce dramatic errors in  
the computation of the optimum.

*Theorem 13:* Consider the dynamics defined by Algorithm 1  
with possibly nonzero initial  $\xi^y$  and  $\xi^z$  but generic  $x_i(0)$ 's.  
Under Assumptions 1, 5, and 9 there exist positive scalars  $a, r$ ,  
 $\bar{\varepsilon}$  and a continuously differentiable function  $\Psi : \mathbb{R}^n \times \mathbb{R}^{n \times n} \mapsto$   
 $\mathbb{R}^n$  satisfying

$$\|\Psi(\xi^y, \xi^z) - x^*\| \leq a(\|\xi^y\| + \|\xi^z\|)$$

s.t. the local estimates exponentially converge to it, i.e.,

$$\lim_{k \rightarrow \infty} x_i(k) = \Psi(\xi^y, \xi^z) \quad \forall i = 1, \dots, N$$

for all  $\varepsilon \in (0, \bar{\varepsilon})$ , initial conditions  $x_i(0) \in \mathbb{R}^n$  and  $(\|\xi^y\| +$   
 $\|\xi^z\|) \leq r$ .

We notice that Theorem 13 ensures global convergence  
properties w.r.t. the initial conditions  $x_i(0)$ 's by requiring  
Assumptions 1, 5, and 9, while for the same convergence  
properties Theorem 12 requires only Assumptions 1 and 5. The  
difference is that Theorem 13 considers a non-null perturbation  
 $\xi$  and Assumption 9 is needed to cope with this additional  
perturbation term.

The Assumptions 1, 5, and 9 are not needed if only local  
convergence is sought. In fact, local differentiability, and there-  
fore local Lipschitzianity, of the cost functions  $f_i(x)$  at the  
minimizer  $x^*$  is sufficient to guarantee that Assumptions 5 and 9  
are locally valid. As so, the proof that the equilibrium point is  
a locally exponentially stable point is exactly the same, with  
the difference that all bounds and inequalities are local. This  
observation is summarized in the following theorem.

**Theorem 14:** Consider the dynamics defined by Algorithm 1 with possibly nonzero initial conditions. Under the assumptions that the  $f_i$ 's are  $\mathcal{C}^3$  and that  $\nabla^2 \bar{f}(x^*) \geq cI$ , there exist positive scalars  $a, r, \bar{\varepsilon}$  and a continuously differentiable function  $\Psi : \mathbb{R}^n \times \mathbb{R}^{n \times n} \mapsto \mathbb{R}^n$  s.t.

$$\lim_{k \rightarrow \infty} x_i(k) = \Psi(\xi^y, \xi^z) \quad \forall i = 1, \dots, N$$

and satisfying

$$\|\Psi(\xi^y, \xi^z) - x^*\| \leq a(\|\xi^y\| + \|\xi^z\|)$$

for all  $\varepsilon \in (0, \bar{\varepsilon})$  and initial conditions

$$\begin{aligned} \|x_i(0) - x^*\| &\leq r, \quad \|y_i - \bar{g}(x^*)\| \leq r, \quad \|z_i - \bar{h}(x^*)\| \leq r \\ \|g_i(x_i(-1)) - \bar{g}(x^*)\| &\leq r, \quad \|h_i(x_i(-1)) - \bar{h}(x^*)\| \leq r. \end{aligned}$$

Numerical simulations suggest that the algorithm is robust w.r.t. numerical errors and quantization noise. We also notice that Theorem 12 guarantees the existence of a critical value  $\bar{\varepsilon}$  but does not provide indications on its value. This is a known issue in all the systems dealing with separation of time scales. A standard rule of thumb is then to let the rate of convergence of the fast dynamics be sufficiently faster than the one of the slow dynamics, typically 2–10 times faster. In our algorithm the fast dynamics inherits the rate of convergence of the consensus matrix  $P$ , given by its spectral gap  $\sigma(P)$ , i.e., its spectral radius  $\rho(P) = 1 - \sigma(P)$ . The rate of convergence of the slow dynamics is instead governed by (18), which is nonlinear and therefore possibly depending on the initial conditions. However, close to the equilibrium point the dynamic behavior is approximately given by  $\dot{\bar{x}}(t) \approx -(\bar{x}(t) - x^*)$ , thus, since  $x_i(k) \approx \bar{x}(\varepsilon k)$ , then the convergence rate of the algorithm approximately given by  $1 - \varepsilon$ .

Thus we aim to let  $1 - \rho(P) \gg 1 - (1 - \varepsilon)$ , which provides the rule of thumb

$$\varepsilon \ll \sigma(P) \quad (32)$$

which is suitable for generic cost functions. We then notice that, although the spectral gap  $\sigma(P)$  might not be known in advance, it is possible to distributedly estimate it, see, e.g., [53]. However, such rule of thumb might be very conservative. In fact, if all the  $f_i$ 's are quadratic and are, w.l.o.g. s.t.  $\nabla^2 f_i \geq cI$ , then one can set  $\varepsilon = 1$  and neglect the thresholding  $[\cdot]_c$ , so that the procedure reduces to

$$\begin{aligned} x(k+1) &= \frac{y(k)}{z(k)} \\ y(k+1) &= (P \otimes I_n)y(k) \\ z(k+1) &= (P \otimes I_n)z(k) \end{aligned} \quad (33)$$

where  $x(k) := [x_1^T(k), \dots, x_N^T(k)]^T$ ,  $y(k) := [y_1^T(k), \dots, y_N^T(k)]^T$ ,  $z(k) := [z_1(k), \dots, z_N(k)]^T$ . Thus:

**Theorem 15:** Consider Algorithm 1 with arbitrary initial conditions  $x_i(0)$ , quadratic cost functions  $f_i = (1/2)(x - d_i)^T A_i (x - d_i)$  with  $A_i > 0$  and  $\varepsilon = 1$ . Then  $\|x_i(k) - x^*\| \leq \alpha(\rho(P))^k$  for all  $k, i$  and for a suitable positive  $\alpha$ .

Thus, if the cost functions are close to be quadratic then the overall rate of convergence is limited by the rate of convergence of the embedded consensus algorithm. Moreover, the values of  $\varepsilon$  that still guarantee convergence can be much larger than those dictated by the rule of thumb (32).

TABLE I  
COMPUTATIONAL, COMMUNICATION, AND MEMORY COSTS  
OF NRC, JC, GDC PER SINGLE UNIT AND SINGLE STEP

Choice	NRC, $h_i(x) = \nabla^2 f_i(x)$	JC, $h_i(x) = \text{diag}[\nabla^2 f_i(x)]$	GDC, $h_i(x) = I$
Computational Cost	$O(n^3)$	$O(n)$	$O(n)$
Communication Cost	$O(n^2)$	$O(n)$	$O(n)$
Memory Cost	$O(n^2)$	$O(n)$	$O(n)$

#### A. On the Selection of the Structure of $h(x)$

549

As introduced in Section III-D, by selecting different structures for  $h_i(x)$  one can obtain different procedures with different convergence properties and different computational/communication requirements. Plausible choices for  $h_i$  are the ones in (13c), and the correspondences are the following:

- $h_i(x) = \nabla^2 f_i(x) \rightarrow$  Newton-Raphson Consensus (NRC): in this case it is possible to rewrite the main algorithm and show that, for sufficiently small  $\varepsilon$ ,  $x_i(k) \approx \bar{x}(\varepsilon k)$ , where  $\bar{x}(t)$  evolves according to the continuous-time Newton-Raphson dynamics

$$\dot{\bar{x}}(t) = -[\nabla^2 \bar{f}(\bar{x}(t))]^{-1} \nabla \bar{f}(\bar{x}(t)).$$

- $h_i(x) = \text{diag}[\nabla^2 f_i(x)] \rightarrow$  Jacobi Consensus (JC): choice  $h_i(x) = \nabla^2 f_i(x)$  requires agents to exchange information on  $O(n^2)$  scalars, and this could pose problems under heavy communication bandwidth constraints and large  $n$ 's. Choice  $h_i(x) = \text{diag}[\nabla^2 f_i(x)]$  instead reduces the amount of information to be exchanged via the underlying diagonalization process, also called Jacobi approximation.<sup>3</sup> In this case, for sufficiently small  $\varepsilon$ ,  $x_i(k) \approx \bar{x}(\varepsilon k)$ , where  $\bar{x}(t)$  evolves according to the continuous-time dynamics

$$\dot{\bar{x}}(t) = -(\text{diag}[\nabla^2 \bar{f}(\bar{x}(t))])^{-1} \nabla \bar{f}(\bar{x}(t))$$

which can be shown to converge to the global optimum  $x^*$  with a convergence rate that in general is slower than the Newton-Raphson when the global cost function is skewed.

- $h_i(x) = I \rightarrow$  Gradient Descent Consensus (GDC): this choice is motivated in frameworks where the computation of the local second derivatives  $(\partial^2 f_i / \partial x_m^2)|_x$  is expensive (with  $x_m$  indicating here the  $m$ -th component of  $x$ ), or where the second derivatives simply might not be continuous. With this choice the main algorithm reduces to a distributed gradient-descent procedure. In fact, for sufficiently small  $\varepsilon$ ,  $x_i(k) \approx \bar{x}(\varepsilon k)$  with  $\bar{x}(t)$  evolving according to the continuous-time dynamics

$$\dot{\bar{x}}(t) = -\nabla \bar{f}(\bar{x}(t))$$

which one again is guaranteed to converge to the global optimum  $x^*$ .

The following Table I summarizes the various costs of the previously proposed strategies.

We remark that  $\bar{\varepsilon}$  in Theorem 12 depends also on the particular choice for  $h_i$ . The list of choices for  $h_i$  given above

<sup>3</sup>In centralized approaches, nulling the Hessian's off-diagonal terms is a well-known procedure, see, e.g., [54]. See also [36], [55] for other Jacobi algorithms with different communication structures.



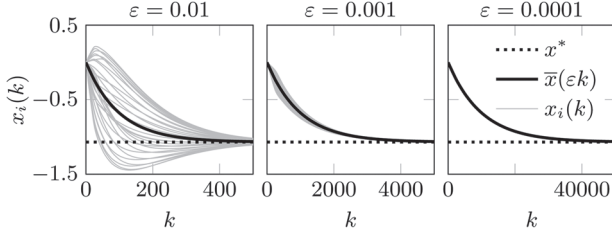


Fig. 1. Temporal evolution of system (43) for different values of  $\varepsilon$ , with  $N = 30$ . The black dotted line indicates  $x^*$ . The black solid line indicates the slow dynamics  $\bar{x}(\varepsilon k)$  of Equation (18). As  $\varepsilon$  decreases, the difference between the time scale of the slow and fast dynamics increases, and the local states  $x_i(k)$  converge to the manifold of  $\bar{x}(\varepsilon k)$ .

588 is not exhaustive. For example, future directions are to imple-  
589 ment distributed quasi-Newton procedures. To this regard, we  
590 recall that approximations of the Hessians that do not maintain  
591 symmetry and positive definiteness or are bad conditioned  
592 require additional modification steps, e.g., through Cholesky  
593 factorizations [56].

594 Finally, we notice that in scalar scenarios JC and NRC are  
595 equivalent, while GDC corresponds to algorithms requiring just  
596 the knowledge of first derivatives.

## V. NUMERICAL EXAMPLES

598 In Section V-A we analyze the effects of different choices  
599 of  $\varepsilon$  on the NRC on regular graphs and exponential cost  
600 functions. We then propose two machine learning problems in  
601 Section V-B, used in Section V-C and D, and numerically com-  
602 pare the convergence performance of the NRC, JC, GDC algo-  
603 rithms and other distributed convex optimization algorithms on  
604 random geometric graphs.

605 Notice that we will use cost functions that may not satisfy  
606 Assumptions 1, 5, and 9 to highlight the fact that the algorithm  
607 seems to have favorable numerical properties and large basins  
608 of stability even if the assumptions needed for global stability  
609 are not satisfied.

### 610 A. Effects of the Choice of $\varepsilon$

611 Consider a ring network of  $S = 30$  agents that communicate  
612 only to their left and right neighbors through the consensus  
613 matrix

$$P = \begin{bmatrix} 0.5 & 0.25 & & & 0.25 \\ 0.25 & 0.5 & 0.25 & & \\ & \ddots & \ddots & \ddots & \\ & & 0.25 & 0.5 & 0.25 \\ 0.25 & & & 0.25 & 0.5 \end{bmatrix} \quad (34)$$

614 so that the spectral radius  $\rho(P) \approx 0.99$ , implying a spectral gap  
615  $\sigma(P) \approx 0.01$ . Consider also scalar costs of the form  $f_i(x) =$   
616  $c_i e^{a_i x} + d_i e^{-b_i x}$ ,  $i = 1, \dots, N$ , with  $a_i, b_i \sim \mathcal{U}[0, 0.2]$ ,  $c_i,$   
617  $d_i \sim \mathcal{U}[0, 1]$  and where  $\mathcal{U}$  indicates the uniform distribution.

618 Fig. 1 compares the evolution of the local states  $x_i$  of the  
619 continuous system (43) for different values of  $\varepsilon$ . When  $\varepsilon$  is  
620 not sufficiently small, then the trajectories of  $x_i(t)$  are different  
621 even if they all start from the same initial condition  $x_i(0) = 0$ .  
622 As  $\varepsilon$  decreases, the difference between the two time scales be-  
623 comes more evident and all the trajectories  $x_i(k)$  become closer  
624 to the trajectory given by the slow NR dynamics  $\bar{x}(\varepsilon k)$  given in  
625 (18) and guaranteed to converge to the global optimum  $x^*$ .

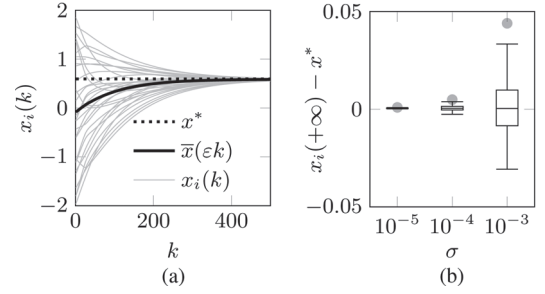


Fig. 2. Characterization of the dependency of the performance of Algorithm 1 on the initial conditions. In all the experiments  $\varepsilon = 0.01$  and  $N = 30$ . (a) Time evolution of the local states  $x_i(k)$  with  $v(0) = w(0) = y(0) = z(0) = 0$  and  $x_i(0) \sim \mathcal{U}[-2, 2]$ . (b) Empirical distribution of the errors  $x_i(+\infty) - x^*$  under artificially perturbed initial conditions  $\alpha(0), \beta(0) \sim \mathcal{U}[-\sigma, \sigma]$  for different values of  $\sigma$ .

In Fig. 2 we address the robustness of the proposed algorithm 626 w.r.t. the choice of the initial conditions. In particular, Fig. 2(a) 627 shows that if  $\alpha = \beta = 0$  then the local states  $x_i(t)$  converge to 628 the optimum  $x^*$  for arbitrary initial conditions  $x_i(0)$ . Fig. 2(b) 629 considers, besides different initial conditions  $x_i(0)$ , also per- 630 turbed initial conditions  $v(0), w(0), y(0), z(0)$  leading to non 631 null  $\alpha$ 's and  $\beta$ 's. More precisely we apply Algorithm 1 to dif- 632 ferent random initial conditions s.t.  $\alpha, \beta \sim \mathcal{U}[-\sigma, \sigma]$ . Fig. 2(b) 633 shows the boxplots of the errors  $x_i(+\infty) - x^*$  for different  $\sigma$ 's 634 based on 300 Monte Carlo runs with  $\varepsilon = 0.01$  and  $N = 30$ . 635

### B. Optimization Problems

The first problem considered is the distributed training of a 637 Binomial-Deviance based classifier, to be used, e.g., for spam- 638 nonspam classification tasks [57, Ch. 10.5]. More precisely, we 639 consider a database of emails  $E$ , where  $j$  is the email index, 640  $y_j = -1, 1$  denotes if the email  $j$  is considered spam or not, 641  $\chi_j \in \mathbb{R}^{n-1}$  numerically summarizes the  $n - 1$  features of the 642  $j$ -th email (how many times the words “money”, “dollars”, etc., 643 appear). If the  $E$  emails come from different users that do not 644 want to disclose their private information, then it is meaningful 645 to exploit the distributed optimization algorithms described in 646 the previous sections. More specifically, letting  $x = (x', x_0) \in$  647  $\mathbb{R}^{n-1} \times \mathbb{R}$  represents a generic classification hyperplane, train- 648 ing a Binomial-Deviance based classifier corresponds to solve a 649 distributed optimization problem where the local cost functions 650 are given by 651

$$f_i(x) := \sum_{j \in E_i} \log(1 + \exp(-y_j(\chi_j^T x' + x_0))) + \gamma \|x'\|_2^2 \quad (35)$$

where  $E_i$  is the set of emails available to agent  $i$ ,  $E = \cup_{i=1}^N E_i$ , 652 and  $\gamma$  is a global regularization parameter. In the following 653 numerical experiments we consider  $|E| = 5000$  emails from 654 the spam-nonspam UCI repository, available at <http://archive.ics.uci.edu/ml/datasets/Spambase>, randomly assigned to 30 dif- 655 ferent users communicating as in graph of Fig. 4. For each email 657 we consider 3 features (the frequency of words “make”, “ad- 658 dress”, “all”) so that the corresponding optimization problem is 659 4-dimensional. 660

The second problem considered is a regression problem 661 inspired by the UCI Housing dataset available at <http://archive.ics.uci.edu/ml/datasets/Housing>. In this task, an example  $\chi_j \in$  663  $\mathbb{R}^{n-1}$  is a vector representing some features of a house (e.g., 664



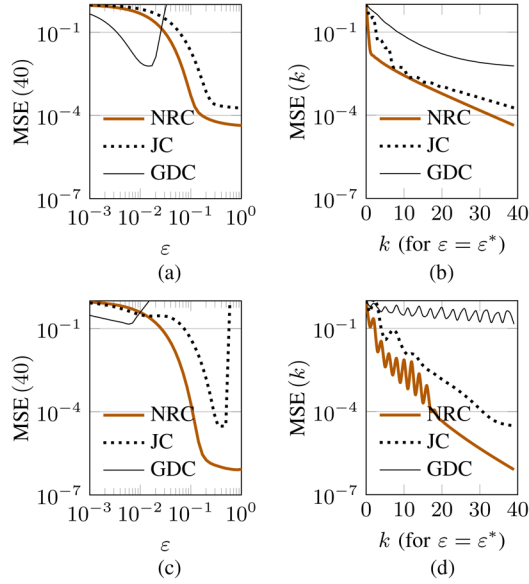


Fig. 3. Convergence properties of Algorithm 1 for the problems described in Section V-B and for different choices of  $h_i(\cdot)$ . Choice  $h_i(x) = \nabla^2 f_i(x)$  corresponds to the NRC algorithm,  $h_i(x) = \text{diag}[\nabla^2 f_i(x)]$  to the JC,  $h_i(x) = I$  to the GDC. (a) Relative MSE at a given time  $k$  as a function of the parameter  $\varepsilon$  for classification problem (35). (b) Relative MSE as a function of the time  $k$ , with the parameter  $\varepsilon$  chosen as the best from Fig. 3(a) for classification problem (35). (c) Relative MSE at a given time  $k$  as a function of the parameter  $\varepsilon$  for regression problem (36). (d) Relative MSE as a function of the time  $k$ , with the parameter  $\varepsilon$  chosen as the best from Fig. 3(c) for regression problem (36).

per capita crime rate by town, index of accessibility to radial highways, etc.), and  $y_j \in \mathbb{R}$  denotes the corresponding median monetary value of the house. The objective is to obtain a predictor of house value based on these data. Similarly as the previous example, if the datasets come from different users that do not want to disclose their private information, then it is meaningful to exploit the distributed optimization algorithms described in the previous sections. This problem can be formulated as a convex regression problem on the local costs

$$f_i(x) := \sum_{j \in E_i} \frac{(y_j - \chi_j^T x' - x_0)^2}{|y_j - \chi_j^T x' - x_0| + \beta} + \gamma \|x'\|_2^2 \quad (36)$$

where  $x = (x', x_0^T) \in \mathbb{R}^{n-1} \times \mathbb{R}$  is the vector of coefficient for the linear predictor  $\hat{y} = \chi^T x' + x_0$  and  $\gamma$  is a common regularization parameter. The loss function  $(\cdot)^2/(|\cdot| + \beta)$  responds to a smooth  $\mathcal{C}^2$  version of the Huber robust loss, a loss that is usually employed to minimize the effects of outliers. In our case  $\beta$  dictates for which arguments the loss is pseudo-linear or pseudo-quadratic and has been manually chosen to minimize the effects of outliers. In our experiments we used four features,  $\beta = 50$ ,  $\gamma = 1$ , and  $|E| = 506$  total number of examples in the dataset randomly assigned to the  $N = 30$  users communicating as in the graph of Fig. 4.

In both the previous problems the optimum, in the following indicated for simplicity with  $x^*$ , has been computed with a centralized NR with the termination rule “stop when in the last 5 steps the norm of the guessed  $x^*$  changed less than  $10^{-9}\%$ ”.

### C. Comparison of the NRC, JC and GDC Algorithms

In Fig. 3 we analyze the performance of the three proposed NRC, JC and GDC algorithms defined by the various choices

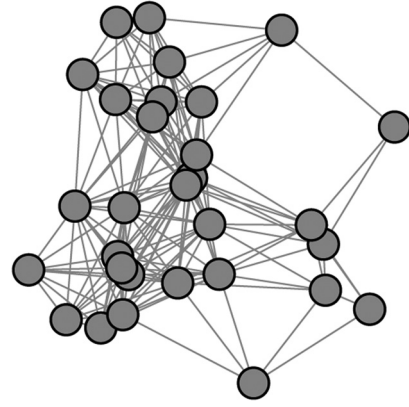


Fig. 4. Random geometric graph exploited in the simulations relative to the optimization problem (35). For this graph  $\rho(P) \approx 0.9338$ , with  $P$  the matrix of Metropolis weights.

for  $h_i(x)$  in Algorithm 1 in terms of the relative MSE

$$\text{MSE}(k) := \frac{1}{N} \sum_{i=1}^N \|x_i(k) - x^*\|^2 / \|x^*\|^2$$

for the classification and regression optimization problem described above. The consensus matrix  $P$  has been by selecting the Metropolis-Hastings weights which are consistent with the communication graph [58]. Panels 3(a) and 3(c) report the MSE obtained at a specific iteration ( $k = 40$ ) by the various algorithms, as a function of  $\varepsilon$ . These plots thus inspect the sensitivity w.r.t. the choice of the tuning parameters. Consistently with the theorems in the previous section, the GDC and JC algorithms are stable only for  $\varepsilon$  sufficiently small, while NRC exhibit much larger robustness and best performance for  $\varepsilon = 1$ . Panels 3(b) and 3(d) instead report the evolutions of the relative MSE as a function of the number of iterations  $k$  for the optimally tuned algorithms.

We notice that the differences between NRC and JC are evident but not resounding, due to the fact that the Jacobi approximations are in this case a good approximation of the analytical Hessians. Conversely, GDC presents a slower convergence rate which is a known drawback of gradient descent algorithms.

### D. Comparisons With Other Distributed Convex Optimization Algorithms

We now compare Algorithm 1 and its accelerated version, referred as Fast Newton-Raphson Consensus (FNRC) and described in detail below in Algorithm 2), with three popular distributed convex optimization methods, namely the DSM, the Distributed Control Method (DCM) and the ADMM, described respectively in Algorithm 3, 4, and 5. The following discussion provides some details about these strategies.

#### Algorithm 2 Fast Newton-Raphson Consensus

- 1: storage allocation, constraints on the parameters and initialization as in Algorithm 1
- 2: **for**  $k = 1, 2, \dots$  **do**
- 3:   **for**  $i = 1, \dots, N$  **do**
- 4:      $x_i(k) = (1 - \varepsilon)x_i(k-1) + \varepsilon[z_i(k-1)]_c^{-1}y_i(k-1)$
- 5:      $\tilde{y}_i(k) = y_i(k-1) + (1/\varphi)g_i(x_i(k-1)) - g_i(x_i(k-2)) - ((1 - \varphi)/\varphi)g_i(x_i(k-3))$

---

```

728 6:  $\tilde{z}_i(k) = z_i(k-1) + (1/\varphi)h_i(x_i(k-1)) - h_i(x_i(k-1))$ 
729  $2)) - ((1-\varphi)/\varphi)h_i(x_i(k-3))$ 
730 7:  $y_i(k) = \varphi \sum_{j=1}^N (p_{ij}\tilde{y}_j(k)) + (1-\varphi)y_i(k-2)$ 
731 8:  $z_i(k) = \varphi \sum_{j=1}^N (p_{ij}\tilde{z}_j(k)) + (1-\varphi)z_i(k-2)$ 
732 9: end for
733 10: end for

```

---



---

#### 734 Algorithm 3 DSM [30]

---

```

735 (storage allocation and constraints on parameters)
736 1:  $x_i(k) \in \mathbb{R}^n$  for all  $i$ .  $\varrho \in \mathbb{R}_+$ 
737 (initialization)
738 2:  $x_i(0) = 0$ 
739 (main algorithm)
740 3: for  $k = 0, 1, \dots$  do
741 4:   for  $i = 1, \dots, N$  do
742 5:      $x_i(k+1) = \sum_{j=1}^N p_{ij}(x_j(k) - (\varrho/k)\nabla f_j(x_j(k)))$ 
743 6:   end for
744 7: end for

```

---



---

#### 745 Algorithm 4 DCM [42]

---

```

746 (storage allocation and constraints on parameters)
747 1:  $x_i(k), z_i(k) \in \mathbb{R}^n$ , for all  $i$ .  $\mu, \nu \in \mathbb{R}_+$ 
748 (initialization)
749 2:  $x_i(0) = z_i(0) = 0$  for all  $i$ 
750 (main algorithm)
751 3: for  $k = 0, 1, \dots$  do
752 4:   for  $i = 1, \dots, N$  do
753 5:      $z_i(k+1) = z_i(k) + \mu \sum_{j \in \mathcal{N}_i} (x_i(k) - x_j(k))$ 
754 6:      $x_i(k+1) = x_i(k) + \mu \sum_{j \in \mathcal{N}_i} (x_j(k) - x_i(k)) +$ 
755  $\mu \sum_{j \in \mathcal{N}_i} (z_j(k) - z_i(k)) - \mu\nu \nabla f_i(x_i(k))$ 
756 7:   end for
757 8: end for

```

---



---

#### 758 Algorithm 5 ADMM [7, pp. 253–261]

---

```

759 (storage allocation and constraints on parameters)
760 1:  $x_i(k), z_{(i,j)}(k), y_{(i,j)}(k) \in \mathbb{R}^n$ ,  $\delta \in (0, 1)$ 
761 (initialization)
762 2:  $x_i(k) = z_{(i,j)}(k) = y_{(i,j)}(k) = 0$ 
763 (main algorithm)
764 3: for  $k = 0, 1, \dots$  do
765 4:   for  $i = 1, \dots, N$  do
766 5:      $x_i(k+1) = \arg \min_{x_i} L_i(x_i, k)$ 
767 6:     for  $j \in \mathcal{N}_i$  do
768 7:        $z_{(i,j)}(k+1) = (1/2\delta)(y_{(i,j)}(k) + y_{(j,i)}(k)) +$ 
769  $(1/2)(x_i(k+1) + x_j(k+1))$ 
770 8:        $y_{(i,j)}(k+1) = y_{(i,j)}(k) + \delta(x_i(k+1) -$ 
771  $z_{(i,j)}(k+1))$ 
772 9:     end for
773 10:   end for
774 11: end for

```

---

775 • FNRC is an accelerated version of Algorithm 1 that in-  
776 herits the structure of the so called *second order diffusive*

*schedules*, see, e.g., [59], and exploits an additional level 777  
of memory to speed up the convergence properties of the 778  
consensus strategy. Here the weights multiplying the  $g_i$ 's 779  
and  $h_i$ 's are necessary to guarantee exact tracking of the 780  
current average, i.e.,  $\sum_i y_i(k) = \sum_i g_i(x(k-1))$  for all  $k$ . 781  
As suggested in [59], we set the  $\varphi$  that weights the gra- 782  
dient and the memory to  $\varphi = 2/(1 + \sqrt{1 - \rho(P)^2})$ . This 783  
guarantees second order diffusive schedules to be faster 784  
than first order ones (even if this does not automatically 785  
imply the FNRC to be faster than the NRC). This setting 786  
can be considered a valid heuristic to be used when  $\rho(P)$  787  
is known. For the graph in Fig. 4,  $\varphi \approx 1.4730$ . 788

- DSM, as proposed in [30], alternates consensus steps on 789  
the current estimated global minimum  $x_i(k)$  with subgra- 790  
dient updates of each  $x_i(k)$  towards the local minimum. To 791  
guarantee the convergence, the amplitude of the local sub- 792  
gradient steps should appropriately decrease. Algorithm 3 793  
presents a synchronous DSM implementation, where  $\varrho$  is 794  
a tuning parameter and  $P$  is the matrix of Metropolis- 795  
Hastings weights. 796
- DCM, as proposed in [42], differentiates from the gradient 797  
searching because it forces the states to the global opti- 798  
mum by controlling the subgradient of the global cost. 799  
This approach views the subgradient as an input/output 800  
map and uses small gain theorems to guarantee the conver- 801  
gence property of the system. Again, each agents  $i$  locally 802  
computes and exchanges information with its neighbors, 803  
collected in the set  $\mathcal{N}_i := \{j \mid (i, j) \in \mathcal{E}\}$ . DCM is sum- 804  
marized in Algorithm 4, where  $\mu, \nu > 0$  are parameters 805  
to be designed to ensure the stability property of the 806  
system. Specifically,  $\mu$  is chosen in the interval  $0 < \mu < 807$   
 $2/(2 \max_{i \in \{1, \dots, N\}} |\mathcal{N}_i| + 1)$  to bound the induced gain 808  
of the subgradients. Also here the parameters have been 809  
manually tuned for best convergence rates. 810
- ADMM, instead, requires the augmentation of the system 811  
through additional constraints that do not change the op- 812  
timal solution but allow the Lagrangian formalism. There 813  
exist different implementations of ADMM in distributed 814  
contexts, see, e.g., [7], [12, pp. 253–261], [60]. For sim- 815  
plicity we consider the following formulation: 816

$$\begin{aligned}
& \min_{x_1, \dots, x_N} \sum_{i=1}^N f_i(x_i) \\
& \text{s.t. } z_{(i,j)} = x_i, \quad \forall i \in \mathcal{N}, \quad \forall (i, j) \in \mathcal{E}
\end{aligned}$$

where the auxiliary variables  $z_{(i,j)}$  correspond to the dif- 817  
ferent links in the network, and where the local Aug- 818  
mented Lagrangian is given by 819

$$L_i(x_i, k) := f_i(x_i) + \sum_{j \in \mathcal{N}_i} y_{(i,j)}(x_i - z_{(i,j)}) + \sum_{\substack{(i,j) \in \mathcal{E} \\ i \in \mathcal{N}_i}} \frac{\delta}{2} \|x_i - z_{(i,j)}\|^2$$

with  $\delta$  a tuning parameter (see [61] for a discussion on how 820  
to tune it) and the  $y_{(i,j)}$ 's Lagrange multipliers. 821

The computational, communication and memory costs of 822  
these algorithms is reported in Table II. Notice that the com- 823  
putational and memory costs of ADMM algorithms depends on 824  
how nodes minimize the local augmented Lagrangian  $L_i(x_i, k)$ . 825  
E.g., in our simulations the step has been performed through a 826  
dedicated Newton-Raphson procedure with associated  $O(n^3)$  827  
computational costs and  $O(n^2)$  memory costs. 828

TABLE II  
COMPUTATIONAL, COMMUNICATION AND MEMORY COSTS OF DSM,  
DCM, AND ADMM PER SINGLE UNIT AND SINGLE STEP

Choice	DSM	DCM	ADMM
Computational Cost	$O(n)$	$O(n)$	from $O(n)$ to $O(n^3)$
Communication Cost	$O(n)$	$O(n)$	$O(n)$
Memory Cost	$O(n)$	$O(n)$	from $O(n)$ to $O(n^2)$

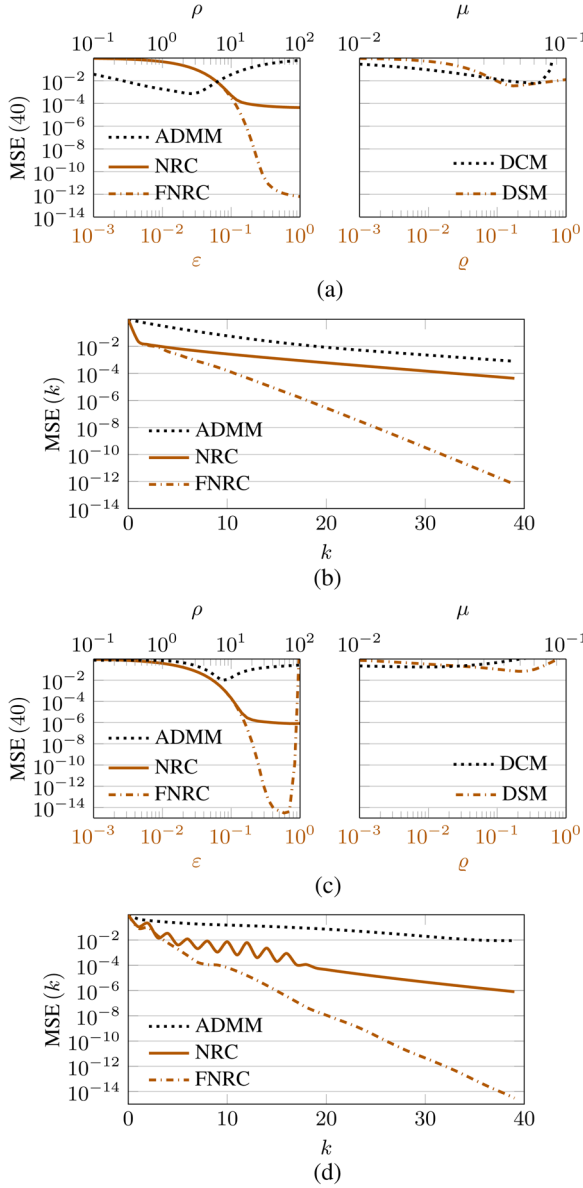


Fig. 5. Convergence properties of the various algorithms for the problems described in Section V-B. (a) Relative MSE at a given time  $k$  as a function of the algorithms parameters for problem (35). For the DCM,  $\nu = 1.7$ . (b) Relative MSE as a function of the time  $k$  for the three fastest algorithms for problem (35). Their parameters are chosen as the best ones from Fig. 5(a). (c) Relative MSE at a given time  $k$  as a function of the algorithms parameters for problem (36). For the DCM,  $\nu = 1.7$ . (d) Relative MSE as a function of the time  $k$  for the three fastest algorithms for problem (36). Their parameters are chosen as the best ones from Fig. 5(c).

Fig. 5 then compares the previously cited algorithms as did in Fig. 3. The first panel thus reports the relative MSE of the various algorithms at a given number of iterations ( $k = 40$ ) as a function of the parameters. The second panel instead reports the temporal evolution of the relative MSE for the case of optimal tuning.

We notice that the DCM and the DSM are both much slower, in terms of communications iterations, than the NRC, FNRC and ADMM. Moreover, both the NRC and its accelerated version converge faster than the ADMM, even if not tuned at their best. These numerical examples seem to indicate that the proposed NRC might be a viable alternative to the ADMM, although further comparisons are needed to strengthen this claim. Moreover, a substantial potential advantage of NRC compared to ADMM is that the former can be readily adapted to asynchronous and time-varying graphs, as preliminary made in [62]. Moreover, as in the case of the FNRC, the strategy to implement any improved linear consensus algorithm.

## VI. CONCLUSION

We proposed a novel distributed optimization strategy suitable for convex, unconstrained, multidimensional, smooth and separable cost functions. The algorithm does not rely on Lagrangian formalisms and acts as a distributed Newton-Raphson optimization strategy by repeating the following steps: agents first locally compute and update second order Taylor expansions around the current local guesses and then they suitably combine them by means of average consensus algorithms to obtain a sort of approximated Taylor expansion of the global cost. This allows each agent to infer a local Newton direction used to locally update the guess of the global minimum.

Importantly, the average consensus protocols and the local updates steps have different time-scales, and the whole algorithm is proved to be convergent only if the step-size is sufficiently slow. Numerical simulations based on real-world databases show that, if suitably tuned, the proposed algorithm is faster than ADMMs in terms of number of communication iterations, although no theoretical proof is provided.

The set of open research paths is extremely vast. We envisage three main avenues. The first one is to study how the agents can dynamically and locally tune the speed of the local updates w.r.t. the consensus process, namely how to tune their local step-size  $\epsilon_i$ . In fact large values of  $\epsilon$  gives faster convergence but might lead to instability. A second one is to let the communication protocol be asynchronous: in this regard we notice that some preliminary attempts can be found in [62]. A final branch is about the analytical characterization of the rate of convergence of the proposed strategies, a theoretical comparison with ADMMs, and the extensions to non-smooth convex functions.

## APPENDIX

*Proof (of Theorem 2): Proof of a):* integrating (5a) twice implies

$$\frac{1}{2}a_1\|x\|^2 \leq V(x) \leq \frac{1}{2}a_2\|x\|^2$$

that, jointly with (5b), immediately guarantee global exponential stability for (4) [46, Thm. 4.10].

*Proof of b):* consider

$$\Delta V(x(k)) := V(x(k+1)) - V(x(k)). \quad (37)$$

To prove the claim we show that  $\Delta V(x(k)) \leq -d\|x(k)\|^2$  for some positive scalar  $d$ . To this aim, expand  $V(x(k+1))$  with a



second order Taylor expansion around  $x(k)$  with remainder in Lagrange form, to obtain

$$V(x + \varepsilon \phi(x)) = V(x) + \varepsilon \frac{\partial V}{\partial x} \phi(x) + \frac{1}{2} \varepsilon^2 \phi^T(x) \nabla^2 V(x_\varepsilon) \phi(x)$$

with  $x_\varepsilon = x + \varepsilon' \phi(x)$  for  $\varepsilon' \in [0, \varepsilon]$ . Using inequalities (5) we then obtain

$$\begin{aligned} \Delta V(x(k)) &= V(x(k+1)) - V(x(k)) \\ &\leq -\varepsilon a_3 \|x(k)\|^2 + \frac{1}{2} \varepsilon^2 a_2 a_4^2 \|x(k)\|^2 \\ &= -\varepsilon \left( a_3 - \varepsilon \frac{1}{2} a_2 a_4^2 \right) \|x(k)\|^2. \end{aligned}$$

Thus, for all  $\varepsilon < \bar{\varepsilon} = 2a_3/a_2a_4^2$  the origin is globally exponentially stable. ■

*Proof (of Theorem 3): Proof of a):* Assumption 1 guarantees that  $V_{\text{NR}}(0) = 0$  and  $V_{\text{NR}}(x) > 0$  for  $x \neq 0$ . Moreover, for  $x \neq 0$

$$\begin{aligned} \frac{\partial V_{\text{NR}}}{\partial x} \phi_{\text{PNR}}(x) &= -(\nabla \bar{f}'(x))^T \bar{h}'(x)^{-1} \nabla \bar{f}'(x) \\ &= -\left\| \bar{h}'(x)^{-\frac{1}{2}} \nabla \bar{f}'(x) \right\|^2 < 0. \end{aligned}$$

*Proof of b):* Assumption 1 guarantees that (11a) is satisfied with  $b_1 = c$  and  $b_2 = m$ . To prove (11c) we start by considering that (11a) guarantees  $c\|x\| \leq \|\nabla \bar{f}'(x)\| \leq m\|x\|$ . This in its turn implies

$$\begin{aligned} \|\phi_{\text{NR}}(x)\| &= \left\| \bar{h}'^{-1}(x) \nabla \bar{f}'(x) \right\| \leq \frac{1}{c} \|\nabla \bar{f}'(x)\| \leq \frac{m}{c} \|x\| \\ &= b_4 \|x\|. \end{aligned}$$

To prove (11b) eventually consider then that (11c) implies

$$\begin{aligned} \frac{\partial V_{\text{NR}}}{\partial x} \phi_{\text{NR}}(x) &= -(\nabla \bar{f}'(x))^T \bar{h}'(x)^{-1} \nabla \bar{f}'(x) \\ &\leq -\frac{c^2}{m} \|x\|^2 = -b_3 \|x\|^2. \end{aligned}$$

■

*Proof (of Theorem 6):* In the interest of clarity we analyze the case where the local costs  $f'_i$  are scalar, i.e.,  $n = 1$ . The multivariable case is indeed a straightforward extension with just a more involved notation. We also recall the following equivalences:

$$\begin{aligned} x &= x^\parallel + x^\perp, \quad (x^\perp)^T x^\parallel = 0 \\ \|x\|^2 &= \|x^\parallel\|^2 + \|x^\perp\|^2 = N|\bar{x}|^2 + \|x^\perp\|^2. \end{aligned}$$

*Proof of a):*  $V_{\text{PNR}}(0) = 0$  and  $V_{\text{PNR}}(x) > 0$  for  $x \neq 0$  follow immediately from the fact that  $V_{\text{NR}}(0) = 0$  and  $V_{\text{NR}}(\bar{x}) > 0$  for  $\bar{x} \neq 0$ .  $\dot{V}_{\text{PNR}} < 0$  is instead proved by proving (22b).

*Proof of Inequality (22a):* given (21)

$$\frac{\partial^2 V_{\text{PNR}}(x)}{\partial x^2} = \frac{\partial^2 (V_{\text{NR}}(\bar{x}) + \frac{1}{2} \eta \|x^\perp\|^2)}{\partial x^2}.$$

Since  $0 \leq \|x^\perp\|^2 \leq \|x\|^2$  and

$$\frac{\partial^2 V_{\text{NR}}(\bar{x})}{\partial x^2} = \frac{1}{N^2} \mathbb{1}^T \nabla^2 V_{\text{NR}}(\bar{x})$$

thanks to (11a) it follows immediately that (22a) holds with

$$b_5 := \min \left\{ \frac{b_1}{N}, \eta \right\}, \quad b_6 := \max \left\{ \frac{b_2}{N}, \eta \right\}.$$

*Proof of Inequality (22c):* since the origin of  $\bar{f}'$  is a minimum, it follows that  $\nabla \bar{f}'(0) = 0$ , and thus  $\bar{g}'(0) = 0$  [cf. (14)]. Thus also  $\psi(0) = 0$ , that in turn implies  $\|\psi(x)\| \leq a_\psi \|x\|$  by Assumption 5. Therefore

$$\|\phi_{\text{PNR}}(x)\| \leq \|x\| + N \|\psi(x)\| \leq (1 + Na_\psi) \|x\| = b_8 \|x\|.$$

*Proof of Inequality (22b):* since

$$\frac{\partial \bar{x}}{\partial x} = \frac{1}{N} \mathbb{1}_N^T, \quad \frac{\partial x^\perp}{\partial x} = I - \frac{1}{N} \mathbb{1}_N \mathbb{1}_N^T =: \Pi$$

it follows that:

$$\begin{aligned} \frac{\partial V_{\text{PNR}}}{\partial x} \phi_{\text{PNR}}(x) &= \left( \frac{\partial V_{\text{PNR}}}{\partial \bar{x}} \frac{\partial \bar{x}}{\partial x} + \frac{\partial V_{\text{PNR}}}{\partial x^\perp} \frac{\partial x^\perp}{\partial x} \right) \phi_{\text{PNR}}(x) \\ &= \left( \frac{\partial V_{\text{NR}}(\bar{x})}{\partial \bar{x}} \frac{1}{N} \mathbb{1}_N^T + \eta (x^\perp)^T \Pi \right) \phi_{\text{PNR}}(x). \end{aligned}$$

Considering then (17), the definition of  $\bar{x}$  and  $x^\perp$ , and the fact that  $\Pi \mathbb{1}_N = 0$ , it follows that:

$$\frac{\partial V_{\text{PNR}}}{\partial x} \phi_{\text{PNR}}(x) = \frac{\partial V_{\text{NR}}(\bar{x})}{\partial \bar{x}} (-\bar{x} + \psi(x)) + \eta (x^\perp)^T (-x^\perp)$$

Adding and subtracting  $(\partial V_{\text{NR}}(\bar{x})/\partial \bar{x}) \psi(x^\parallel)$ , and recalling definition (7) and equivalence (16c), since  $(-\bar{x} + \psi(x^\parallel)) = \phi_{\text{NR}}(\bar{x})$  it then follows that:

$$\begin{aligned} \frac{\partial V_{\text{PNR}}}{\partial x} \phi_{\text{PNR}}(x) &= \frac{\partial V_{\text{NR}}(\bar{x})}{\partial \bar{x}} \phi_{\text{NR}}(\bar{x}) - \eta \|x^\perp\|^2 \\ &\quad + \frac{\partial V_{\text{NR}}(\bar{x})}{\partial \bar{x}} (\psi(x) - \psi(x^\parallel)) \\ &\leq -b_3 \bar{x}^2 - \eta \|x^\perp\|^2 + b_2 |\bar{x}| a_\psi \|x - x^\parallel\| \\ &= -b_3 \bar{x}^2 - \eta \|x^\perp\|^2 + b_2 a_\psi |\bar{x}| \|x^\perp\| \\ &\leq -\frac{b_3 + \eta}{2} (|\bar{x}|^2 + \|x^\perp\|^2) \\ &\leq -\frac{b_3 + \eta}{2} (N|\bar{x}|^2 + \|x^\perp\|^2) \\ &= -\frac{b_3 + \eta}{2N} \|x\|^2 = -b_7 \|x\|^2 \end{aligned}$$

where for obtaining the various inequalities we used the various assumptions and where the second inequality is valid for  $\eta > b_2^2 a_\psi^2 / b_3$ . ■

*Proof (of Lemma 7): Proof of (26a):* notice that  $\phi_x(x, \chi)$  is globally defined since  $[\cdot]_c$  ensures that the matrix inverse exists. Also note that, since  $h'(x) \geq cI > (c/2)I$  by Assumption 5, then there exists  $r > 0$  such that, for  $\|x\| + \|\chi\| \leq r$

$$\phi_x(x, \chi) = -x - \mathbb{1}_N \otimes x^* + \frac{\chi^y + \mathbb{1}_N \otimes (\bar{g}'(x) + \bar{h}'(x)x^*)}{\chi^z + \mathbb{1}_N \otimes \bar{h}'(x)}.$$

The differentiability of the elements defining  $\phi_x$ , plus the fact that  $[\cdot]_c$  acts as the identity in the neighborhood under consideration implies that  $\phi_x$  is locally differentiable in  $\|x\| + \|\chi\| \leq r$ .

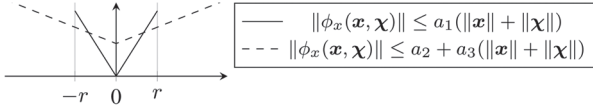


Fig. 6. Inequality (38) represents a proper cone defined in the neighborhood of radius  $r$ , while inequality (39) represents an improper cone defined in the whole domain.

In addition to this local differentiability, also observe that  $\phi_x(\mathbf{0}, \mathbf{0}) = 0$ , therefore there must exist  $a_1 > 0$  s.t.

$$\|\phi_x(\mathbf{x}, \boldsymbol{\chi})\| \leq a_1 (\|\mathbf{x}\| + \|\boldsymbol{\chi}\|), \quad \forall (\|\mathbf{x}\| + \|\boldsymbol{\chi}\|) \leq r. \quad (38)$$

To extend the linear inequality (38) for  $(\mathbf{x}, \boldsymbol{\chi})$  s.t.  $(\|\mathbf{x}\| + \|\boldsymbol{\chi}\|) \geq r$  we then prove that  $\phi_x(\mathbf{x}, \boldsymbol{\chi})$  cannot grow more than linearly globally. In fact

$$\begin{aligned} \|\phi_x(\mathbf{x}, \boldsymbol{\chi})\| &\leq \|\mathbf{x}\| + N\|x^*\| + \frac{2}{c} \|\boldsymbol{\chi}^y + \mathbf{1}_N \otimes (\bar{g}'(\mathbf{x}) + \bar{h}'(\mathbf{x})x^*)\| \\ &\leq \|\mathbf{x}\| + N\|x^*\| + \frac{2}{c} \|\boldsymbol{\chi}\| \\ &\quad + \frac{2N}{c} (\|\bar{g}'(\mathbf{x})\| + \|x^*\| \|\bar{h}'(\mathbf{x})\|) \\ &\leq \|\mathbf{x}\| + N\|x^*\| + \frac{2}{c} \|\boldsymbol{\chi}\| + \frac{2N}{c} a_g \|\mathbf{x}\| \\ &\quad + \frac{2N}{c} \|x^*\| (a_h \|\mathbf{x}\| + \|\bar{h}'(0)\|) \\ &\leq a_2 + a_3 (\|\mathbf{x}\| + \|\boldsymbol{\chi}\|), \quad \forall \mathbf{x}, \boldsymbol{\chi} \end{aligned} \quad (39)$$

where we used Assumption 5 and where  $a_2, a_3$  are suitable positive scalars. In particular inequality (39) is valid for  $(\|\mathbf{x}\| + \|\boldsymbol{\chi}\|) > r$ . As depicted in Fig. 6, inequalities (38) and (39) define two cones, one affine (shifted by  $a_2$ ) and one proper.

Therefore, combining the geometry of the two cones leads to an inequality that is defined in the whole domain. In other words, it follows that:

$$\|\phi_x(\mathbf{x}, \boldsymbol{\chi})\| \leq a_x (\|\mathbf{x}\| + \|\boldsymbol{\chi}\|) \quad \forall \mathbf{x}, \boldsymbol{\chi}$$

where

$$a_x := \max \left\{ a_1, \frac{a_2 + a_3 r}{r} \right\}.$$

*Proof of (26b):* Let  $\Delta(\mathbf{x}, \boldsymbol{\chi}) := \phi_x(\mathbf{x}, \boldsymbol{\chi}) - \phi_{\text{PNR}}(\mathbf{x})$ , with  $\phi_{\text{PNR}}$  as in (17). Then there exists a positive scalar  $r > 0$  such that, for all  $\|\boldsymbol{\chi}\| + \|\mathbf{x}\| \leq r$

$$\begin{aligned} \Delta(\mathbf{x}, \boldsymbol{\chi}) &= -\mathbf{1}_N \otimes x^* + \frac{\boldsymbol{\chi}^y + \mathbf{1}_N \otimes (\bar{g}'(\mathbf{x}) + \bar{h}'(\mathbf{x})x^*)}{\boldsymbol{\chi}^z + \mathbf{1}_N \otimes \bar{h}'(\mathbf{x})} \\ &\quad - \mathbf{1}_N \otimes \psi(\mathbf{x}) \\ &= \frac{\boldsymbol{\chi}^y + \mathbf{1}_N \otimes (\bar{g}'(\mathbf{x}) + \bar{h}'(\mathbf{x})x^*)}{\boldsymbol{\chi}^z + \mathbf{1}_N \otimes \bar{h}'(\mathbf{x})} \\ &\quad - \frac{\mathbf{1}_N \otimes (\bar{g}'(\mathbf{x}) + \bar{h}'(\mathbf{x})x^*)}{\mathbf{1}_N \otimes \bar{h}'(\mathbf{x})}. \end{aligned}$$

Considerations similar to the ones that led us claim the differentiability of  $\phi_x$  in the proof of Lemma 7 imply that  $\Delta(\mathbf{x}, \boldsymbol{\chi})$  is continuously differentiable for  $\|\boldsymbol{\chi}\| + \|\mathbf{x}\| \leq r$ . Moreover, since  $\Delta(\mathbf{x}, \mathbf{0}) = 0$ , then there exists a positive scalar  $a_4 > 0$  s.t.

$$\|\Delta(\mathbf{x}, \boldsymbol{\chi})\| \leq a_4 \|\boldsymbol{\chi}\| \quad \|\boldsymbol{\chi}\| + \|\mathbf{x}\| \leq r. \quad (40)$$

By using (19a) and (19b) we can then show that  $\Delta(\mathbf{x}, \boldsymbol{\chi})$  cannot grow more than linearly in the variable  $\boldsymbol{\chi}$ , since

$$\begin{aligned} \|\Delta(\mathbf{x}, \boldsymbol{\chi})\| &= \left\| \frac{\boldsymbol{\chi}^y + \mathbf{1}_N \otimes (\bar{g}'(\mathbf{x}) + \bar{h}'(\mathbf{x})x^*)}{[\boldsymbol{\chi}^z + \mathbf{1}_N \otimes \bar{h}'(\mathbf{x})]_c} \right. \\ &\quad \left. - \mathbf{1}_N \otimes \left( x^* + \frac{\bar{g}'(\mathbf{x})}{\bar{h}'(\mathbf{x})} \right) \right\| \\ &\leq \frac{2}{c} (\|\boldsymbol{\chi}\| + 2N\|\bar{g}'(\mathbf{x})\| + N\|x^*\| \|\bar{h}'(\mathbf{x})\|) + N\|x^*\| \\ &\leq a_5 + a_6 \|\boldsymbol{\chi}\|, \quad \forall \mathbf{x}, \boldsymbol{\chi} \end{aligned} \quad (41)$$

for suitable positive scalars  $a_5$  and  $a_6$ . Repeating the same geometrical arguments used above we then obtain

$$\|\Delta(\mathbf{x}, \boldsymbol{\chi})\| \leq a_\Delta \|\boldsymbol{\chi}\|, \quad \forall \mathbf{x}, \boldsymbol{\chi}$$

with

$$a_\Delta := \max \left\{ a_3, \frac{a_5 + a_6 r}{r} \right\}.$$

■

*Proof (of Theorem 8):* For notational brevity we omit the dependence on  $\xi$ , i.e., let  $\mathbf{x}^{\text{eq}} = \mathbf{x}^{\text{eq}}(\xi)$  and  $x^{\text{eq}} = x^{\text{eq}}(\xi)$ .

We start by assuming that there exists a  $\mathbf{x}^{\text{eq}}(\xi)$  satisfying (27) for  $\|\xi\| \leq r$  and prove that  $\mathbf{x}^{\text{eq}}(\xi)$  must satisfy  $\mathbf{x}^{\text{eq}}(\xi) = \mathbf{1}_N \otimes x^{\text{eq}}(\xi)$  and (28). Consider then  $r$  sufficiently small. Then, since  $\bar{h}'(\mathbf{x}) > cI$  by Assumption 1

$$[\boldsymbol{\xi}^z + \mathbf{1}_N \otimes \bar{h}'(\mathbf{x})]_c = \boldsymbol{\xi}^z + \mathbf{1}_N \otimes \bar{h}'(\mathbf{x}) = \mathbf{1}_N \otimes (\bar{h}'(\mathbf{x}) + \boldsymbol{\xi}^z).$$

This implies that for  $\|\xi\| \leq r$  we have

$$\begin{aligned} \phi_x(\mathbf{x}^{\text{eq}}, \boldsymbol{\xi}) &= -\mathbf{x}^{\text{eq}} \\ &\quad - \mathbf{1}_N \otimes \left( x^* - (\boldsymbol{\xi}^z + \bar{h}'(\mathbf{x}^{\text{eq}}))^{-1} (\boldsymbol{\xi}^y + \bar{g}'(\mathbf{x}^{\text{eq}}) + \bar{h}'(\mathbf{x}^{\text{eq}})x^*) \right). \end{aligned}$$

Therefore  $\phi_x(\mathbf{x}^{\text{eq}}, \boldsymbol{\xi}) = 0$  if and only if

$$x_i^{\text{eq}} = -x^* + (\boldsymbol{\xi}^z + \bar{h}'(\mathbf{x}^{\text{eq}}))^{-1} (\boldsymbol{\xi}^y + \bar{g}'(\mathbf{x}^{\text{eq}}) + \bar{h}'(\mathbf{x}^{\text{eq}})x^*).$$

Since the right-hand-side is independent of  $i$ , this implies both that the  $\mathbf{x}^{\text{eq}}(\xi)$  satisfying (27) must satisfy  $\mathbf{x}^{\text{eq}} = \mathbf{1}_N \otimes x^{\text{eq}}$  and that its expression is given by (28) (indeed (28) can be retrieved immediately from the equivalence above since  $-x^* = (\boldsymbol{\xi}^z + \bar{h}'(\mathbf{x}^{\text{eq}}))^{-1} (-\boldsymbol{\xi}^z x^* - \bar{h}'(\mathbf{x}^{\text{eq}})x^*)$ ).

We now prove (27) by exploiting the Implicit Function Theorem [63]. If we indeed substitute the necessary condition  $\mathbf{x}^{\text{eq}} = \mathbf{1}_N \otimes x^{\text{eq}}$  into the definition of  $\phi_x(\mathbf{x}^{\text{eq}}, \boldsymbol{\xi})$ , we obtain the parallelization of  $N$  equivalent equations of the form

$$x^{\text{eq}} + x^* = (\bar{h}'(\mathbf{x}^{\text{eq}}) + \boldsymbol{\xi}^z)^{-1} (\bar{g}'(\mathbf{x}^{\text{eq}}) + \boldsymbol{\xi}^y + \bar{h}'(\mathbf{x}^{\text{eq}})x^*)$$

where we used properties (16a) and (16b) that lead to  $\bar{h}'(\mathbf{1}_N \otimes x) = \bar{h}'(x)$  and  $\bar{g}'(\mathbf{1}_N \otimes x) = \bar{g}'(x)$ .

Moreover, Assumption 5 ensures that  $\bar{h}'(x^*) \geq cI$ . Thus, for the continuity assumptions in Assumption 1, there exists a sufficiently small  $r > 0$  s.t. if  $\|\xi^z\| \leq \|\xi\| \leq r$  then  $\bar{h}'(x^*) + \xi^z$  is still invertible. Therefore

$$\bar{g}'(\mathbf{x}^{\text{eq}}) + \boldsymbol{\xi}^y + \bar{h}'(\mathbf{x}^{\text{eq}})x^* = \bar{h}'(\mathbf{x}^{\text{eq}})(x^{\text{eq}} + x^*) + \xi^z(x^{\text{eq}} + x^*).$$

Exploiting now the equivalence  $\bar{g}'(x^{\text{eq}}) = \bar{h}'(x^{\text{eq}})x^{\text{eq}} - \nabla \bar{f}'(x^{\text{eq}})$ , it follows that  $x^{\text{eq}}$  must satisfy the following condition:

$$\nabla \bar{f}'(x^{\text{eq}}) - \xi^y + \xi^z(x^{\text{eq}} + x^*) = 0.$$

Given Assumption 1, the left-hand side of the previous inequality is a continuously differentiable function, since

$$\frac{\partial (\nabla \bar{f}'(x^{\text{eq}}) - \xi^y + \xi^z(x^{\text{eq}} + x^*))}{\partial x^{\text{eq}}} = \nabla^2 \bar{f}'(x^{\text{eq}}) + \xi^z.$$

Notice moreover that if  $r$  is sufficiently small (i.e.,  $\|\xi^z\|$  is sufficiently small) then the differentiation is an invertible matrix, since once again  $\nabla^2 \bar{f}'(x^*) \geq cI$  by assumption. Therefore, by the Implicit Function Theorem,  $x^{\text{eq}}(\xi)$  exists, is unique and continuously differentiable. ■

*Proof (of Theorem 10): Proof of a):*  $V_{\text{PNR}}(\mathbf{0}) = 0$  and  $V_{\text{PNR}}(\mathbf{x}) > 0$  for  $\mathbf{x} \neq \mathbf{0}$  have been proved before.  $\dot{V}_{\text{PNR}} < 0$  is instead proved by proving (31a).

*Proof of b):* as for (31a), consider that,  $\forall \mathbf{x} \in \mathbb{R}^{nN}$

$$\begin{aligned} \frac{\partial V_{\text{PNR}}}{\partial \mathbf{x}} \phi'_x(\mathbf{x}, \xi) &= \frac{\partial V_{\text{PNR}}}{\partial \mathbf{x}} \phi'_x(\mathbf{x}, 0) + \frac{\partial V_{\text{PNR}}}{\partial \mathbf{x}} \\ &\quad \times (\phi'_x(\mathbf{x}, \xi) - \phi'_x(\mathbf{x}, 0)) \\ &\leq \frac{\partial V_{\text{PNR}}}{\partial \mathbf{x}} \phi_{\text{PNR}}(\mathbf{x}) + \left\| \frac{\partial V_{\text{PNR}}}{\partial \mathbf{x}} \right\| \\ &\quad \times \|\phi'_x(\mathbf{x}, \xi) - \phi'_x(\mathbf{x}, 0)\| \\ &\leq -b_7 \|\mathbf{x}\|^2 + b_6 \|\mathbf{x}\| a_\xi \|\xi\| \|\mathbf{x}\| \\ &\leq -(b_7 - b_6 a_\xi r) \|\mathbf{x}\|^2 \leq -b'_7 \|\mathbf{x}\|^2. \end{aligned}$$

Notice that this inequality is meaningful for  $r < (b_7/b_6 a_\xi)$ .

As for (31b), consider that,  $\forall \mathbf{x} \in \mathbb{R}^{nN}$

$$\begin{aligned} \|\phi'_x(\mathbf{x}, \xi)\| &\leq \|\phi'_x(\mathbf{x}, 0)\| + \|\phi'_x(\mathbf{x}, \xi) - \phi'_x(\mathbf{x}, 0)\| \\ &\leq (b_8 + a_\xi r) \|\mathbf{x}\| \leq b'_8 \|\mathbf{x}\|. \end{aligned}$$

■

*Proof (of Theorem 11):* The minimizer of the global cost function is easily seen to be  $x^* = (\sum_i A_i)^{-1} (\sum_i A_i d_i)$  from which it follows that  $\bar{f}'(x) = (1/N)x^T A x$ . Clearly  $\bar{f}(x)$  satisfies Assumption 1 since  $\nabla^2 \bar{f}(x) = (1/N)A > 0$  is independent of  $x$ . Considering then  $h'_i(x) = \nabla^2 f'_i(x) = A_i$  it follows after some suitable simplifications that:

$$\begin{aligned} \bar{h}'(x) &= \frac{1}{N} A \\ g'_i(x) &= A_i x - A_i(x + x^* - d_i) = A_i(d_i - x^*) \\ g'(x) - g'(x') &= 0 \\ \bar{g}'(x) &= \frac{1}{N} \left( \sum_i A_i d_i - \sum_i A_i x^* \right) = 0 \\ h'(x) - h'(x') &= 0 \\ \psi(x) &= \bar{h}^{-1}(x) \bar{g}(x) = 0 \\ x^{\text{eq}}(\xi) &= \left( \frac{1}{N} A + \xi^z \right)^{-1} (\xi^y - \xi^z x^*) \\ \phi'_x(x, \xi) &= \phi'_x(x, 0) = -x \end{aligned}$$

where in the last equivalence we exploited definition (28). Thus also the other assumptions are satisfied. ■

*Proof (of Theorem 12):* The proof considers the system as an autonomous singularly perturbed system, and proceeds as follows: a) show that  $x^*$  is an equilibrium; b) perform a change of variables; c) construct a Lyapunov function for the boundary layer system; d) construct a Lyapunov function for the reduced system; e) join the two Lyapunov functions into one, and show (by cascading the previously introduced Lemmas and Theorems) that the complete system (43) converges to  $x^*$  while satisfying the hypotheses of Theorem 2. By doing so it follows that (42), i.e., Algorithm 1, is exponentially stable.

For notational simplicity we let  $x^* := \mathbf{1}_N \otimes x^*$ . We also use all the notation collected in Section II.

- *Discrete to continuous dynamics)* The dynamics of Algorithm 1 can be written in state space as

$$\begin{cases} v(k) = g(x(k-1)) \\ w(k) = h(x(k-1)) \\ y(k) = P[y(k-1) + g(x(k-1)) - v(k-1)] \\ z(k) = P[z(k-1) + h(x(k-1)) - w(k-1)] \\ x(k) = (1-\varepsilon)x(k-1) + \varepsilon \frac{y(k-1)}{[z(k-1)]_c} \end{cases} \quad (42)$$

with suitable initial conditions. (42) can then be interpreted as the forward-Euler discretization of

$$\begin{cases} \varepsilon \dot{v}(t) = -v(t) + g(x(t)) \\ \varepsilon \dot{w}(t) = -w(t) + h(x(t)) \\ \varepsilon \dot{y}(t) = -K y(t) + (I - K)[g(x(t)) - v(t)] \\ \varepsilon \dot{z}(t) = -K z(t) + (I - K)[h(x(t)) - w(t)] \\ \dot{x}(t) = -x(t) + \frac{y(t)}{[z(t)]_c} \end{cases} \quad (43)$$

with null initial conditions, where  $\varepsilon$  is the discretization time interval and  $K := I - P$ . Notice that, as for  $P$ , if  $n$  is the dimension of the local costs then  $P = P' \otimes I_n$  with  $P'$  a doubly-stochastic average consensus matrix. Nonetheless for brevity we will omit the superscripts.

- b) let

$$\begin{aligned} v' &:= v - g(x) \\ w' &:= w - h(x) \\ y' &:= y - v' - \mathbf{1}_N \otimes \bar{g}(x) \\ z' &:= z - w' - \mathbf{1}_N \otimes \bar{h}(x) \\ x' &:= x - x^* \end{aligned}$$

$$\begin{aligned} \phi_g(x') &:= \frac{\partial g}{\partial x'} - \mathbf{1}_N \otimes \frac{\partial \bar{g}}{\partial x'} \\ \phi_h(x') &:= \frac{\partial h}{\partial x'} - \mathbf{1}_N \otimes \frac{\partial \bar{h}}{\partial x'} \\ \phi_x(x', \chi) &:= -x'(t) - x^* + \\ &\quad + \frac{y'(t) + v'(t) + \mathbf{1}_N \otimes \bar{g}(x'(t) + x^*)}{[z'(t) + w'(t) + \mathbf{1}_N \otimes \bar{h}(x'(t) + x^*)]_c} \end{aligned}$$

with  $\chi := (v', w', y', z')$ , so that (43) becomes

$$\begin{cases} \varepsilon \dot{v}'(t) = -v'(t) - \varepsilon \frac{\partial g}{\partial x'} \dot{x}'(t) \\ \varepsilon \dot{w}'(t) = -w'(t) - \varepsilon \frac{\partial h}{\partial x'} \dot{x}'(t) \\ \varepsilon \dot{y}'(t) = -K y'(t) + \varepsilon \phi_g(x') \dot{x}'(t) \\ \varepsilon \dot{z}'(t) = -K z'(t) + \varepsilon \phi_h(x') \dot{x}'(t) \\ \dot{x}'(t) = \phi_x(x', \chi') \end{cases} \quad (44)$$



with initial conditions

$$\begin{cases} \mathbf{v}'(0) = \mathbf{v}(0) - g(\mathbf{x}(0)) \\ \mathbf{w}'(0) = \mathbf{w}(0) - h(\mathbf{x}(0)) \\ \mathbf{y}'(0) = \mathbf{y}(0) - \mathbf{v}(0) + g^\perp(\mathbf{x}(0)) \\ \mathbf{z}'(0) = \mathbf{z}(0) - \mathbf{w}(0) + h^\perp(\mathbf{x}(0)) \\ \mathbf{x}'(0) = \mathbf{x}(0) - \mathbf{x}^* \end{cases}$$

where  $g^\perp(\mathbf{x}) := g(\mathbf{x}) - \mathbb{1}_N \otimes \bar{g}(\mathbf{x})$  (equivalent definition for  $h^\perp$ ). Notice that (44) has the origin as an equilibrium point. Moreover this dynamics exploits the function  $\phi_x$  defined in (24), with  $\chi^y = \mathbf{y}' + \mathbf{v}'$ , and  $\chi^z = \mathbf{z}' + \mathbf{w}'$ .

The next step is to exploit the structure of  $K$  (more precisely, the fact that it contains an average consensus matrix) to reduce the dynamics, i.e., to eliminate the dynamics of the average since the latter does not change in time. To this aim, we analyze the behavior of the average of the  $y'_i$ s, i.e., the behavior of  $(\mathbb{1}_N^T \otimes I_n) \dot{\mathbf{y}}'$ . To this point, consider the third equation in (44). Recalling that  $(A \otimes B)(C \otimes D) = AB \otimes CD$ , and exploiting the fact that  $\mathbb{1}_N^T P' = 0$ , we notice that  $(\mathbb{1}_N^T \otimes I_n) K = 0$ . Moreover, from the definitions of  $g$  and  $\bar{g}$

$$(\mathbb{1}_N^T \otimes I_n) \frac{\partial g(\mathbf{x}')}{\partial \mathbf{x}'} = N \frac{\partial \bar{g}(\mathbf{x}')}{\partial \mathbf{x}'}.$$

Since  $N = \mathbb{1}_N^T \mathbb{1}_N$ , it follows also that:

$$(\mathbb{1}_N^T \otimes I_n) \phi_g(\mathbf{x}') = 0$$

for all  $t \geq 0$ , i.e.,  $\mathbb{1}^T \mathbf{y}'(t) = \mathbb{1}^T \mathbf{y}'(0) \equiv 0$ . Similarly it is possible to show that  $\mathbf{z}'(t) \equiv 0$ . This eventually implies that

$$\mathbf{y}^{\parallel}(t) = 0 \quad \mathbf{z}^{\parallel}(t) = 0 \quad \forall t$$

that means, recalling that  $\mathbf{y}' = \mathbf{y}^{\parallel} + \mathbf{y}^{\perp}$  and  $\mathbf{z}' = \mathbf{z}^{\parallel} + \mathbf{z}^{\perp}$ , that (44) can be equivalently rewritten as

$$\begin{cases} \varepsilon \dot{\mathbf{v}}'(t) = -\mathbf{v}'(t) - \varepsilon \frac{\partial g}{\partial \mathbf{x}'} \phi_x(\mathbf{x}', \chi') \\ \varepsilon \dot{\mathbf{w}}'(t) = -\mathbf{w}'(t) - \varepsilon \frac{\partial h}{\partial \mathbf{x}'} \phi_x(\mathbf{x}', \chi') \\ \varepsilon \dot{\mathbf{y}}^{\perp}(t) = -K \mathbf{y}^{\perp}(t) + \varepsilon \phi_g(\mathbf{x}') \phi_x(\mathbf{x}', \chi') \\ \varepsilon \dot{\mathbf{z}}^{\perp}(t) = -K \mathbf{z}^{\perp}(t) + \varepsilon \phi_h(\mathbf{x}') \phi_x(\mathbf{x}', \chi') \\ \dot{\mathbf{x}}'(t) = \phi_x(\mathbf{x}', \chi') \end{cases} \quad (45)$$

where now  $\chi' := (\mathbf{v}, \mathbf{w}, \mathbf{y}^{\perp}, \mathbf{z}^{\perp})$  and where the novel initial conditions for the changed variables are

$$\begin{cases} \mathbf{y}^{\perp}(0) = \mathbf{y}^{\perp}(0) - \mathbf{v}^{\perp}(0) + g^\perp(\mathbf{x}(0)) \\ \mathbf{z}^{\perp}(0) = \mathbf{z}^{\perp}(0) - \mathbf{w}^{\perp}(0) + h^\perp(\mathbf{x}(0)) \end{cases}$$

• c) the boundary layer of (45) is computed by setting  $\mathbf{x}'(t) = \mathbf{x}'$ . Since a constant  $\mathbf{x}'$  implies  $\dot{\mathbf{x}}' = \phi_x = 0$ , this boundary layer reduces to a linear system globally exponentially converging to the origin. Notice that this implies that, in the original coordinates system

$$\mathbf{v} = g(\mathbf{x}), \quad \mathbf{w} = h(\mathbf{x}), \quad \mathbf{y} = \mathbb{1}_N \otimes \bar{g}(\mathbf{x}), \quad \mathbf{z} = \mathbb{1}_N \otimes \bar{h}(\mathbf{x}).$$

In the novel coordinates system we thus consider, as a Lyapunov function,  $(1/2)\|\chi'\|^2$ .

• d) the reduced system of (45) is computed by plugging  $\chi' = \mathbf{0}$  into the equations (i.e., by setting  $\mathbf{v}'(t) = \mathbf{0}$ ,  $\mathbf{w}'(t) = \mathbf{0}$ ,  $\mathbf{y}^{\perp}(t) = \mathbf{0}$ ,  $\mathbf{z}^{\perp}(t) = \mathbf{0}$ ). Defining then

$$f'_i(\mathbf{x}') := f_i(\mathbf{x}' + \mathbf{x}^*), \quad h'_i(\mathbf{x}') := h_i(\mathbf{x}' + \mathbf{x}^*)$$

we obtain

$$\begin{aligned} \dot{\mathbf{x}}'(t) &= -\mathbf{x}'(t) - \mathbf{x}^* + \mathbb{1}_N \otimes \frac{\bar{g}'(\mathbf{x}'(t))}{\bar{h}'(\mathbf{x}'(t))} \\ &= -\mathbf{x}'(t) - \mathbf{x}^* + \mathbb{1}_N \\ &\quad \otimes \frac{\bar{h}'(\mathbf{x}'(t))(\mathbf{x}'(t) + \mathbf{x}^*) - \nabla f'(\mathbf{x}'(t))}{\bar{h}'(\mathbf{x}'(t))} \\ &= -\mathbf{x}'(t) + \mathbb{1}_N \otimes \frac{\bar{h}'(\mathbf{x}'(t))\mathbf{x}'(t) - \nabla f'(\mathbf{x}'(t))}{\bar{h}'(\mathbf{x}'(t))} \\ &= -\mathbf{x}'(t) + \mathbb{1}_N \otimes \psi(\mathbf{x}'(t)) \\ &= \phi_{\text{PNR}}(\mathbf{x}') \end{aligned}$$

where  $\psi$  and  $\phi_{\text{PNR}}$  are the functions defined in (15) and (17), respectively. Thus the reduced system, thanks to Theorem 6, admits  $\mathbf{x}^*$  as a global exponentially stable equilibrium, and admits  $V_{\text{PNR}}$  in (21) as a Lyapunov function.

• e) we now notice that the interconnection of the boundary layer and reduced systems maintains the global stability, since their Lyapunov functions are quadratic type. Thus (see [46, pp. 453]) the global system is asymptotically globally stable. To check that forward-Euler discretizations of the system preserve these stability properties we then consider as a global Lyapunov function the function

$$V(\mathbf{x}', \chi') = (1-d)V_{\text{PNR}}(\mathbf{x}') + d \frac{1}{2} \|\chi'\|^2$$

that is clearly positive definite for every  $d \in (0, 1)$ , and prove that inequalities (5) of Theorem 2 are satisfied.

*Proof That (5a) Holds:* from (22a) and the structure of  $V$  it follows immediately that:

$$((1-d)b_5 + d)I \leq \nabla^2 V(\mathbf{x}', \chi') \leq ((1-d)b_6 + d)I.$$

*Proof That (5c) Holds:* applying (20) and (26a) to (45) it follows that (5c) holds with

$$a_4 = a_V := \max\{1 + 2\varepsilon a_g a_x, 1 + 2\varepsilon a_h a_x, a_x\}.$$

*Proof That (5b) Holds:* the part relative to the slow dynamics is already characterized by (31a). For the part relative to the fast dynamics, since  $(\partial(1/2)\|\chi'\|^2)/\partial \chi = \chi'^T$  to check that (5b) corresponds to check the negativity of the terms

$$\begin{aligned} & -\mathbf{v}'^T \mathbf{v}' - \varepsilon \mathbf{v}'^T \frac{\partial g}{\partial \mathbf{x}'} \phi_x(\mathbf{x}', \chi') \\ & -\mathbf{w}'^T \mathbf{w}' - \varepsilon \mathbf{w}'^T \frac{\partial h}{\partial \mathbf{x}'} \phi_x(\mathbf{x}', \chi') \\ & -(\mathbf{y}^{\perp})^T K \mathbf{y}' \perp + \varepsilon \mathbf{y}^{\perp T} \phi_g(\mathbf{x}') \phi_x(\mathbf{x}', \chi') \\ & -(\mathbf{z}^{\perp})^T K \mathbf{z}' \perp + \varepsilon \mathbf{z}^{\perp T} \phi_h(\mathbf{x}') \phi_x(\mathbf{x}', \chi'). \end{aligned}$$

These terms can then be majorized using (20) and (26a). E.g., the third term can be majorized with

$$-\sigma(P)\|\mathbf{y}^{\perp}\|^2 + 2\varepsilon a_g a_x \|\mathbf{y}^{\perp}\|(\|\mathbf{x}\| + \|\chi\|)$$

where  $\sigma(P)$  is the spectral gap of  $P$ . Applying similar concepts also to the other terms it follows that (5b) holds with

$$a_3 = \min\{\sigma(P) - 2\varepsilon a_g a_x, \sigma(P) - 2\varepsilon a_h a_x\}.$$

1089 *Proof (of Theorem 13):* The proof is identical to the one  
 1090 of Theorem 12 with the exception that the substitution is  
 1091 now  $\mathbf{x}'' = \mathbf{x} - \mathbf{x}^* - \mathbb{1}_N \otimes \Psi(\xi^y, \xi^z)$ . Indeed one can prove  
 1092 the stability of the novel system using the same Lyapunov  
 1093 function of Theorem 12. Notice that we are ensured that there  
 1094 exists a sufficiently small neighborhood of the origin for which  
 1095 the function  $\Psi$  exists due to the smoothness conditions in  
 1096 Assumption 1. ■

1097 *Proof (of Theorem 14):* The proof is the local version of the  
 1098 one in Theorem 13. Indeed the local versions of Assumptions 1,  
 1099 5, and 9 always hold, i.e., they hold when considering  $\mathbf{x}$  s.t.  
 1100  $\|\mathbf{x}\| \leq r'$ , and one can thus repeat that reasonings using local  
 1101 perspectives. ■

1102 *Proof (of Theorem 15):* Consider for simplicity the scalar  
 1103 case. Let  $y^* := (1/N) \sum_i A_i d_i$  and  $z^* := (1/N) \sum_i A_i$ , so  
 1104 that  $x^* = y^*/z^*$ . Since  $\mathbf{y}(k+1) = P\mathbf{y}(k)$  and  $\mathbf{z}(k+1) =$   
 1105  $P\mathbf{z}(k)$ , given the assumptions on  $P$ , there exist positive  $\alpha_y$ ,  
 1106  $\alpha_z$  independent of  $\mathbf{x}(0)$  s.t.  $|y_i(k) - y^*| \leq \alpha_y(\rho(P))^k$  and  
 1107  $|z_i(k) - z^*| \leq \alpha_z(\rho(P))^k$ . The claim thus follows considering  
 1108 that  $x_i(k) = y_i(k)/[z_i(k)]_c$  and that, since the elements of  $P$   
 1109 are non negative, all the  $z_i(k)$  are non smaller than  $c$  for all  
 1110  $k \geq 0$  (i.e., the operator  $[\cdot]_c$  is always performing as the identity  
 1111 operator). ■

## 1112 REFERENCES

1113 [1] F. Zanella, D. Varagnolo, A. Cenedese, G. Pillonetto, and L. Schenato,  
 1114 "Newton-Raphson consensus for distributed convex optimization," in  
 1115 *Proc. IEEE Conf. Decision Control & Eur. Control Conf.*, Dec. 2011,  
 1116 pp. 5917–5922.  
 1117 [2] F. Zanella, D. Varagnolo, A. Cenedese, G. Pillonetto, and L. Schenato,  
 1118 "Multidimensional Newton-Raphson consensus for distributed convex op-  
 1119 timization," in *Proc. Amer. Control Conf.*, 2012, pp. 1079–1084.  
 1120 [3] N. Z. Shor, *Minimization Methods for Non-Differentiable Functions*.  
 1121 New York: Springer-Verlag, 1985.  
 1122 [4] D. P. Bertsekas, A. Nedić, and A. E. Ozdaglar, *Convex Analysis and*  
 1123 *Optimization*. Belmont, MA: Athena Scientific, 2003.  
 1124 [5] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.:  
 1125 Cambridge Univ. Press, 2004.  
 1126 [6] J. N. Tsitsiklis, "Problems in Decentralized Decision Making and Com-  
 1127 putation," Ph.D. dissertation, MIT, Cambridge, MA, 1984.  
 1128 [7] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation:*  
 1129 *Numerical Methods*. Belmont, MA: Athena Scientific, 1997.  
 1130 [8] D. P. Bertsekas, *Network Optimization: Continuous and Discrete Models*.  
 1131 Belmont, MA, USA: Athena Scientific, 1998.  
 1132 [9] M. Bürger, G. Notarstefano, F. Bullo, and F. Allgöwer, "A distributed  
 1133 simplex algorithm for degenerate linear programs and multi-agent assign-  
 1134 ments," *Automatica*, vol. 48, no. 9, pp. 2298–2304, 2012.  
 1135 [10] D. P. Bertsekas, *Constrained Optimization and Lagrange Multiplier Meth-*  
 1136 *ods*. Boston, MA: Academic Press, 1982.  
 1137 [11] M. R. Hestenes, "Multiplier and gradient methods," *J. Optim. Theory*  
 1138 *Appl.*, vol. 4, no. 5, pp. 303–320, 1969.  
 1139 [12] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed Opti-  
 1140 mization and Statistical Learning via the Alternating Direction Method of  
 1141 Multipliers," Statistics Dept., Stanford Univ., Stanford, CA, Tech. Rep.,  
 1142 2010.  
 1143 [13] T. Erseghe, D. Zennaro, E. Dall'Anese, and L. Vangelista, "Fast consen-  
 1144 sus by the alternating direction multipliers method," *IEEE Trans. Signal*  
 1145 *Processing*, vol. 59, no. 11, pp. 5523–5537, Nov. 2011.  
 1146 [14] B. He and X. Yuan, "On the  $O(1/t)$  convergence rate of alternating di-  
 1147 rection method," *SIAM J. Numer. Anal.*, 2011. [Online]. Available: [http://](http://www.optimization-online.org/DB_FILE/2011/09/3157.pdf)  
 1148 [www.optimization-online.org/DB\\_FILE/2011/09/3157.pdf](http://www.optimization-online.org/DB_FILE/2011/09/3157.pdf)  
 1149 [15] W. Deng and W. Yin, "On the Global and Linear Convergence of the Gen-  
 1150 eralized Alternating Direction Method of Multipliers," DTIC Document,  
 1151 Tech. Rep., 2012.  
 1152 [16] E. Wei and A. Ozdaglar, "Distributed alternating direction method of  
 1153 multipliers," in *Proc. IEEE Conf. Decision Control*, 2012, pp. 5445–5450.  
 1154 [17] J. a. Mota, J. Xavier, P. Aguiar, and M. Püschel, "Distributed ADMM for  
 1155 model predictive control and congestion control," in *Proc. IEEE Conf.*  
 1156 *Decision Control*, 2012, pp. 5110–5115.

[18] D. Jakovetić, J. a. Xavier, and J. M. F. Moura, "Cooperative convex op- 1157  
 timization in networked systems: Augmented lagrangian algorithms with 1158  
 directed gossip communication," *IEEE Trans. Signal Processing*, vol. 59, 1159  
 no. 8, pp. 3889–3902, Aug. 2011. 1160  
 [19] V. F. Dem'yanov and L. V. Vasil'ev, *Nondifferentiable Optimization*. 1161  
 New York: Springer-Verlag, 1985. 1162  
 [20] B. Johansson, "On Distributed Optimization in Networked Systems," 1163  
 Ph.D. dissertation, KTH Royal Institute of Technology, Stockholm, 1164  
 Sweden, 2008. 1165  
 [21] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip al- 1166  
 gorithms," *IEEE Trans. Inform. Theory/ACM Trans. Netw.*, vol. 52, no. 6, 1167  
 pp. 2508–2530, Jun. 2006. 1168  
 [22] A. Ribeiro, "Ergodic stochastic optimization algorithms for wireless com- 1169  
 munication and networking," *IEEE Trans. Signal Processing*, vol. 58, 1170  
 no. 12, pp. 6369–6386, Dec. 2010. 1171  
 [23] A. Nedić and D. P. Bertsekas, "Incremental subgradient methods for 1172  
 nondifferentiable optimization," *SIAM J. Optim.*, vol. 12, no. 1, pp. 109– 1173  
 138, 2001. 1174  
 [24] A. Nedić, D. Bertsekas, and V. Borkar, "Distributed asynchronous incre- 1175  
 mental subgradient methods," *Studies Computat. Math.*, vol. 8, pp. 381– 1176  
 407, 2001. 1177  
 [25] A. Nedić and A. Ozdaglar, "Approximate primal solutions and rate anal- 1178  
 ysis for dual subgradient methods," *SIAM J. Optim.*, vol. 19, no. 4, 1179  
 pp. 1757–1780, 2008. 1180  
 [26] K. C. Kiwiel, "Convergence of approximate and incremental subgradi- 1181  
 ent methods for convex optimization," *SIAM J. Optim.*, vol. 14, no. 3, 1182  
 pp. 807–840, 2004. 1183  
 [27] D. Blatt, A. Hero, and H. Gauchman, "A convergent incremental gradient 1184  
 method with a constant step size," *SIAM J. Optim.*, vol. 18, no. 1, pp. 29– 1185  
 51, 2007. 1186  
 [28] L. Xiao and S. Boyd, "Optimal scaling of a gradient method for distributed 1187  
 resource allocation," *J. Optim. Theory Appl.*, vol. 129, no. 3, pp. 469–488, 1188  
 2006. 1189  
 [29] S. S. Ram, A. Nedić, and V. V. Veeravalli, "Incremental stochastic subgra- 1190  
 dient algorithms for convex optimization," *SIAM J. Optim.*, vol. 20, no. 2, 1191  
 pp. 691–717, 2009. 1192  
 [30] A. Nedić and A. Ozdaglar, "Distributed subgradient methods for multi- 1193  
 agent optimization," *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48– 1194  
 61, 2009. 1195  
 [31] B. Johansson, M. Rabi, and M. Johansson, "A randomized incremental 1196  
 subgradient method for distributed optimization in networked systems," 1197  
*SIAM J. Optim.*, vol. 20, no. 3, pp. 1157–1170, 2009. 1198  
 [32] I. Lobel, A. Ozdaglar, and D. Feijer, "Distributed multi-agent optimiza- 1199  
 tion with state-dependent communication," *Math. Programming*, vol. 129, 1200  
 no. 2, pp. 255–284, 2011. 1201  
 [33] A. Nedić, "Asynchronous broadcast-based convex optimization over a 1202  
 network," *IEEE Trans. Autom. Control*, vol. 56, no. 6, pp. 1337–1351, 1203  
 Jun. 2010. 1204  
 [34] E. Ghadimi, I. Shames, and M. Johansson, "Accelerated gradient meth- 1205  
 ods for networked optimization," in *Proc. Amer. Control Conf.*, 2011, 1206  
 pp. 1668–1673. 1207  
 [35] A. Jadbabaie, A. Ozdaglar, and M. Zargham, "A distributed newton 1208  
 method for network optimization," in *Proc. IEEE Conf. Decision Control*, 1209  
 2009, pp. 2736–2741. 1210  
 [36] M. Zargham, A. Ribeiro, A. Ozdaglar, and A. Jadbabaie, "Accelerated 1211  
 dual descent for network optimization," *IEEE Trans. Autom. Control*, 1212  
 vol. 59, no. 4, pp. 905–920, 2013. 1213  
 [37] E. Wei, A. Ozdaglar, and A. Jadbabaie, "A distributed newton method 1214  
 for network utility maximization," in *Proc. IEEE Conf. Decision Control*, 1215  
 2010, pp. 1816–1821. 1216  
 [38] R. S. Dembo, S. C. Eisenstat, and T. Steihaug, "Inexact newton methods," 1217  
*SIAM J. Numer.*, vol. 19, no. 2, pp. 400–408, 1982. 1218  
 [39] A. Nedić, A. Ozdaglar, and P. A. Parrilo, "Constrained consensus and op- 1219  
 timization in multi-agent networks," *IEEE Trans. Autom. Control*, vol. 55, 1220  
 no. 4, pp. 922–938, Apr. 2010. 1221  
 [40] M. Zhu and S. Martínez, "On distributed convex optimization under in- 1222  
 equality and equality constraints," *IEEE Trans. Autom. Control*, vol. 57, 1223  
 no. 1, pp. 151–164, 2012. 1224  
 [41] C. Fischione, "F-lipschitz optimization with wireless sensor networks 1225  
 applications," *IEEE Trans. Autom. Control*, vol. 56, no. 10, pp. 2319– 1226  
 2331, 2011. 1227  
 [42] J. Wang and N. Elia, "Control approach to distributed optimization," in 1228  
*Proc. 48th Annu. Allerton Conf.*, Allerton, IL, USA, Sep. 2010, vol. 1, 1229  
 no. 1, pp. 557–561, IEEE. 1230  
 [43] N. Freris and A. Zouzias, "Fast distributed smoothing for network 1231  
 clock synchronization," in *Proc. IEEE Conf. Decision Control*, 2012, 1232  
 pp. 1411–1416. 1233

- [44] I. Necoara and V. Nedelcu, "Distributed dual gradient methods and error bound conditions," 2014. [Online]. Available: <http://arxiv.org/abs/1401.4398>
- [45] F. Garin and L. Schenato, *A Survey on Distributed Estimation and Control Applications Using Linear Consensus Algorithms*, vol. 406. New York: Springer, 2011, pp. 75–107.
- [46] H. K. Khalil, *Nonlinear Systems*, 3rd ed. Englewood Cliffs, NJ: Prentice Hall, 2001.
- [47] A. Nedic and A. Olshevsky, "Distributed optimization over time-varying directed graphs," in *Proc. IEEE 52nd Annu. Conf. Decision Control (CDC)*, 2013, pp. 6855–6860.
- [48] F. Fagnani and S. Zampieri, "Randomized consensus algorithms over large scale networks," *IEEE J. Selected Areas Commun.*, vol. 26, no. 4, pp. 634–649, May 2008.
- [49] A. D. Domínguez-García, C. N. Hadjicostis, and N. H. Vaidya, "Distributed Algorithms for Consensus and Coordination in the Presence of Packet-Dropping Communication Links Part I: Statistical Moments Analysis Approach," Coordinated Sciences Laboratory, University of Illinois at Urbana-Champaign, Tech. Rep., 2011.
- [50] P. Kokotović, H. K. Khalil, and J. O'Reilly, *Singular Perturbation Methods in Control: Analysis and Design*. Philadelphia, PA: SIAM, 1999, ser. Classics in applied mathematics.
- [51] K. Tanabe, "Global analysis of continuous analogues of the Levenberg-Marquardt and Newton-Raphson methods for solving nonlinear equations," *Annals Inst. Stat. Math.*, vol. 37, no. 1, pp. 189–203, 1985.
- [52] R. Hauser and J. Nedić, "The continuous newton-raphson method can look ahead," *SIAM J. Optim.*, vol. 15, pp. 915–925, 2005.
- [53] T. Sahai, A. Speranzon, and A. Banaszuk, "Hearing the clusters of a graph: A distributed algorithm," *Automatica*, vol. 48, no. 1, pp. 15–24, Jan. 2012.
- [54] S. Becker and Y. Le Cun, "Improving the Convergence of Back-Propagation Learning With Second Order Methods," Univ. Toronto, Tech. Rep., Sep. 1988.
- [55] S. Athuraliya and S. H. Low, "Optimization flow control with newton like algorithm," *Telecommun. Syst.*, vol. 15, pp. 345–358, 2000.
- [56] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 3rd ed. Baltimore, MD: John Hopkins Univ. Press, 1996.
- [57] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, Prediction*, 2nd ed. New York: Springer, 2001.
- [58] L. Xiao, S. Boyd, and S.-J. Kim, "Distributed average consensus with least-mean-square deviation," *J. Parallel Distrib. Comp.*, vol. 67, no. 1, pp. 33–46, Jan. 2007.
- [59] S. Muthukrishnan, B. Ghosh, and M. H. Schultz, "First and second order diffusive methods for rapid, coarse, distributed load balancing," *Theory Comp. Syst.*, vol. 31, no. 4, pp. 331–354, 1998.
- [60] I. D. Schizas, A. Ribeiro, and G. B. Giannakis, "Consensus in Ad Hoc WSNs with noisy links—part I: Distributed estimation of deterministic signals," *IEEE Trans. Signal Processing*, vol. 56, pp. 350–364, Jan. 2008.
- [61] E. Ghadimi, A. Teixeira, I. Shames, and M. Johansson, "On the optimal step-size selection for the alternating direction method of multipliers," *Necsys*, 2012. [Online]. Available: [https://people.kth.se/~andretei/papers/gtsj\\_NECSYS12.pdf](https://people.kth.se/~andretei/papers/gtsj_NECSYS12.pdf)
- [62] F. Zanella, D. Varagnolo, A. Cenedese, G. Pillonetto, and L. Schenato, "Asynchronous newton-raphson consensus for distributed convex optimization," *Necsys*, 2012. [Online]. Available: [http://automatica.dei.unipd.it/tl\\_files/utenti/lucaschenato/Papers/Conference/Necsys12\\_Asynchrous\\_Newton\\_Raphson\\_Consensus.pdf](http://automatica.dei.unipd.it/tl_files/utenti/lucaschenato/Papers/Conference/Necsys12_Asynchrous_Newton_Raphson_Consensus.pdf)
- [63] L. Kudryavtsev, *Encyclopedia of Mathematics*. New York: Springer, 2001.



**Damiano Varagnolo** (M'11) received the M.S. degree in automation engineering and the Ph.D. degree in information engineering from the University of Padova, Italy, in 2005 and 2011, respectively. He was a Research Engineer at Tecnogamma S.p.A., Treviso, Italy, from 2006 to 2007 and visited UC Berkeley as a Scholar Researcher in 2010. From March 2012 to December 2013 he worked as a Post-Doctoral Scholar at the KTH, Royal Institute of Technology, Stockholm, Sweden. Currently, he is Associate Senior Lecturer at LTU, Luleå University of Technology. His interests include distributed optimization, distributed estimation, identification and control of HVAC systems.



Student Branch of the University of Padova from 2006 to 2008.

**Filippo Zanella** (M'06) was born in Treviso, Italy, in 1983. He received the B.S. and M.S. degrees in automation engineering and the Ph.D. degree in information engineering from the University of Padova, Padova, Italy, in 2005, 2008, and 2013, respectively. He has been a Visiting Student Researcher at UC Berkeley in 2011 and at UC Santa Barbara in 2012. His research interests are in wireless cameras/sensors networks and mobile networks with emphasis on distributed control, estimation and optimization. Dr. Zanella was a Staff Member of the IEEE Student Branch of the University of Padova from 2006 to 2008.



and National government institutions and industries, with different roles participant and/or principal investigator. He coauthored more than 90 papers and holds three patents in the area of sensor/actor networks and videosurveillance. His research interests include system modeling, control theory and its applications, sensor and actuator networks, home automation systems.

**Angelo Cenedese** (M'12) received the M.S. and the Ph.D. degrees from the University of Padova, Padova, Italy, in 1999 and 2004, respectively. He is currently an Assistant Professor with the Department of Information Engineering and member of the Human Inspired Technologies Research Center and the Research Center on Fusion. He has been and he is currently involved in several projects on distributed systems (sensor and actor networks, camera networks), control of complex systems (adaptive optics systems, fusion devices), funded by European and National government institutions and industries, with different roles participant and/or principal investigator. He coauthored more than 90 papers and holds three patents in the area of sensor/actor networks and videosurveillance. His research interests include system modeling, control theory and its applications, sensor and actuator networks, home automation systems.



Control and Dynamic Systems at the Department of Information Engineering, University of Padova where he currently serves as an Associate Professor. He is an Associate Editor of *Automatica* and *Systems and Control Letters*. His research interests are in the field of system identification and machine learning.

**Gianluigi Pillonetto** (M'03) was born in Montebelluna (TV), Italy on January 21, 1975. He received the Doctoral degree in computer science engineering (with highest honors) from the University of Padova, Padova, Italy, in 1998 and the PhD degree in bioengineering from the Polytechnic of Milan, Milan, Italy, in 2002. In 2000 and 2002, he was Visiting Scholar and Visiting Scientist, respectively, at the Applied Physics Laboratory, University of Washington, Seattle. In 2005, he became Assistant Professor of Control and Dynamic Systems at the Department of Information Engineering, University of Padova where he currently serves as an Associate Professor. He is an Associate Editor of *Automatica* and *Systems and Control Letters*. His research interests are in the field of system identification and machine learning.



smart grids and cooperative robotics. Dr. Schenato received the 2004 Researchers Mobility Fellowship from the Italian Ministry of Education, University and Research (MIUR), and the 2006 Eli Jury Award in U.C. Berkeley and the EUCA European Control Award in 2014. He served as Associate Editor for the IEEE TRANSACTIONS ON AUTOMATIC CONTROL from 2010 to 2014.

**Luca Schenato** (SM'15) received the Dr. Eng. degree in electrical engineering from the University of Padova, Padova, Italy, in 1999 and the Ph.D. degree in electrical engineering and computer sciences from UC Berkeley, in 2003. He held a post-doctoral position in 2004 and a Visiting Professor position from 2013 to 2014 at U.C. Berkeley. Currently, he is Associate Professor at the Information Engineering Department, University of Padova. His interests include networked control systems, multi-agent systems, wireless sensor networks, smart grids and cooperative robotics. Dr. Schenato received the 2004 Researchers Mobility Fellowship from the Italian Ministry of Education, University and Research (MIUR), and the 2006 Eli Jury Award in U.C. Berkeley and the EUCA European Control Award in 2014. He served as Associate Editor for the IEEE TRANSACTIONS ON AUTOMATIC CONTROL from 2010 to 2014.