# Distributed Gradient Methods with Variable Number of Working Nodes

Dušan Jakovetić, Dragana Bajović, Nataša Krejić, and Nataša Krklec-Jerinkić

*Abstract*—We consider distributed optimization where $N$ nodes in a connected network minimize the sum of their local costs subject to a common constraint set. We propose a distributed projected gradient method where each node, at each iteration $k$, performs an update (is active) with probability $p_k$, and stays idle (is inactive) with probability $1 - p_k$. Whenever active, each node performs an update by weight-averaging its solution estimate with the estimates of its active neighbors, taking a negative gradient step with respect to its local cost, and performing a projection onto the constraint set; inactive nodes perform no updates. Assuming that nodes' local costs are strongly convex, with Lipschitz continuous gradients, we show that, as long as activation probability $p_k$ grows to one asymptotically, our algorithm converges in the mean square sense (MSS) to the same solution as the standard distributed gradient method, i.e., as if all the nodes were active at all iterations. Moreover, when $p_k$ grows to one linearly, with an appropriately set convergence factor, the algorithm has a linear MSS convergence, with practically the same factor as the standard distributed gradient method. Simulations demonstrate that, when compared with the standard distributed gradient method, the proposed algorithm significantly reduces the overall number of per-node communications and per-node gradient evaluations (computational cost) for the same required accuracy.

*Index Terms*—Distributed optimization, distributed gradient method, variable number of working nodes, convergence rate, consensus.

## I. INTRODUCTION

We consider distributed optimization where $N$ nodes constitute a generic, connected network, each node $i$ has a convex cost function $f_i : \mathbb{R}^d \mapsto \mathbb{R}$ known only by $i$, and the nodes want to solve the following problem:

$$\begin{aligned} \text{minimize} \quad & \textstyle\sum_{i=1}^N f_i(x) =: f(x) \\ \text{subject to} \quad & x \in \mathcal{X}, \end{aligned} \tag{1}$$

where $x \in \mathbb{R}^d$ is the optimization variable common to all nodes, and $\mathcal{X} \subset \mathbb{R}^d$ is a closed, convex constraint set, known by all. The above and related problems arise frequently, e.g., in big data analytics in cluster or cloud environments, e.g., [1]-[3], distributed estimation in wireless sensor networks (WSNs), e.g., [4]-[8], and distributed control applications, e.g., [9], [10]. With all the above applications, data is split across multiple networked nodes (sensors, cluster machines, etc.), and $f_i(x) = f_i(x; D_i)$ represents a loss with respect to data $D_i$ stored locally at node $i$.

A popular approach to solve (1) is via distributed (projected) (sub)gradient methods, e.g., [11], [12], [13]. With these methods, each node $i$, at each iteration $k$, updates its solution estimate by weight-averaging it with the estimates of its neighbors, taking a negative gradient step with respect to its local cost, and projecting the result onto the constraint set $\mathcal{X}$. Distributed gradient methods are attractive as they do not require centralized coordination, have inexpensive iterations (provided that projections onto $\mathcal{X}$ are computationally light), and exhibit resilience to inter-node communication failures and delays; however, they have a drawback of slow convergence rate.

Several techniques to improve convergence rates of distributed (projected) gradient methods have been proposed, including Newton-like methods, e.g., [14], [15], and Nesterov-like methods, e.g., [16], [17]. In this paper, we make distributed (projected) gradient methods more efficient by proposing a novel method with a variable number of working nodes. Each node $i$, at each iteration $k$, performs an update (is active) with probability $p_k$, and stays idle (is inactive) with probability $1 - p_k$. Whenever active, each node $i$ performs the same update as with the standard distributed gradient method, while inactive nodes perform no updates.

Our main results are as follows. Assuming that the costs $f_i$'s are strongly convex and their gradients are Lipschitz continuous, we show that, whenever the activation probability $p_k$ grows asymptotically to one, our method converges in the mean square sense to the same point as the standard distributed gradient method.[1] Moreover, when $p_k$ grows to one linearly, with the convergence factor $\delta \in (0, 1)$, our algorithm has a linear convergence rate (in the sense of the expected distance to the solution). When, in addition, quantity $\delta$ is set in accordance with the $f_i$'s condition number and the underlying network's spectral gap, we show that the proposed algorithm converges practically with the same linear convergence factor as the standard distributed gradient method (albeit with a larger hidden constant). Hence, interestingly, our algorithm achieves practically the same rate in iterations $k$ as the standard distributed gradient method, but with the reduced "work" per iteration $k$ (overall communication and computational cost), thus making distributed gradient methods more efficient. Simulation examples on $l_2$-regularized logistic losses confirm that our method significantly reduces the communication and computational cost with respect to the standard distributed

D. Jakovetić and D. Bajović are with University of Novi Sad, BioSense Center, Novi Sad, Serbia. N. Krejić and N. Krklec-Jerinkić are with Department of Mathematics and Informatics, University of Novi Sad, Novi Sad, Serbia. Authors' e-mails: [djakovet,dbajovic]@uns.ac.rs, natasak@uns.ac.rs, natasa.krklec@dmi.uns.ac.rs.

---

[1]Under a constant step-size $\alpha$, the standard (projected) distributed gradient method converges to a point in a neighborhood of the solution of (1), where the corresponding squared distance is $O(\alpha)$; see ahead Theorem 1 and, e.g., [18], [19].

gradient method, for the same desired accuracy.

The communication and computational savings are highly relevant with applications like WSNs and distributed learning in cluster or cloud environments. With WSNs, the reduced communication and computational cost to retrieve the result translate into energy saving of the sensor motes' batteries and the increase of the network lifetime. With distributed learning in cluster or cloud environments, less amount of communication and computation for a specific application/task means that the saved resources can be re-allocated to another concurrent tasks. For example, at times when a node with our method is idle, the resources allocated to it (e.g., a virtual cloud machine) can be released and re-allocated to other tasks.

We explain intuitively the above results that we achieve. Namely, standard distributed gradient method exhibits, in a sense, two sources of redundancy–the first corresponds to the inter-node communications aspect, while the second corresponds to an optimization aspect (number of gradient evaluations per iteration) of the algorithm. It turns out that, as we show here, a careful simultaneous exploitation of these two redundancies allows us to match the rate of the standard distributed gradient method with a reduced "work." The two sources of redundancy have been already noted in the literature, but have not been exploited simultaneously before. The communication redundancy, e.g., [20] means that the inter-node communications can be "sparsified," e.g., through the intermittent link failures, so that the algorithm still remains convergent. In other words, it is not necessary to utilize communications through all the available links at all iterations for the algorithm to converge. The optimization redundancy has been previously studied only in the context of centralized optimization, e.g., [21]. The core idea is that, under certain assumptions on the cost functions, a (centralized) stochastic-type gradient method with an appropriately increasing sample size matches the convergence rate of the standard gradient method with the full sample size at all iterations, as shown in [21].

We now briefly review existing work relevant to our contributions to help us further contrast our work from the literature. We divide the literature into two classes: 1) distributed gradient methods for multi-agent optimization; and 2) centralized stochastic approximation methods with variable sample sizes. The former class relates to our work through the communication redundancy, while the latter considers the optimization redundancy.

**Distributed gradient methods for multi-agent optimization**. Distributed methods of this type date back at least to the 80s, e.g., [22], and have received renewed interest in the past decade, e.g., [11]. Reference [11] proposes the distributed (sub)gradient method with a constant step-size, and analyzes its performance under time-varying communication networks. Reference [20] considers distributed (sub)gradient method under random communication networks with failing links and establishes almost sure convergence under a diminishing step-size rule. A major difference of our paper from the above works is that, in [22], [11], [20], only inter-node communications over iterations are "sparsified," while each node performs gradient evaluations at each iteration $k$. In [13],

the authors propose a gossip-like scheme where, at each $k$, only two neighboring nodes in the network wake up and perform weight-averaging (communication among them) and the negative gradient step with respect to their respective local costs, while the remaining nodes stay idle. The key difference with respect to our paper is that, with our method, the number of active nodes over iterations $k$ (on average) is increasing, while in [13] it remains equal to two for all $k$. Consequently, the established convergence properties of the two methods are very different.

Our paper is also related to reference [23], which considers diffusion algorithms with two types of nodes – informed and uninformed. The informed nodes both: 1) acquire measurements and perform in-network processing (which translates into computing gradients in our scenario); and 2) perform consultation with neighbors (which translates into weight-averaging the estimates across neighborhoods), while the uninformed nodes only perform the latter task. The authors study the effect of the proportion of informed nodes and their distribution in space. A key difference with respect to our work is that the uninformed nodes in [23] still perform weight-averaging, while the idle nodes here perform no processing. Finally, we comment on reference [24] which introduces an adaptive policy for each node to decide whether it will communicate with its neighbors or not and demonstrate significant savings in communications with respect to the always-communicating scenario. A major difference of [24] from our paper is that, with [24], nodes always perform local gradients, i.e., they do not stay idle (in the sense defined here).

**Centralized stochastic approximation methods with variable sample sizes** have been studied for a long time. We distinguish two types of methods: the ones that assume unbounded sample sizes (where the cost function is in the form of a mathematical expectation) and the methods with bounded sample sizes (where the cost function is of the form in (1).) Our work contrasts with both of these threads of works by considering distributed optimization over an arbitrary connected network, while they consider centralized methods.

Unbounded sample sizes have been studied, e.g., in [25], [26], [27], [28], [29]. Reference [25] uses a Bayesian scheme to determine the sample size at each iteration within the trust region framework, and it shows almost sure convergence to a problem solution. Reference [26] shows almost sure convergence as long as the sample size grows sufficiently fast along iterations. In [27], the variable sample size strategy is obtained as the solution of an associated auxiliary optimization problem. Further references on careful analyses of the increasing sample sizes are, e.g., [28], [29].

References [30], [31] consider a trust region framework and assume bounded sample sizes, but, differently from our paper and [25], [27], [28], [29], [21], they allow the sample size both increase and decrease at each iteration. The paper chooses a sample size at each iteration such that a balance is achieved between the decrease of the cost function and the width of an associated confidence interval. Reference [32] proposes a schedule sequence in the monotone line search framework which also allows the sample size both increase and decrease at each iteration; paper [33] extends the results in [32] to a

non-monotone line search.

Reference [21] is closest to our paper within this thread of works, and our work mainly draws inspiration from it. The authors consider a bounded sample size, as we do here. They consider both deterministic and stochastic sampling and determine the increase of the sample size along iterations such that the algorithm attains (almost) the same rate as if the full sample size was used at all iterations. A major difference of [21] with respect to the current paper is that they are not concerned with the networked scenario, i.e., therein a central entity works with the variable (increasing) sample size. This setup is very different from ours as it has no problem dimension of propagating information across the networked nodes – the dimension present in distributed multi-agent optimization.

**Paper organization**. The next paragraph introduces notation. Section II explains the model that we assume and presents our proposed distributed algorithm. Section III states our main results which we prove in Section IV. Section V provides numerical examples. Finally, we conclude in Section VI.

**Notation**. We denote by: $\mathbb{R}$ the set of real numbers; $\mathbb{R}^d$ the $d$-dimensional real coordinate space; $A_{ij}$ the entry in the $i$-th row and $j$-th column of a matrix $A$; $A^\top$ the transpose of a matrix $A$; $\odot$ and $\otimes$ the Hadamard (entry-wise) and Kronecker product of matrices, respectively; $I$, $0$, $\mathbf{1}$, and $e_i$, respectively, the identity matrix, the zero matrix, the column vector with unit entries, and the $i$-th column of $I$; $J$ the $N \times N$ matrix $J := (1/N)\mathbf{1}\mathbf{1}^\top$; $A \succ 0\,(A \succeq 0)$ means that the symmetric matrix $A$ is positive definite (respectively, positive semi-definite); $\| \cdot \|_l$ the vector (respectively, matrix) $l$-norm of its vector (respectively, matrix) argument; $\| \cdot \| = \| \cdot \|_2$ the Euclidean (respectively, spectral) norm of its vector (respectively, matrix) argument; $\lambda_i(\cdot)$ the $i$-th largest eigenvalue, $\mathrm{Diag}\,(a)$ the diagonal matrix with the diagonal equal to the vector $a$; $|\cdot|$ the cardinality of a set; $\nabla h(w)$ the gradient evaluated at $w$ of a function $h : \mathbb{R}^d \to \mathbb{R}$, $d \geq 1$; $\mathbb{P}(\mathcal{A})$ and $\mathbb{E}[u]$ the probability of an event $\mathcal{A}$ and expectation of a random variable $u$, respectively. For two positive sequences $\eta_n$ and $\chi_n$, we have: $\eta_n = O(\chi_n)$ if $\limsup_{n\to\infty} \frac{\eta_n}{\chi_n} < \infty$.

## II. Model and algorithm

Subsection II-A describes the optimization and network models that we assume, while Subsection II-B presents our proposed distributed algorithm with variable number of working nodes.

### A. Problem model

**Optimization model**. We consider optimization problem (1), and we impose the following assumptions on (1).

*Assumption 1 (Optimization model)* (a) For all $i$, $f_i : \mathbb{R}^d \mapsto \mathbb{R}$ is strongly convex with modulus $\mu > 0$, i.e.:

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2}\|y - x\|^2, \ \forall x, y \in \mathbb{R}^d.$$

(b) For all $i$, $f_i : \mathbb{R}^d \mapsto \mathbb{R}$ has Lipschitz continuous gradient with constant $L$, $0 < \mu \leq L < \infty$, i.e.:

$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|, \ \ \forall x, y \in \mathbb{R}^d.$$

(c) The set $\mathcal{X} \subset \mathbb{R}^d$ is nonempty, closed, convex, and bounded.

We denote by $D := \max\{\|x\| : x \in \mathcal{X}\}$ the diameter of $\mathcal{X}$. Note that, as $\mathcal{X}$ is compact, the gradients $\nabla f_i(x)$'s are bounded over $x \in \mathcal{X}$, i.e., there exists $G > 0$, such that, for all $i$, for all $x \in \mathcal{X}$, $\|\nabla f_i(x)\| \leq G$. The constant $G$ can be taken as $LD + \max_{i=1,\dots,N} \|\nabla f_i(0)\|$. Indeed, for any $x \in \mathcal{X}$, we have:

$$
\begin{aligned}
\|\nabla f_i(x)\| &= \|\nabla f_i(x) - \nabla f_i(0)\| + \|\nabla f_i(0)\| \\
&\leq L\|x\| + \|\nabla f_i(0)\| \\
&\leq LD + \max_{i=1,\dots,N} \|\nabla f_i(0)\|.
\end{aligned}
$$

Similarly, there exist constants $-\infty < m_f \leq M_f < \infty$, such that $m_f \leq f_i(x) \leq M_f$, $\forall i$, $\forall x \in \mathcal{X}$. Constants $m_f$ and $M_f$ can be taken as $M_f = -m_f = GD + \max_{i=1,\dots,N} |f_i(0)|$. Under Assumption 1, (1) is solvable and has a unique solution, which we denote by $x^\star$.

**Network model**. Nodes are connected in a generic undirected network $\mathcal{G} = (\mathcal{V}, E)$, where $\mathcal{V}$ is the set of $N$ nodes and $E$ is the set of edges – all node (unordered) pairs $\{i, j\}$ that can exchange messages through a communication link. We impose the following assumption.

*Assumption 2 (Network connectedness)* The network $\mathcal{G} = (\mathcal{V}, E)$ is connected, undirected, and simple (no self-loops nor multiple links).

Both Assumptions 1 and 2 hold throughout the paper. We denote by $\Omega_i$ the neighborhood set of node $i$ (excluding $i$). We associate with $\mathcal{G}$ a $N \times N$ symmetric weight matrix $C$, which is also stochastic (rows sum to one and all the entries are non-negative). We let $C_{ij}$ be strictly positive for each $\{i, j\} \in E$, $i \neq j$; $C_{ij} = 0$ for $\{i, j\} \notin E$, $i \neq j$; and $C_{ii} = 1 - \sum_{j \neq i} C_{ij}$. As we will see, the weights $C_{ij}$'s will play a role in our distributed algorithm. The quantities $C_{ij}$, $j \in \Omega_i$, are assumed available to node $i$ before execution of the distributed algorithm. We assume that matrix $C$ has strictly positive diagonal entries (each node assigns a non-zero weight to itself) and is positive definite, i.e., $\lambda_N(C) > 0$. For a given arbitrary stochastic, symmetric weight matrix $C'$ with positive diagonal elements, positive definitness may not hold. However, such arbitrary $C'$ can be easily adapted to generate matrix $C$ that obeys all the required properties (symmetric, stochastic, positive diagonal elements), and, in addition, is positive definite. Namely, letting, for some $\kappa \in (0, 1)$, $C := \frac{\kappa+1}{2}I + \frac{1-\kappa}{2}C'$, we obtain that $\lambda_N(C) > \kappa$. It can be shown that, under the above assumptions on $C$, $\lambda_1(C) = 1$, and $\lambda_2(C) < 1$.

### B. Proposed distributed algorithm

We now describe the distributed algorithm to solve (1) that we propose. We assume that all nodes are synchronized according to a global clock and simultaneously (in parallel) perform iterations $k = 0, 1, \dots$ At each iteration $k$, each node $i$ updates its solution estimate $x_i^{(k)} \in \mathcal{X}$, with arbitrary initialization $x_i^{(0)} \in \mathcal{X}$. To avoid notational clutter, we will

assume that $x_i^{(0)} = x_j^{(0)}$, $\forall i,j$. Besides, each node has an internal Bernoulli state variable $z_i^{(k)}$. If $z_i^{(k)} = 1$, node $i$ updates $x_i(k)$ at iteration $k$; we say that, in this case, node $i$ is working at $k$. If $z_i^{(k)} = 0$, node $i$ keeps its current state $x_i(k)$ and does not perform an update; we say that, in this case, node $i$ is idle. At each $k$, each node $i$ generates $z_i^{(k)}$ independently from the previous iterations, and independently from other nodes. We denote by $p_k := \mathbb{P}\left(z_i(k) = 1\right)$. The quantity $p_k$ is our algorithm's tuning parameter, and is common for all nodes. We assume that, for all $k$, $p_k \geq p_{\min}$, for a positive constant $p_{\min}$.

Denote by $\Omega_i^{(k)}$ the set of working neighbors of node $i$ at $k$, i.e., all nodes $j \in \Omega_i$ with $z_j^{(k)} = 1$. The update of node $i$ is as follows. If $z_i^{(k)} = 0$, node $i$ is idle and sets $x_i^{(k+1)} = x_i^{(k)}$. Otherwise, if $z_i^{(k)} = 1$, node $i$ broadcasts its state to all its working neighbors $j \in \Omega_i^{(k)}$. The non-working (idle) neighbors do not receive $x_i^{(k)}$; for example, with WSNs, this corresponds to switching-off the receiving antenna of a node. Likewise, node $i$ receives $x_j^{(k)}$ from all $j \in \Omega_i^{(k)}$. Upon reception, node $i$ updates $x_i^{(k)}$ as follows:

$$x_i^{(k+1)} = \mathcal{P}_{\mathcal{X}}\left\{ \left(1 - \sum_{j\in\Omega_i^{(k)}} C_{ij}\right) x_i^{(k)} \right. \tag{2}$$
$$\left. + \sum_{j\in\Omega_i^{(k)}} C_{ij}\, x_j^{(k)} - \frac{\alpha}{p_k}\nabla f_i(x_i^{(k)}) \right\}.$$

In (2), $\mathcal{P}_{\mathcal{X}}(y) = \arg\min_{v\in\mathcal{X}}\|v - y\|$ denotes the Euclidean projection of point $y$ on $\mathcal{X}$, and $\alpha > 0$ is a constant; we let $\alpha \leq \lambda_N(C)/L$. (See ahead Remark 2.) In words, (2) means that node $i$ makes a convex combination of its own estimate with the estimates of its working neighbors, takes a step in the negative direction of its local gradient, and projects the resulting value onto the constraint set. As we will see, multiplying the step-size in (2) by $1/p_k$ compensates for non-working (idle) nodes over iterations.

*Remark 1* Setting $p_k = 1$, $\forall k$, corresponds to the standard distributed (sub)gradient method in [34].

**Compact representation**. We present (2) in a compact form. Denote by $x^{(k)} := \left((x_1^{(k)})^\top, ..., (x_N^{(k)})^\top\right)^\top$, and $z^{(k)} := (z_1^{(k)}, ..., z_N^{(k)})^\top$. Further, introduce $F : \mathbb{R}^{Nd} \mapsto \mathbb{R}$, with

$$F(x) = F(x_1, ..., x_N) := \sum_{i=1}^{N} f_i(x_i).$$

Also, denote by $\mathcal{X}^N \subset \mathbb{R}^{Nd}$ the Cartesian product $\mathcal{X} \times ... \times \mathcal{X}$, where $\mathcal{X}$ is repeated $N$ times. Next, introduce the $N \times N$ random matrix $W^{(k)}$, defined as follows:

$$W_{ij}^{(k)} = \begin{cases} C_{ij} z_i^{(k)} z_j^{(k)} & \text{for } \{i,j\} \in E,\ i \neq j \\ 0 & \text{for } \{i,j\} \notin E,\ i \neq j \\ 1 - \sum_{s\neq i} W_{is}^{(k)} & \text{for } i = j. \end{cases}$$

Then, it is easy to see verify that, for $k = 0, 1, ...$, update

rule (2) can be written as:

$$\begin{aligned} x^{(k+1)} = \ & \mathcal{P}_{\mathcal{X}^N}\left\{ (W^{(k)} \otimes I)\, x^{(k)} \right. \tag{3} \\ & \left. - \frac{\alpha}{p_k}\left(\nabla F(x^{(k)}) \odot (z^{(k)} \otimes \mathbf{1})\right) \right\}, \end{aligned}$$

where $W^{(k)} \otimes I$ denotes the Kronecker product of $W^{(k)}$ and the $d\times d$ identity matrix, $\mathbf{1}$ in (3) is of size $d\times 1$, and $\odot$ denotes the Hadamard (entry-wise) product. Note that sequence $\{x^{(k)}\}$ is a sequence of random vectors, due to the randomness of the $z^{(k)}$'s. The case $p_k \equiv 1$, $\forall k$, corresponds to standard distributed (sub)gradient method in [11], in which case (3) becomes:

$$x^{(k+1)} = \mathcal{P}_{\mathcal{X}^N}\left\{ (C \otimes I)\, x^{(k)} - \alpha\,\nabla F(x^{(k)}) \right\}. \tag{4}$$

## III. STATEMENT OF MAIN RESULTS

We now present our main results on the proposed distributed method (2). For benchmarking of (2), we first present a result on the convergence of standard distributed gradient algorithm (4). All the results in the current section, together with needed auxiliary results, are proved in Section IV. Recall that $x^\star \in \mathbb{R}^d$ is the solution to (1).

*Theorem 1* Consider standard distributed gradient algorithm (4) with step-size $\alpha \leq \lambda_N(C)/L$. Then, $x^{(k)}$ converges to a point $x^\bullet = \left((x_1^\bullet)^\top, ..., (x_N^\bullet)^\top\right)^\top \in \mathcal{X}^N$ that satisfies, for all $i = 1, ..., N$:

$$\|x_i^\bullet - x^\star\|^2 \leq \|x^\bullet - \mathbf{1}\otimes x^\star\|^2 \leq \alpha\,\mathcal{C}_\Psi \tag{5}$$
$$\mathcal{C}_\Psi := \frac{4N(M_f - m_f)}{1 - \lambda_2(C)} + \frac{2N^2 G^2}{\mu\,(1 - \lambda_2(C))}. \tag{6}$$

Furthermore:

$$\|x^{(k)} - x^\bullet\| \leq 2\sqrt{N}\, D\,(1 - \alpha\mu)^k = O\left((1 - \alpha\mu)^k\right). \tag{7}$$

Theorem 1 says that, with algorithm (4), each node's estimate $x_i^{(k)}$ converges to a point $x_i^\bullet$ in the neighborhood of the true solution $x^\star$; the distance of the limit $x_i^\bullet$ from $x^\star$ is controlled by step-size $\alpha$ – the smaller the step-size, the closer the limit to the true solution. Furthermore, $x_i^{(k)}$ converges to a solution neighborhood (to $x_i^\bullet$) at a globally linear rate, equal to $1-\alpha\mu$. Hence, there is a tradeoff with respect to the choice of $\alpha$: a small $\alpha$ means a higher precision in the limit, but a slower rate to reach this precision. Note also that, for $\alpha \leq \lambda_N(C)/L$, the convergence factor $(1-\alpha\mu)$ does not depend on the underlying network, but the distance $\|x_i^\bullet - x^\star\|$ between arbitrary node $i$'s limit $x_i^\bullet$ and the solution $x^\star$ depends on the underlying network – through the number of nodes $N$ and the second largest eigenvalue of matrix $C$.

*Remark 2* It is possible to extend Theorem 1 to allow also for the step-sizes $\alpha \in (\lambda_N(C)/L,\ (1 + \lambda_N(C))/L]$, in which case the convergence factor $(1 - \alpha\mu)$ in (5) is replaced with $\max\{\alpha L - 1, 1 - \alpha\mu\}$. We restrict ourselves to the case $\alpha \leq \lambda_N(C)/L$, both for simplicity and due to the fact that step-sizes $\alpha$ – needed to achieve sufficient accuracies in practice – are usually much smaller than $1/L$. (See also Section V.)

*Remark 3* For $\alpha \le \lambda_N(C)/L$, the convergence factor $(1 - \alpha\mu)$ is an exact (tight) worst case convergence factor, in the following sense: given an arbitrary network and matrix $C$, and given an arbitrary step-size $\alpha \le \lambda_N(C)/L$, there exists a specific choice of functions $f_i$'s, set $\mathcal{X}$, and initial point $x^{(0)} \in \mathcal{X}^N$, such that $\|x^{(k+1)} - x^\bullet\| = (1-\alpha\mu)\|x^{(k)} - x^\bullet\|$, for all $k = 0, 1, ...$[2]

We benchmark the proposed method against the standard distributed gradient method by checking: 1) whether it converges to the same point $x^\bullet$; 2) if so, whether it converges linearly; and 3) if the convergence is linear, how the corresponding convergence factor compares with $(1 - \alpha\mu)$ – the convergence factor of the standard distributed gradient method.

References [18], [19] also analyze the convergence rate of the standard distributed gradient method, allowing for step-size ranges wider than $\alpha \in (0, \lambda_N(C)/L)$. They establish bounds on quantity $\|x^{(k)} - \mathbf{1} \otimes x^\star\|$ which are in general different than (7), and they are not directly concerned with quantity $\|x^{(k)} - x^\bullet\|$, i.e., precise characterization of convergence rate of $x^{(k)}$ towards its limit. We adopt here (7) as it gives an exact worst-case characterization of the convergence rate towards $x^\bullet$ for $\alpha \in (0, \lambda_N(C)/L]$ (see Remark 3).

We now state our main results on the proposed algorithm (2). The first result deals with a more generic sequence of the $p_k$'s that converge to one; the second result is for the $p_k$'s that converge to one geometrically.

*Theorem 2* Consider algorithm (2) with step-size $\alpha \le \lambda_N(C)/L$. Further, suppose that $p_k \ge p_{\min}$, $\forall k$, for some $p_{\min} > 0$, and let $p_k \to 1$ as $k \to \infty$. Then, with algorithm (2), the iterates $x^{(k)}$ converge, in the mean square sense, to the same point $x^\bullet$ as the standard distributed gradient method (4), i.e., $\mathbb{E}\left[\|x^{(k)} - x^\bullet\|^2\right] \to 0$ as $k \to \infty$.

*Theorem 3* Consider algorithm (2) with step-size $\alpha \le \lambda_N(C)/L$. Further, suppose that $p_k = 1 - \delta^{k+1}$, $k = 0, 1, ...$, for some $\delta \in (0, 1)$, and let $\eta := \max\{1 - \alpha\mu, \delta^{1/2}\}$. Then, in the mean square sense, algorithm (2) converges to the same point $x^\bullet$ as the standard distributed gradient method (4), and, moreover:

$$\mathbb{E}\left[\|x^{(k)} - x^\bullet\|\right] = O\left(k\eta^k\right) = O\left((\eta + \epsilon)^k\right),$$

for arbitrarily small positive $\epsilon$. Furthermore, if $\sqrt{\delta} \le 1 - \alpha\mu$:

$$\mathbb{E}\left[\|x^{(k)} - x^\bullet\|\right] = O\left(k(1-\alpha\mu)^k\right) = O\left((1-\alpha\mu+\epsilon)^k\right).$$

Theorem 2 states that, provided that the $p_k$'s are uniformly bounded away from zero, from below, and $p_k \to 1$, the method (4) converges (in the mean square) to the same point as the standard distributed method (4). That is, the random idling schedule governed by the $p_k$'s does not affect the method's limit. Theorem 3 furthermore suggests that, provided that $\sqrt{\delta} \le 1 - \alpha\mu$, (2) converges at the practically same rate as

---

[2]Consider $f_i : \mathbb{R} \to \mathbb{R}$, $f_i(x) = x^2$, $\forall i$, $\mathcal{X} = \{x \in \mathbb{R} : |x| \le 2\}$, and $x^{(0)} = 1$. Note that, in this case, $x^\bullet = 0$, and $\mu = L = 1$. For this example, it is easy to show that $\|x^{(k+1)} - x^\bullet\| = (1 - \alpha)\|x^{(k)} - x^\bullet\|$, for all $k = 0, 1, ...$, and so the convergence factor equals $1 - \alpha\mu$.

the standard method (4), i.e., as if all the nodes were working all the time. (albeit with a larger hidden constant). Hence, we may expect that the proposed method (2) achieves the same desired accuracy as (4) with less amount of resources spent (smaller number of the overall node works–activations). This indeed occurs in practice, as confirmed by simulations in Section V. The hidden convergence constant is dependent on the underlying network, the sequence $\{p_k\}$, and step-size $\alpha$, and is given explicitly in Remark 5.

## IV. INTERMEDIATE RESULTS AND PROOFS

Subsection IV-A gives intermediate results on the random matrices $W^{(k)}$ and provides the disagreement estimates – how far apart are the estimates $x_i^{(k)}$ of different nodes in the network. Subsection IV-B introduces a penalty-like interpretation of algorithm (4) and proves Theorem 1. Finally, Subsection IV-C proves our main results, Theorems 2 and 3, by applying the penalty-like interpretation on algorithm (3). For notational simplicity, this section presents auxiliary results and all the proofs for the case $d = 1$, but all these extend to a generic $d > 1$. Throughout this Section, all the claims (equalities and inequalities) which deal with random quantities hold either: 1) surely, for any random realization; or 2) in expectation. It is clear from notation which of the two cases is in force.

### A. Matrices $W^{(k)}$ and disagreement estimates

**Matrices $W^{(k)}$.** Recall that $J := (1/N)\mathbf{1}\mathbf{1}^\top$. We have the following Lemma on the matrices $W^{(k)}$. Lemma 4 follows from simple arguments and standard results on symmetric, stochastic matrices (see, e.g., [35]). Hence, we omit the proof for brevity.

*Lemma 4 (Matrices $W^{(k)}$)* (a) The sequence $\{W^{(k)}\}$ is a sequence of independent random matrices.
(b) For all $k$, $W^{(k)}$ is symmetric and stochastic (rows sum to one and all the entries are nonnegative).
(c) For all $k$, $0 \prec W^{(k)} \preceq I$.
(d) There exists a constant $\beta \in (0, 1)$ such that, $\forall k$, $\mathbb{E}\left[\|W^{(k)} - J\|^2\right] < \beta^2$.

It can be shown that $\beta$ can be taken as $\beta^2 = 1 - (p_{\min})^N\left[1 - (\lambda_2(C))^2\right]$; see, e.g., [35].

*Remark 4* Quantities $\mathbb{E}\left[\|W^{(k)} - J\|^2\right]$ clearly depend on $k$, and, more specifically, on $p_k$. We adopt here a (possibly loose) uniform bound $\beta$ (independent of $k$) as this suffices to establish conclusions about convergence rates of algorithm (3) while simplifying the presentation.

**Disagreement estimate.** Denote by $\overline{x}^{(k)} := \frac{1}{N}\sum_{i=1}^N x_i^{(k)}$ the global average of the nodes' estimates, and by $\widetilde{x}_i^{(k)} = x_i^{(k)} - \overline{x}^{(k)}$. Note that both quantities are random. The quantity $\widetilde{x}_i^{(k)}$ measures how far is the node $i$'s estimate from the global average. Denote by $\widetilde{x}^{(k)} := (\widetilde{x}_1^{(k)}, ..., \widetilde{x}_N^{(k)})^\top$. The next Lemma shows that $\mathbb{E}\left[\|\widetilde{x}^{(k)}\|^2\right]$ is uniformly bounded, $\forall k$, and that the bound is $O(\alpha^2)$, i.e., the disagreement size is controlled

by the step-size. (The smaller the step-size, the smaller the disagreements are.)

*Lemma 5 (Disagreements bound)* For all $k$, there holds:

$$\mathbb{E}\left[\|\widetilde{x}^{(k)}\|^2\right] \le \left(\frac{3\alpha\sqrt{N}G}{p_{\min}(1-\beta)}\right)^2.$$

*Proof:* Consider (2), and denote by:

$$
\begin{aligned}
y_i^{(k)} &:= \left(1 - \sum_{j\in\Omega_i^{(k)}} C_{ij}\right)x_i^{(k)} + \sum_{j\in\Omega_i^{(k)}} C_{ij}x_j^{(k)} \\
&\quad - \frac{\alpha}{p_k}\nabla f_i(x_i^{(k)}) \\
&= \sum_{j\in\Omega_i\cup\{i\}} W_{ij}^{(k)}x_j^{(k)} - \frac{\alpha}{p_k}\nabla f_i(x_i^{(k)}).
\end{aligned}
$$

Also, let $\epsilon^{(k)} = (\epsilon_1^{(k)},...,\epsilon_N^{(k)})^\top$, where $\epsilon_i^{(k)} := \mathcal{P}_{\mathcal{X}}\left\{y_i^{(k)}\right\} - y_i^{(k)}$. Then, (3) can be written in the following equivalent form:

$$x^{(k+1)} = W^{(k)}x^{(k)} - \frac{\alpha}{p_k}\left(\nabla F(x^{(k)})\odot z^{(k)}\right) + \epsilon^{(k)}. \quad (8)$$

We first upper bound $\|\epsilon^{(k)}\|$. Consider $\epsilon_i^{(k)}$. We have:

$$
\begin{aligned}
|\epsilon_i^{(k)}| &= \left|\mathcal{P}_{\mathcal{X}}\left\{\sum_{j\in\Omega_i\cup\{i\}}W_{ij}^{(k)}x_j^{(k)} - \frac{\alpha}{p_k}\nabla f_i(x_i^{(k)})\right\}\right. \\
&\quad \left. - \left(\sum_{j\in\Omega_i\cup\{i\}}W_{ij}^{(k)}x_j^{(k)} - \frac{\alpha}{p_k}\nabla f_i(x_i^{(k)})\right)\right| \\
&= \left|\mathcal{P}_{\mathcal{X}}\left\{\sum_{j\in\Omega_i\cup\{i\}}W_{ij}^{(k)}x_j^{(k)} - \frac{\alpha}{p_k}\nabla f_i(x_i^{(k)})\right\}\right. \\
&\quad \left. - \mathcal{P}_{\mathcal{X}}\left\{\sum_{j\in\Omega_i\cup\{i\}}W_{ij}^{(k)}x_j^{(k)}\right\} + \frac{\alpha}{p_k}\nabla f_i(x_i^{(k)})\right| \quad (9) \\
&\le \left|\mathcal{P}_{\mathcal{X}}\left\{\sum_{j\in\Omega_i\cup\{i\}}W_{ij}^{(k)}x_j^{(k)} - \frac{\alpha}{p_k}\nabla f_i(x_i^{(k)})\right\}\right. \\
&\quad \left. - \mathcal{P}_{\mathcal{X}}\left\{\sum_{j\in\Omega_i\cup\{i\}}W_{ij}^{(k)}x_j^{(k)}\right\}\right| \\
&\quad + \left|\frac{\alpha}{p_k}\nabla f_i(x_i^{(k)})\right| \quad (10) \\
&\le \frac{2\alpha}{p_k}|\nabla f_i(x_i^{(k)})| \le \frac{2\alpha G}{p_{\min}}. \quad (11)
\end{aligned}
$$

Equality (9) holds because $\sum_{j\in\Omega_i\cup\{i\}}W_{ij}^{(k)}x_j^{(k)} \in \mathcal{X}$, as a convex combination of the $x_j^{(k)}$ that belong to $\mathcal{X}$ by construction, and due to convexity of $\mathcal{X}$. Inequality (10) is by the triangle inequality. Finally, (11) is by the non-expansiveness property of the Euclidean projection: $|\mathcal{P}_{\mathcal{X}}\{u\} - \mathcal{P}_{\mathcal{X}}\{v\}| \le |u - v|$, $\forall u, v \in \mathbb{R}$. Therefore, we obtain the following bound on $\|\epsilon^{(k)}\|$:

$$\|\epsilon^{(k)}\| \le \frac{2\sqrt{N}\,\alpha\,G}{p_{\min}}. \quad (12)$$

We now return to (8). Note that $\widetilde{x}^{(k)} = (I - J)x^{(k)}$. Also, $W^{(k)}J = JW^{(k)} = J$, by Lemma 4 (b). Thus, we have that: $(I-J)W^{(k)} = W^{(k)} - J$. Also, $(W^{(k)} - J)(I - J) = W^{(k)} - J - JW^{(k)} + J^2 = W^{(k)} - J$, because $JW^{(k)} = J$ and $J^2 = J$. Using the latter, and multiplying (8) from the left by $(I - J)$, we obtain:

$$
\begin{aligned}
\widetilde{x}^{(k+1)} &= (W^{(k)} - J)\,\widetilde{x}^{(k)} \\
&\quad - \frac{\alpha}{p_k}(I - J)\left(\nabla F(x^{(k)})\odot z^{(k)}\right) + (I - J)\epsilon^{(k)}.(13)
\end{aligned}
$$

Taking the norm and using sub-additive and sub-multiplicative properties:

$$
\begin{aligned}
\|\widetilde{x}^{(k+1)}\| &\le \|W^{(k)} - J\|\,\|\widetilde{x}^{(k)}\| \\
&\quad + \frac{\alpha}{p_k}\|(I - J)\left(\nabla F(x^{(k)})\odot z^{(k)}\right)\| \\
&\quad + \|(I - J)\epsilon^{(k)}\|. \quad (14)
\end{aligned}
$$

It is easy to see that $\|\left(\nabla F(x^{(k)})\odot z^{(k)}\right)\| \le \sqrt{N}G$. Hence, using the sub-multiplicative property of norms and the fact that $\|I - J\| = 1$, there holds: $\|(I - J)\left(\nabla F(x^{(k)})\odot z^{(k)}\right)\| \le \sqrt{N}G$. Similarly, from (12): $\|(I - J)\epsilon^{(k)}\| \le \frac{2\sqrt{N}\,\alpha\,G}{p_{\min}}$. Combining the latter conclusions with (14):

$$\|\widetilde{x}^{(k+1)}\| \le \|W^{(k)} - J\|\,\|\widetilde{x}^{(k)}\| + \frac{3\sqrt{N}\,\alpha\,G}{p_{\min}}. \quad (15)$$

Squaring the latter inequality, we obtain:

$$
\begin{aligned}
\|\widetilde{x}^{(k+1)}\|^2 &\le \|W^{(k)} - J\|^2\,\|\widetilde{x}^{(k)}\|^2 \\
&\quad + \frac{6\sqrt{N}\,\alpha\,G}{p_{\min}}\|W^{(k)} - J\|\,\|\widetilde{x}^{(k)}\| \\
&\quad + \left(\frac{3\sqrt{N}\,\alpha\,G}{p_{\min}}\right)^2. \quad (16)
\end{aligned}
$$

Taking expectation, using the independence of $W^{(k)}$ and $\widetilde{x}^{(k)}$, and the inequality $\mathbb{E}[|u|] \le \left(\mathbb{E}[u^2]\right)^{1/2}$, for a random variable $u$, we obtain:

$$
\begin{aligned}
\mathbb{E}\left[\|\widetilde{x}^{(k+1)}\|^2\right] &\le \mathbb{E}\left[\|W^{(k)} - J\|^2\right]\mathbb{E}\left[\|\widetilde{x}^{(k)}\|^2\right] \\
&\quad + \frac{6\sqrt{N}\,\alpha\,G}{p_{\min}}\left(\mathbb{E}\left[\|W^{(k)} - J\|^2\right]\right)^{1/2}\left(\mathbb{E}\left[\|\widetilde{x}^{(k)}\|^2\right]\right)^{1/2} \\
&\quad + \left(\frac{3\sqrt{N}\,\alpha\,G}{p_{\min}}\right)^2. \quad (17)
\end{aligned}
$$

Denote by $\xi^{(k)} := \left(\mathbb{E}\left[\|\widetilde{x}^{(k)}\|^2\right]\right)^{1/2}$. Applying Lemma 4 (d), writing the right-hand side of (17) as a complete square, and taking the square root of the resulting inequality:

$$\xi^{(k+1)} \le \beta\,\xi^{(k)} + \frac{3\sqrt{N}\,\alpha\,G}{p_{\min}}. \quad (18)$$

Unwinding recursion (18), and using the bound $1 + \beta + ... + \beta^k \le \frac{1}{1-\beta}$, we obtain that, for all $k = 0, 1, ...$, there holds:

$$\xi^{(k)} \le \frac{3\alpha\sqrt{N}G}{p_{\min}(1-\beta)}.$$

Squaring the last inequality, the desired result follows. ∎

### B. Analysis of the standard distributed gradient method through a penalty-like reformulation

We analyze the proposed method (3) through a penalty-like interpretation, to our best knowledge first introduced in [36]. Introduce an auxiliary function $\Psi_\alpha : \mathbb{R}^N \mapsto \mathbb{R}$, defined by: $\Psi_\alpha(x) := \sum_{i=1}^N f_i(x_i) + \frac{1}{2\alpha} x^\top (I - C)x = F(x) + \frac{1}{2\alpha} x^\top (I - C)x$, and the associated optimization problem:

$$\text{minimize } \Psi_\alpha(x) = \sum_{i=1}^N f_i(x_i) + \frac{1}{2\alpha} x^\top (I - C)x \\ \text{subject to } x \in \mathcal{X}^N. \quad (19)$$

Function $\Psi_\alpha$ and (19) will be very useful in the analysis of (2). In fact, we will show that (2) is an inexact version of the (projected) gradient method on function $\Psi_\alpha$. Clearly, (19) is solvable, and it has a unique solution, which we denote by $x^\bullet$.[3]

We start by showing that standard distributed (sub)gradient method in [11] is an exact (projected) gradient method on $\Psi_\alpha$. Indeed, the derivative $\nabla \Psi_\alpha(x) = \nabla F(x) + \frac{1}{\alpha}(I - C)x$. The projected gradient method on $\Psi_\alpha$ with step-size $\alpha$ then takes the form:

$$\begin{aligned} x^{(k+1)} &= \mathcal{P}_{\mathcal{X}^N} \left\{ x^{(k)} - \alpha \nabla \Psi_\alpha(x^{(k)}) \right\} \quad (20) \\ &= \mathcal{P}_{\mathcal{X}^N} \left\{ x^{(k)} - \right. \\ &\quad \left. \alpha \left( \nabla F(x^{(k)}) + \frac{1}{\alpha}(I - C)x^{(k)} \right) \right\}, \end{aligned}$$

which, after rearranging terms, is precisely (4).

It is easy to see that $\Psi_\alpha$ is strongly convex on $\mathbb{R}^N$, with modulus $\mu' = \mu$ (equal the strong convexity modulus of the $f_i$'s). Further, $\nabla \Psi_\alpha$ is Lipschitz continuous on $\mathbb{R}^N$, with constant $L' = L + \frac{1 - \lambda_N(C)}{\alpha}$. Namely, $\forall x, y \in \mathbb{R}^N$:

$$\begin{aligned} \|\nabla \Psi_\alpha(x) - \nabla \Psi_\alpha(y)\| &\leq \|\nabla F(x) - \nabla F(y)\| \\ &+ \frac{1}{\alpha} \|I - C\| \, \|x - y\| \\ &\leq L\|x - y\| + \frac{1 - \lambda_N(C)}{\alpha} \|x - y\|. \end{aligned}$$

(Note that $\|\nabla F(x) - \nabla F(y)\| \leq L\|x - y\|$ follows after summing the inequalities: $|\nabla f_i(x_i) - \nabla f_i(y_i)|^2 \leq L^2 |x_i - y_i|^2$, $i = 1, ..., N$, and using $\|\nabla F(x)\|^2 = \sum_{i=1}^N |\nabla f_i(x_i)|^2$.) We impose that $\alpha$ satisfies $\alpha \leq \frac{1}{L'}$, which, after simple manipulations, gives: $\alpha \leq \lambda_N(C)/(L)$, as introduced before.

An immediate consequence of the fact that algorithm (4) is precisely the projected gradient method to solve (19) is the following Lemma, first observed in [36].

*Lemma 6 ([36])* Standard distributed gradient algorithm (4) with step-size $\alpha \leq \lambda_N(C)/L$ converges to point $x^\bullet \in \mathcal{X}^N$ – the solution to (19).

We proceed by proving Theorem 1.

*Proof of Theorem 1:* As per Lemma 6, algorithm (4) converges to $x^\bullet$ – the solution to (19). We hence need to prove for the solution to (19) the characterization in (5).

---

[3]The point of convergence of algorithm (4) and the solution to (19) are intentionally denoted by the same symbol because – as we will show – they actually are the same point.

Consider an arbitrary point $x \in \mathcal{X}^N$, and let $\overline{x} := \frac{1}{N} \sum_{i=1}^N x_i$. We first prove the following inequality:

$$f(\overline{x}) - f(x^\star) \leq (\Psi_\alpha(x) - \Psi_\alpha(x^\bullet)) + \frac{\alpha N G^2}{2(1 - \lambda_2(C))}. \quad (21)$$

Indeed, we have that:

$$\begin{aligned} x^\top (I - C)x &= (x - \overline{x}\mathbf{1})^\top (I - C)(x - \overline{x}\mathbf{1}) \\ &\geq \lambda_{N-1}(I - C)\|x - \overline{x}\mathbf{1}\|^2 \\ &= (1 - \lambda_2(C))\|\widetilde{x}\|^2, \end{aligned}$$

where we let $\widetilde{x} := x - \overline{x}\mathbf{1}$. Further,

$$\begin{aligned} \sum_{i=1}^N f_i(x_i) &= \sum_{i=1}^N f_i(\overline{x}) + \left( \sum_{i=1}^N (f_i(x_i) - f_i(\overline{x})) \right) \\ &\geq f(\overline{x}) - G \sum_{i=1}^N |x_i - \overline{x}| \\ &\geq f(\overline{x}) - G\sqrt{N}\|\widetilde{x}\|. \end{aligned}$$

The second from last inequality follows because $f_i(x_i) \geq f_i(\overline{x}) + \nabla f_i(\overline{x})(x_i - \overline{x}) \geq f_i(\overline{x}) - G|x_i - \overline{x}|$. Combining the previous conclusions:

$$\begin{aligned} \Psi_\alpha(x) - \Psi_\alpha(x^\bullet) &\geq f(\overline{x}) - \Psi_\alpha(x^\bullet) - G\sqrt{N}\|\widetilde{x}\| \\ &+ \frac{1}{2\alpha}(1 - \lambda_2(C))\|\widetilde{x}\|^2 \\ &\geq f(\overline{x}) - \Psi_\alpha(x^\bullet) \\ &- \sup_{t \geq 0} \left\{ G\sqrt{N}t - \frac{1}{2\alpha}(1 - \lambda_2(C))t^2 \right\} \\ &\geq f(\overline{x}) - \Psi_\alpha(x^\bullet) - \frac{\alpha N G^2}{2(1 - \lambda_2(C))}. \quad (22) \end{aligned}$$

Next, note that $\Psi_\alpha(x^\bullet) = \min_{x \in \mathcal{X}^N} \Psi_\alpha(x) \leq \Psi_\alpha(x^\star \mathbf{1}) = f(x^\star)$, and so $-\Psi_\alpha(x^\bullet) \geq -f(x^\star)$. Applying this to (22), completes the proof of (21).

We now prove claim (5) in Theorem 1. We have:

$$\begin{aligned} \|x^\bullet - x^\star \mathbf{1}\|^2 &= \|x^\bullet - \overline{x}^\bullet \mathbf{1} + \overline{x}^\bullet \mathbf{1} - x^\star \mathbf{1}\|^2 \\ &\leq 2\|x^\bullet - \overline{x}^\bullet \mathbf{1}\|^2 + 2N|\overline{x}^\bullet - x^\star|^2. \quad (23) \end{aligned}$$

For the second summand in (23), we have:

$$\begin{aligned} |\overline{x}^\bullet - x^\star|^2 &\leq \frac{2}{\mu}(f(\overline{x}^\bullet) - f(x^\star)) \\ &\leq \frac{\alpha N G^2}{\mu(1 - \lambda_2(C))}. \end{aligned}$$

The first inequality above is due to strong convexity of $f$, and the second applies (21) with $x = x^\bullet$ (where $\overline{x}^\bullet = \frac{1}{N} \sum_{i=1}^N x_i^\bullet$). We now upper bound the first summand in (23). We have that:

$$\begin{aligned} \Psi_\alpha(x^\bullet) &= \sum_{i=1}^N f_i(x_i^\bullet) + \frac{1}{2\alpha}(x^\bullet)^\top (I - C)x^\bullet \\ &\geq \frac{1 - \lambda_2(C)}{2\alpha}\|\widetilde{x}^\bullet\|^2 + N m_f, \end{aligned}$$

where $\widetilde{x}^\bullet = x^\bullet - \overline{x}^\bullet \mathbf{1}$. On the other hand,

$$\Psi_\alpha(x^\bullet) \leq f(x^\star) \leq N M_f.$$

Combining the obtained upper and lower bounds on $\Psi_\alpha(x^\bullet)$, we obtain for the first summand in (23):

$$\|x^\bullet - \overline{x}^\bullet \mathbf{1}\|^2 \le \frac{2\alpha N(M_f - m_f)}{1 - \lambda_2(C)}.$$

Combining the bounds on the first and second summands, the claim in (5) follows.

It remains to prove the claim in (7). By standard analysis of gradient methods, we have that:

$$\|x^{(k)} - x^\bullet\| \le (1 - \alpha\mu)^k\|x^{(0)} - x^\bullet\| \le (1 - \alpha\mu)^k 2\sqrt{N}D,$$

where we used that $\|x^{(0)}\| \le \sqrt{N}D$, and the same bound for $x^\bullet$. Thus, the desired result. ∎

### C. Analysis of the proposed method (2)

We now turn our attention to the proposed method (3). It is easy to verify that (3) can be written as:

$$x^{(k+1)} = \mathcal{P}_{\mathcal{X}^N}\left\{x^{(k)} - \alpha\left[\nabla\Psi_\alpha(x^{(k)}) + e^{(k)}\right]\right\}, \quad (24)$$

where $e^{(k)} = (e_1^{(k)}, ..., e_N^{(k)})^\top$ is a random vector, with $i$-th component equal to:

$$
\begin{aligned}
e_i^{(k)} &= \left(\frac{z_i^{(k)}}{p_k} - 1\right)\nabla f_i(x_i^{(k)}) \\
&+ \frac{1}{\alpha}\sum_{j\in\Omega_i} C_{ij}(z_i^{(k)}z_j^{(k)} - 1)\left(x_i^{(k)} - x_j^{(k)}\right). \quad (25)
\end{aligned}
$$

Hence, (2) is an inexact projected gradient method applied to $\Psi_\alpha$, with step-size $\alpha$, where the amount of inexactness is given by vector $e^{(k)}$.

Overall, our strategy in analyzing (24) consists of two main steps: 1) analyzing the inexact projected gradient method (24); and 2) characterizing (upper bounding) the inexactness vector $e^{(k)}$. For the former step, we apply Proposition 3 in [21]. Adapted to our setting, the proposition says the following. Consider minimization of $\phi(y)$ over $y \in \mathcal{Y}$, where $\phi : \mathbb{R}^m \to \mathbb{R}$ is a convex function, and $\mathcal{Y} \subset \mathbb{R}^m$ is a closed convex set. Let $y^\bullet$ be the solution to the above problem. Further, let $\phi$ be strongly convex with modulus $\mu_\phi > 0$, and let $\phi$ have a Lipschitz continuous gradient with constant $L_\phi \ge \mu_\phi$.

*Lemma 7 (Proposition 3, [21])* Consider the algorithm:

$$y^{(k+1)} = \mathcal{P}_{\mathcal{Y}}\left\{y^{(k)} - \frac{1}{L_\phi}\left[\nabla\phi(y^{(k)}) + e_y^{(k)}\right]\right\}, \; k = 0, 1, ...,$$

where $e_y^{(k)}$ is a random vector. Then, $\forall k = 1, 2, ...$:

$$
\begin{aligned}
\|y^{(k)} - y^\bullet\| &\le (1 - \mu_\phi/L_\phi)^k\|y^{(0)} - y^\bullet\| \\
&+ \frac{1}{L_\phi}\sum_{t=1}^k(1 - \mu_\phi/L_\phi)^{k-t}\|e_y^{(t-1)}\|, \quad (26)
\end{aligned}
$$

where $y^{(0)} \in \mathcal{Y}$ is the initial point.

Note that, if $\nabla\phi$ is Lipschitz continuous with constant $L_\phi$, then $\nabla\phi$ is also Lipschitz continuous with constant $1/\alpha \ge L_\phi$. Therefore, for the function $\phi$ and the iterations:

$$y^{(k+1)} = \mathcal{P}_{\mathcal{Y}}\left\{y^{(k)} - \alpha\left[\nabla\phi(y^{(k)}) + e_y^{(k)}\right]\right\}, \; k = 0, 1, ...,$$

there holds:

$$
\begin{aligned}
\|y^{(k)} - y^\bullet\| &\le (1 - \alpha\,\mu_\phi)^k\|y^{(0)} - y^\bullet\| \\
&+ \alpha\sum_{t=1}^k(1 - \alpha\,\mu_\phi)^{k-t}\|e_y^{(t-1)}\|, \; k = 1, .\text{(27)}
\end{aligned}
$$

In other words, the modified claim (27) holds even if we take a step size different (smaller than) $1/L_\phi$.

For analyzing the inexact projected gradient method (24), we will also make use of the following result in [12].

*Lemma 8 (Lemma 3.1, [12])* Consider a sequence $\{u_k\}$ such that $u_k \to 0$ as $k \to \infty$. Then, for any constant $a \in (0, 1)$ there holds:

$$\sum_{t=1}^k a^{k-t}u_t \to 0.$$

**Step 1: gradient inexactness**. We proceed by characterizing the gradient inexactness; Lemma 9 upper bounds quantity $\mathbb{E}\left[\|e^{(k)}\|^2\right]$.

*Lemma 9 (Gradient inexactness)* For all $k = 0, 1, ...$, there holds:

$$
\begin{aligned}
\mathbb{E}\left[\|e^{(k)}\|^2\right] &\le 4(1 - p_k)\frac{N\,G^2}{p_{\min}} \\
&+ 72(1 - p_k^2)\frac{N G^2}{(p_{\min})^2(1 - \beta)^2} \\
&\le \mathcal{C}_e(1 - p_k^2), \quad (28)
\end{aligned}
$$

where

$$\mathcal{C}_e = \frac{4\,N\,G^2}{p_{\min}} + \frac{72\,N G^2}{(p_{\min})^2(1 - \beta)^2}. \quad (29)$$

*Proof:* Consider (25). We have:

$$
\begin{aligned}
|e_i^{(k)}|^2 &\le 2\left|\frac{z_i^{(k)}}{p_k} - 1\right|^2|\nabla f_i(x_i^{(k)})|^2 \quad (30) \\
&+ \frac{2}{\alpha^2}\sum_{j\in\Omega_i}C_{ij}|z_i^{(k)}z_j^{(k)} - 1|^2\left|x_i^{(k)} - x_j^{(k)}\right|^2 \\
&\le 2G^2\left|\frac{z_i^{(k)}}{p_k} - 1\right|^2 + \frac{4}{\alpha^2}\sum_{j\in\Omega_i}C_{ij}|z_i^{(k)}z_j^{(k)} - 1|^2 \\
&\times \left(\left|\widetilde{x}_i^{(k)}\right|^2 + \left|\widetilde{x}_j^{(k)}\right|^2\right). \quad (31)
\end{aligned}
$$

Inequality (30) uses the following bound: $(u + v)^2 \le 2u^2 + 2v^2$. It also uses, with $u_i := (z_i^{(k)}z_j^{(k)} - 1)(x_i^{(k)} - x_j^{(k)})$, the following relation:

$$
\begin{aligned}
\left(\sum_{j\in\Omega_i}C_{ij}u_j\right)^2 &= \left(\sum_{j\in\Omega_i}C_{ij}u_j + C_{ii}\cdot 0\right)^2 \\
&\le \sum_{j\in\Omega_i}C_{ij}u_j^2 + C_{ii}\cdot 0^2 \\
&= \sum_{j\in\Omega_i}C_{ij}u_j^2,
\end{aligned}
$$

which follows due to the fact that $\sum_{j\in\Omega_i}C_{ij}u_j + C_{ii}\cdot 0$ is a convex combination, and $v \mapsto v^2$, $v \in \mathbb{R}$, is convex.

Inequality (31) uses that

$$\begin{aligned}
\left|x_i^{(k)} - x_j^{(k)}\right|^2 &= \left|x_i^{(k)} - \overline{x}^{(k)} + \overline{x}^{(k)} - x_j^{(k)}\right|^2 \\
&\leq 2\left|x_i^{(k)} - \overline{x}^{(k)}\right|^2 + 2\left|\overline{x}^{(k)} - x_j^{(k)}\right|^2.
\end{aligned}$$

Taking expectation, and using independence of $x^{(k)}$ from $z^{(k)}$:

$$\begin{aligned}
\mathbb{E}\left[|e_i^{(k)}|^2\right] &\leq 2G^2\,\mathbb{E}\left[\left|\frac{z_i^{(k)}}{p_k} - 1\right|^2\right] \\
&\quad + \frac{4}{\alpha^2}\sum_{j\in\Omega_i} C_{ij}\,\mathbb{E}\left[|z_i^{(k)} z_j^{(k)} - 1|^2\right] \\
&\quad \times \left(\mathbb{E}\left[\left|\widetilde{x}_i^{(k)}\right|^2\right] + \mathbb{E}\left[\left|\widetilde{x}_j^{(k)}\right|^2\right]\right). \quad (32)
\end{aligned}$$

We proceed by upper bounding $\mathbb{E}\left[\left|\frac{z_i^{(k)}}{p_k} - 1\right|^2\right]$, using the total probability law with respect to the following partition: $\{z_i^{(k)} = 1\}$, and $\{z_i^{(k)} = 0\}$:

$$\begin{aligned}
\mathbb{E}\left[\left|\frac{z_i^{(k)}}{p_k} - 1\right|^2\right] &= \left|\frac{1}{p_k} - 1\right|^2 \mathbb{P}(z_i^{(k)} = 1) \\
&= \left|\frac{1}{p_k} - 1\right|^2 p_k + (1 - p_k) \quad (33) \\
&= \frac{1}{p_k}(1 - p_k)^2 + (1 - p_k) \\
&\leq \frac{1}{p_k}(1 - p_k) + (1 - p_k) \\
&\leq 2(1 - p_k)/p_{\min}. \quad (34)
\end{aligned}$$

We next upper bound $\mathbb{E}\left[|z_i^{(k)} z_j^{(k)} - 1|^2\right]$, using the total probability law with respect to the event $\{z_i^{(k)} = 1,\, z_j^{(k)} = 1\}$ and its complement; we obtain:

$$\begin{aligned}
\mathbb{E}\left[|z_i^{(k)} z_j^{(k)} - 1|^2\right] &= (1 - \mathbb{P}(z_i^{(k)} = 1,\, z_j^{(k)} = 1)) \\
&= (1 - p_k)^2. \quad (35)
\end{aligned}$$

Substituting (34) and (35) in (32):

$$\begin{aligned}
\mathbb{E}\left[|e_i^{(k)}|^2\right] &\leq 4G^2\,(1 - p_k)/p_{\min} \\
&\quad + \frac{4}{\alpha^2}\sum_{j\in\Omega_i} C_{ij}\,(1 - p_k^2) \\
&\quad \times \left(\mathbb{E}\left[\left|\widetilde{x}_i^{(k)}\right|^2\right] + \mathbb{E}\left[\left|\widetilde{x}_j^{(k)}\right|^2\right]\right). \quad (36)
\end{aligned}$$

Summing the above inequalities over $i = 1, ..., N$, using the fact that $\sum_{j\in\Omega_i} C_{ij} \leq 1,\ \forall i$, $\mathbb{E}\left[\|e^{(k)}\|^2\right] = \sum_{i=1}^N \mathbb{E}\left[|e_i^{(k)}|^2\right]$, and $\mathbb{E}\left[\|\widetilde{x}^{(k)}\|^2\right] = \sum_{i=1}^N \mathbb{E}\left[|\widetilde{x}_i^{(k)}|^2\right]$, we obtain:

$$\begin{aligned}
\mathbb{E}\left[\|e^{(k)}\|^2\right] &\leq 4NG^2\,(1 - p_k)/p_{\min} \\
&\quad + \frac{8}{\alpha^2}(1 - p_k^2)\,\mathbb{E}\left[\left\|\widetilde{x}^{(k)}\right\|^2\right].
\end{aligned}$$

Finally, applying Lemma 5 to the last inequality, the claim follows. ∎

**Step 2: Analyzing the inexact projected gradient method**. We first state and prove the following Lemma on algorithm (2).

*Lemma 10* Consider algorithm (2) with step-size $\alpha \leq \lambda_N(C)/(L)$. Then, for the iterates $x^{(k)}$ and $x^\bullet$–the solution to (19), $\forall k = 1, 2, ...$, there holds:

$$\begin{aligned}
\mathbb{E}\left[\|x^{(k)} - x^\bullet\|^2\right] &\leq 8N\,(1 - \alpha\mu)^{2k}D^2 \\
&\quad + \frac{\alpha\,\mathcal{C}_e}{\mu}\sum_{t=1}^k (1 - \alpha\,\mu)^{k-t}(1 - p_{t-1}^2).
\end{aligned}$$

*Proof:* As already established, algorithm (2) is an inexact projected gradient method to solve (19), with the inexactness vector $e^{(k)}$. We now apply (27) to sequence $x^{(k)}$ and iterations (3); we obtain:

$$\begin{aligned}
\|x^{(k)} - x^\bullet\| &\leq (1 - \alpha\,\mu)^k \|x^{(0)} - x^\bullet\| \\
&\quad + \alpha\sum_{t=1}^k (1 - \alpha\mu)^{k-t}\|e^{(t-1)}\|. \quad (37)
\end{aligned}$$

Squaring the latter inequality, using $(u + v)^2 \leq 2u^2 + 2v^2$, and $\|x^{(0)} - x^\bullet\| \leq 2\sqrt{N}D$:

$$\begin{aligned}
\|x^{(k)} - x^\bullet\|^2 &\leq 8(1 - \alpha\,\mu)^{2k}ND^2 \\
&\quad + \alpha^2\left(\sum_{t=0}^k (1 - \alpha\mu)^{k-t}\right) \\
&\quad \times \sum_{t=1}^k (1 - \alpha\mu)^{k-t}\|e^{(t-1)}\|^2. \quad (38)
\end{aligned}$$

In (38), we used the following. Let $\theta_t = (1 - \alpha\mu)^{k-t}$, and $S_t := \sum_{t=1}^k \theta_t$. Then,

$$\begin{aligned}
\left(\sum_{t=1}^k \theta_t\|e^{(t-1)}\|\right)^2 &= S_t^2\left(\sum_{t=1}^k \frac{\theta_t}{S_t}\|e^{(t-1)}\|\right)^2 \\
&\leq S_t^2 \sum_{t=1}^k \frac{\theta_t}{S_t}\|e^{(t-1)}\|^2 \\
&= S_t \sum_{t=1}^k \theta_t\|e^{(t-1)}\|^2,
\end{aligned}$$

where we used convexity of the scalar quadratic function $v \mapsto v^2$. Now, using $\sum_{t=1}^k (1 - \alpha\mu)^{k-t} \leq \frac{1}{1 - (1 - \alpha\mu)} = \frac{1}{\alpha\mu}$, (38) is further upper bounded as:

$$\begin{aligned}
\|x^{(k)} - x^\bullet\|^2 &\leq 8(1 - \alpha\,\mu)^{2k}ND^2 \\
&\quad + \frac{\alpha^2}{\alpha\mu}\sum_{t=1}^k (1 - \alpha\mu)^{k-t}\|e^{(t-1)}\|^2.
\end{aligned}$$

Taking expectation, and applying Lemma 9, we obtain the claimed result. ∎

We are now ready to prove Theorems 2 and 3.

*Proof of Theorem 2:* The proof follows from Lemma 10 by applying Lemma 8. Namely, setting $a := 1 - \alpha\mu$ and $u_t := 1 - p_{t-1}^2$, the desired result follows. ∎

*Proof of Theorem 3:* Consider (37). Taking expectation:

$$
\mathbb{E}\left[\|x^{(k)} - x^{\bullet}\|\right] \leq \sqrt{N}(1 - \alpha\mu)^k 2D \tag{39}
$$

$$
+ \ \alpha \sum_{t=1}^{k} (1 - \alpha\mu)^{k-t} \sqrt{\mathcal{C}_e(1 - p_{t-1}^2)}
$$

$$
\leq \sqrt{N}(1 - \alpha\mu)^k 2D \tag{40}
$$

$$
+ \ \alpha \sum_{t=1}^{k} (1 - \alpha\mu)^{k-t} \sqrt{\mathcal{C}_e} \sqrt{2}(\sqrt{\delta})^t.
$$

The first inequality uses $\mathbb{E}[|u|] \leq (\mathbb{E}[|u|^2])^{1/2}$. The second inequality uses $1 - p_{t-1}^2 = (1 - (1 - \delta^t))^2 \leq 2\delta^t$. Consider the sum in (40). For each $t$, each summand is upper bounded by $\eta^k \sqrt{2\mathcal{C}_e}$, and so the sum is $O(k\eta^k)$. The term $\sqrt{N}(1 - \alpha\mu)^k 2D = O(\eta^k)$. Hence, the overall right-hand-side in (40) is $O(k\eta^k) = O((\eta + \epsilon)^k)$, which completes the proof. ∎

*Remark 5* Proof of Theorem 3 also determines the constant in the convergence rate. From the above proof, substituting the expression for $\mathcal{C}_e$ in (29), it is straightforward to observe that, for all $k = 1, 2, ...$:

$$
\mathbb{E}\left[\|x^{(k)} - x^{\bullet}\|\right] \leq 12 \max\left\{\sqrt{N}\,D, \ \frac{\alpha\sqrt{N}\,G}{p_{\min}(1 - \beta)}\right\} k\,\eta^k.
$$

## V. SIMULATIONS

We provide a simulation example of learning a linear classifier via logistic loss. Simulations demonstrate that our method significantly reduces the total work (communication and computational cost) required to achieve a certain accuracy, when compared with standard distributed gradient method.

We consider distributed learning of a linear classifier via logistic loss, e.g., [37]. Each node $i$ possesses $J = 2$ data samples $\{a_{ij}, b_{ij}\}_{j=1}^{J}$. Here, $a_{ij} \in \mathbb{R}^3$ is a feature vector, and $b_{ij} \in \{-1, +1\}$ is its class label. We want to learn a vector $x = (x_1^\top, x_0)^\top$, $x_1 \in \mathbb{R}^3$, and $x_0 \in \mathbb{R}$, such that the corresponding linear classifier $\text{sign}(H_x(a)) = \text{sign}(x_1^\top a + x_0)$ minimizes the total surrogate loss with $l_2$ regularization:

$$
\sum_{i=1}^{N} \sum_{j=1}^{J} \mathcal{J}_{\text{logis}}(b_{ij} H_x(a_{ij})) + \mathcal{R}\|x\|^2, \tag{41}
$$

subject to a prior knowledge that $\|x\| \leq \mathcal{M}$, where $\mathcal{M} > 0$ is a constant. Here, $\mathcal{J}_{\text{logis}}(\cdot)$ is the logistic loss $\mathcal{J}_{\text{logis}}(\alpha) = \log(1 + e^{-\alpha})$, and $\mathcal{R}$ is a positive regularization parameter; we set $\mathcal{R} = 0.1$. Clearly, problem (41) fits the generic framework in (1) with $f_i(x) = \sum_{j=1}^{J} \mathcal{J}_{\text{logis}}(b_{ij} H_x(a_{ij})) + \frac{\mathcal{R}}{N}\|x\|^2$, $f(x) = \sum_{i=1}^{N} f_i(x)$, and $\mathcal{X} = \{x \in \mathbb{R}^4 : \|x\| \leq \mathcal{M}\}$. A strong convexity constant of the $f_i$'s $\mu$ can be taken as $\mu = \frac{\mathcal{R}}{N}$, while it can be shown that a Lipschitz constant $L$ can be taken as $\frac{1}{4N}\|\sum_{i=1}^{N} \sum_{j=1}^{J} c_{ij} c_{ij}^\top\| + \frac{\mathcal{R}}{N}$, where $c_{ij} = (b_{ij} a_{ij}^\top, b_{ij})^\top$.

We generate $a_{ij}$ independently over $i$ and $j$; each entry of $a_{ij}$ is drawn independently from the standard normal distribution. We generate the "true" vector $x^\star = ((x_1^\star)^\top, x_0^\star)^\top$ by drawing its entries independently from standard normal distribution. Then, the class labels are generated as $b_{ij} = \text{sign}((x_1^\star)^\top a_{ij} + x_0^\star + \epsilon_{ij})$, where $\epsilon_{ij}$'s are drawn independently from normal distribution with zero mean and standard

deviation $0.1$. The obtained corresponding strong convexity parameter $\mu = 0.1$, and the Lipschitz constant $L \approx 0.69$. Finally, we set $\mathcal{M} = 100$.

The network is connected and has $N = 50$ nodes and $214$ links, and it is generated as a random geometric graph: we place nodes randomly (uniformly) on a unit square, and the node pairs whose distance is less than a radius are connected by an edge. With both algorithms, we use the initialization $x_i(0) = 0, \forall i$, and we utilize the Metropolis weights, e.g., [38]. With the proposed method, we set $p_k = 1 - \delta^{k+1}$, $k = 0, 1, ...$, and $\delta = (1 - \alpha\mu)^2$.

As an error metric, we use the relative error in the objective function averaged across nodes:

$$
\frac{1}{N} \sum_{i=1}^{N} \frac{f(x_i^{(k)}) - f^\star}{f(0) - f^\star}, \ f(0) - f^\star > 0.
$$

We evaluate numerically the optimal solution $f^\star$ via the centralized projected gradient method.

We compare the two methods with respect to the total number of works across all nodes, where a single work corresponds to a single node activation at one iteration. We consider three different values of step-sizes, $\alpha \in \{\frac{1}{250\,L}, \frac{1}{50\,L}, \frac{1}{10\,L}\}$, which correspond to different achievable accuracies by both methods. (As shown in Figure 1, $\alpha = \frac{1}{10\,L}$ corresponds to accuracy $\approx 1.5 \times 10^{-2}$; $\alpha = \frac{1}{50\,L}$ corresponds to $\approx 2 \times 10^{-3}$, while $\alpha = \frac{1}{250\,L}$ allows to achieve an accuracy better than $5 \times 10^{-4}$.)
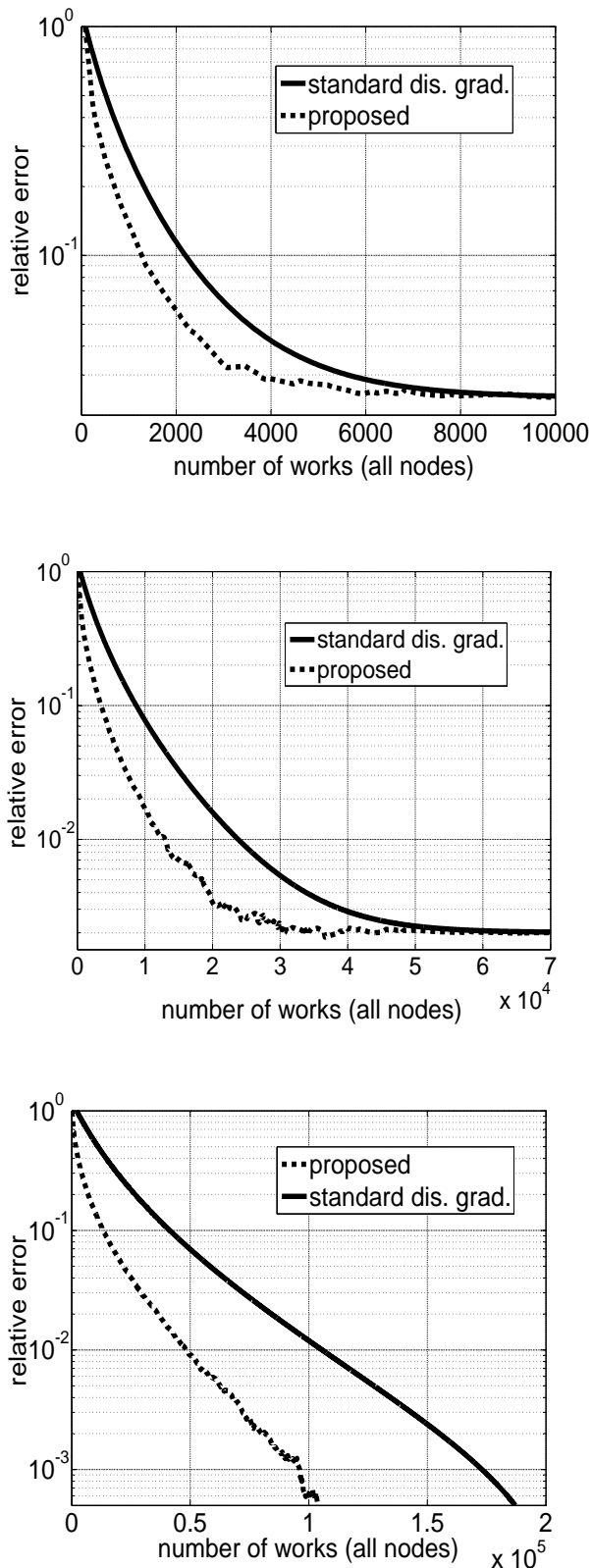
Figure 1 plots the relative error versus total number of works for the two methods (dotted line for the proposed method and solid for the standard distributed gradient method), for the three different values of step-size $\alpha$. We can see that the proposed method significantly reduces the overall work to achieve the target accuracy. For example, consider $\alpha = \frac{1}{50\,L}$. We can see that the proposed method reaches the achievable accuracy $\approx 2 \times 10^{-3}$ in about $30,000$ works ($600$ works per node), while the standard distributed gradient method takes at least $50,000$ works ($1000$ works per node). Hence, our method achieves a reduction in work of at least $40\%$. Also, we can see from the example in Figure 1 that the gain of the proposed method increases with the increase of the desired accuracy (decrease of step-size $\alpha$.)

Figure 2 plots the relative error versus number of iterations $k$ for the two methods (First from top: $\alpha = \frac{1}{10\,L}$; second from top: $\alpha = \frac{1}{50\,L}$; and bottom: $\alpha = \frac{1}{250\,L}$.) We can see that, interestingly, our method incurs almost no loss (or shows a slight gain) in terms of the number of iterations with respect to the standard distributed gradient method. Hence, our method shows a significant gain in terms of the number of works (Figure 1) while at the same time not sacrificing performance in terms of the number of iterations $k$ – even though its iterations are significantly cheaper from the iterations of the standard distributed gradient method (Figure 2).

We can see, e.g., in Figure 2 (second from top), that our method shows a slight gain in terms of $k$. This could be somewhat unexpected but is explained as follows. Namely, the standard distributed gradient method is an *exact* gradient method *on minimizing* $\Psi_\alpha$, while our proposed method is an *inexact* gradient method on minimizing the same function.

Hence, it is natural to expect that the standard distributed gradient method performs better than ours in *terms of the $\Psi_\alpha$'s optimality gap*; however, this does not have to be the case when considering the optimality gap in terms of the desired cost function $f$ in (1), as demonstrated in Figure 2 (second from top).

## VI. CONCLUSION

We explored the effect of two sources of redundancy with distributed projected gradient algorithms. The first redundancy, well-known in the literature on distributed multi-agent optimization, stems from the fact that not all inter-neighbor links need to be utilized at all iterations for the algorithm to converge. The second redundancy, explored before only in centralized optimization, arises when we minimize the sum of cost functions, each summand corresponding to a distinct data sample. In this setting, it is known that performing a gradient method with an appropriately increasing sample size can exhibit convergence properties that essentially match the properties of a standard gradient method, where the full sample size is utilized at all times. We simultaneously explored the two sources of redundancy for the first time to develop a novel distributed gradient method. With the proposed method, each node, at each iteration $k$, is active with a certain probability $p_k$, and is idle with probability $1 - p_k$, where the activation schedule is independent across nodes and across iterations. Assuming that the nodes' local costs are strongly convex and have Lipschitz continuous gradients, we showed that the proposed method essentially matches the linear convergence rate (towards a solution neighborhood) of the standard distributed projected gradient method, where all nodes are active at all iterations. Simulations on $l_2$-regularized logistic losses demonstrate that the proposed method significantly reduces the total communication and computational cost to achieve a desired accuracy, when compared with the standard distributed gradient method. As a future work, we plan to apply the proposed idling nodes strategy to other distributed multi-agent algorithms, including, e.g., the Nesterov gradient variants.

## REFERENCES

[1] A. Daneshmand, F. Facchinei, V. Kungurtsev, and G. Scutari, "Hybrid random/deterministic parallel algorithms for nonconvex big data optimization," *submitted to IEEE Trans. on Signal Processing*, 2014.

[2] K. Slavakis, G. B. Giannakis, and G. Mateos, "Modeling and optimization for big data analytics," *IEEE Signal Processing Magazine*, vol. 31, pp. 18–31, 2014.

[3] K. Slavakis, S.-J. Kim, G. Mateos, and G. B. Giannakis, "Stochastic approximation vis-a-vis online learning for big data analytics," *IEEE Signal Processing Magazine*, vol. 31, no. 11, pp. 124–129, Nov. 2014.

[4] 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. Sig. Process.*, vol. 56, no. 1, pp. 350–364, Jan. 2009.

[5] M. Rabbat and R. Nowak, "Distributed optimization in sensor networks," in *IPSN 2004, 3rd International Symposium on Information Processing in Sensor Networks*, Berkeley, California, USA, April 2004, pp. 20 – 27.

[6] S. Kar, J. M. F. Moura, and K. Ramanan, "Distributed parameter estimation in sensor networks: Nonlinear observation models and imperfect communication," *IEEE Transactions on Information Theory*, vol. 58, no. 6, pp. 3575–3605, June 2012.
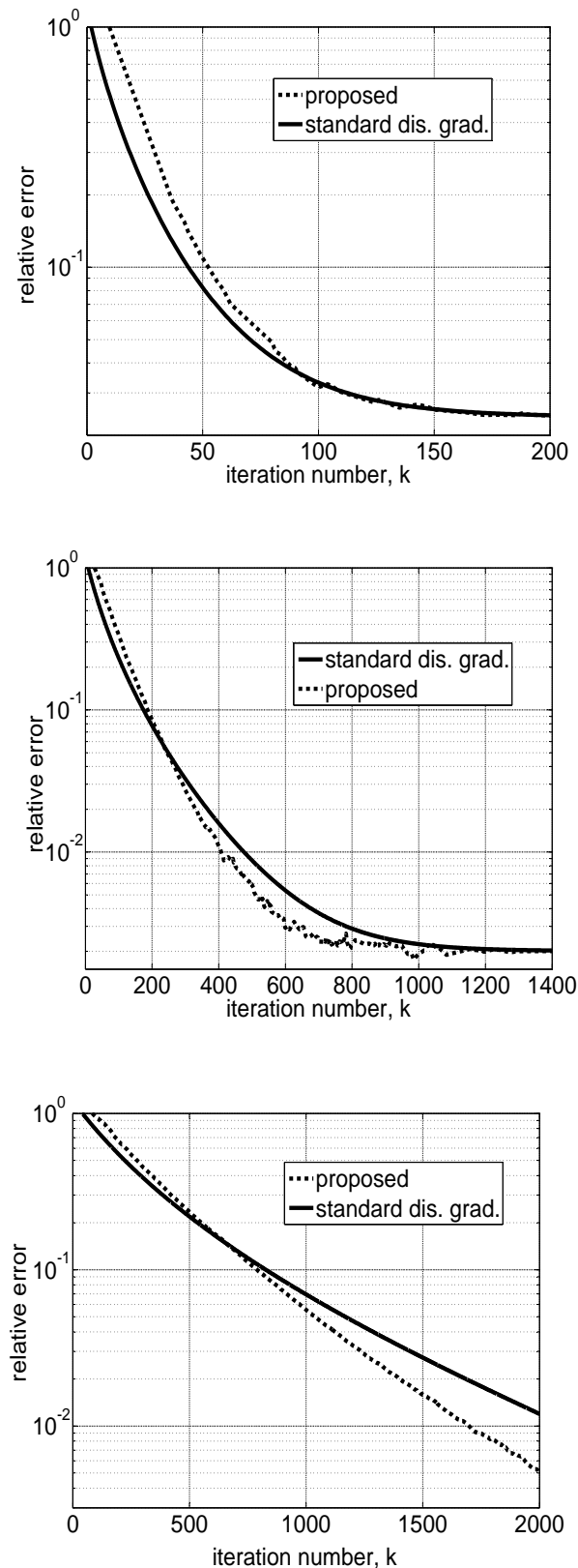


Fig. 1. Relative error versus total number of works for the two methods (dotted line corresponds to the proposed method and solid for the standard distributed gradient method), for three different values of step-size $\alpha$: First from top: $\alpha = \frac{1}{10 L}$; second from top: $\alpha = \frac{1}{50 L}$; and bottom: $\alpha = \frac{1}{250 L}$.

Fig. 2. Relative error versus number of iterations $k$ for the two methods (dotted line corresponds to the proposed method and solid for the standard distributed gradient method), for three different values of step-size $\alpha$: First from top: $\alpha = \frac{1}{10\,L}$; second from top: $\alpha = \frac{1}{50\,L}$; and bottom: $\alpha = \frac{1}{250\,L}$.

[7] C. Lopes and A. H. Sayed, "Adaptive estimation algorithms over distributed networks," in *21st IEICE Signal Processing Symposium*, Kyoto, Japan, Nov. 2006.

[8] F. Cattivelli and A. H. Sayed, "Diffusion LMS strategies for distributed estimation," *IEEE Trans. Sig. Process.*, vol. 58, no. 3, pp. 1035–1048, March 2010.

[9] I. Necoara and J. A. K. Suykens, "Application of a smoothing technique to decomposition in convex optimization," *IEEE Trans. Autom. Contr.*, vol. 53, no. 11, pp. 2674–2679, Dec. 2008.

[10] J. Mota, J. Xavier, P. Aguiar, and M. Püschel, "Distributed optimization with local domains: Applications in mpc and network flows," *to appear in IEEE Trans. Autom. Contr.*, 2015.

[11] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, January 2009.

[12] S. Ram, A. Nedic, and V. Veeravalli, "Distributed stochastic subgradient projection algorithms for convex optimization," *Jour. Opt. Theory and App.*, vol. 147, no. 3, pp. 516–545, 2011.

[13] S. S. Ram, A. Nedic, and V. Veeravalli, "Asynchronous gossip algorithms for stochastic optimization," in *CDC '09, 48th IEEE International Conference on Decision and Control*, Shanghai, China, December 2009, pp. 3581 – 3586.

[14] A. Mokhtari, Q. Ling, and A. Ribeiro, "Network newton," 2014, available at: http://arxiv.org/abs/1412.3740.

[15] M. Zargham, A. Ribeiro, A. Ozdaglar, and A. Jadbabaie, "Accelerated dual descent for network flow optimization," *IEEE Transactions on Automatic Control*, vol. 59, no. 4, pp. 905–920, 2014.

[16] D. Jakovetic, J. Xavier, and J. M. F. Moura, "Fast distributed gradient methods," *IEEE Trans. Autom. Contr.*, vol. 59, no. 5, pp. 1131–1146, May 2014.

[17] I.-A. Chen and A. Ozdaglar, "A fast distributed proximal gradient method," in *Allerton Conference on Communication, Control and Computing*, Monticello, IL, October 2012.

[18] K. Yuan, Q. Ling, and W. Yin, "On the convergence of decentralized gradient descent," 2013, available at: http://arxiv.org/abs/1310.7063.

[19] I. Matei and J. S. Baras, "Performance evaluation of the consensus-based distributed subgradient method under random communication topologies," *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 4, pp. 754–771, 2011.

[20] I. Lobel and A. Ozdaglar, "Convergence analysis of distributed sub-gradient methods over random networks," in *46th Annual Allerton Conference onCommunication, Control, and Computing*, Monticello, Illinois, September 2008, pp. 353 – 360.

[21] M. Schmidt, N. L. Roux, and F. Bach, "Convergence rates of inexact proximal-gradient methods for convex optimization," in *Advances in Neural Information Processing Systems 24*, 2011, pp. 1458–1466.

[22] J. Tsitsiklis, D. Bertsekas, and M. Athans, "Distributed asynchronous deterministic and stochastic gradient optimization algorithms," *IEEE Trans. Autom. Contr.*, vol. 31, no. 9, pp. 803–812, Sep. 1986.

[23] S.-Y. Tu and A. H. Sayed, "On the influence of informed agents on learning and adaptation over networks," *IEEE Trans. Signal Processing*, vol. 61, no. 6, pp. 1339–1356, March 2013.

[24] K. I. Tsianos, S. F. Lawlor, J. Y. Yu, and M. G. Rabbat, "Networked optimization with adaptive communication," in *IEEE GlobalSIP Network Theory Symposium*, Austin, Texas, December 2013.

[25] G. Deng and M. C. Ferris, "Variable-number sample path optimization," *Mathematical Programming*, vol. 117, no. 1–2, pp. 81–109, 2009.

[26] T. H. de Mello, "Variable-sample methods for stochastic optimization," *ACM Transactions on Modeling and Computer Simulation*, vol. 13, no. 2, pp. 108–133, 2003.

[27] E. Polak and J. O. Royset, "Eficient sample sizes in stochastic nonlinear programing," *Journal of Computational and Applied Mathematics*, vol. 217, no. 2.

[28] R. Pasupathy, "On choosing parameters in restrospective-approximation algorithms for simulation-optimization," in *2006 Winter Simulation Conference, L.F. Perrone, F.P. Wieland, J. Liu, B.G. Lawson, D.M. Nicol and R.M. Fujimoto, eds.*, 2006, pp. 208–215.

[29] ——, "On choosing parameters in retrospective-approximation algorithms for stochastic root finding and simulation optimization," *Operations Research*, vol. 58, no. 4, pp. 889–901, 2010.

[30] F. Bastin, "Trust-region algorithms for nonlinear stochastic programming and mixed logit models," 2004, phD Thesis, University of Namur, Belgium.

[31] F. Bastin, C. Cirillo, and P. L. Toint, "An adaptive monte carlo algorithm for computing mixed logit estimators," *Computational Management Science*, vol. 3, no. 1, pp. 55–79, 2006.

[32] N. Krejić and N. Krklec, "Line search methods with variable sample size for unconstrained optimization," *Journal of Computational and Applied Mathematics*, vol. 245, pp. 213–231, 2013.

[33] N. Krejić and N. K. Jerinkić, "Nonmonotone line search methods with variable sample size," *Numerical Algorithms*, vol. 68, no. 4.

[34] A. Nedic, A. Ozdaglar, and A. Parrilo, "Constrained consensus and optimization in multi-agent networks," *IEEE Transactions on Automatic Control*, vol. 55, no. 4, pp. 922–938, April 2010.

[35] D. Bajovic, J. Xavier, J. M. F. Moura, and B. Sinopoli, "Consensus and products of random stochastic matrices: Exact rate for convergence in probability," *IEEE Trans. Sig. Process.*, vol. 61, no. 10, pp. 2557–2571, May 2013.

[36] D. Jakovetic, J. M. F. Moura, and J. Xavier, "Distributed Nesterov-like gradient algorithms," in *CDC'12, 51$^{st}$ IEEE Conference on Decision and Control*, Maui, Hawaii, December 2012, pp. 5459–5464.

[37] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning, Michael Jordan, Editor in Chief*, vol. 3, no. 1, pp. 1–122, 2011.

[38] L. Xiao, S. Boyd, and S. Lall, "A scheme for robust distributed sensor fusion based on average consensus," in *IPSN '05, Information Processing in Sensor Networks*, Los Angeles, California, 2005, pp. 63–70.